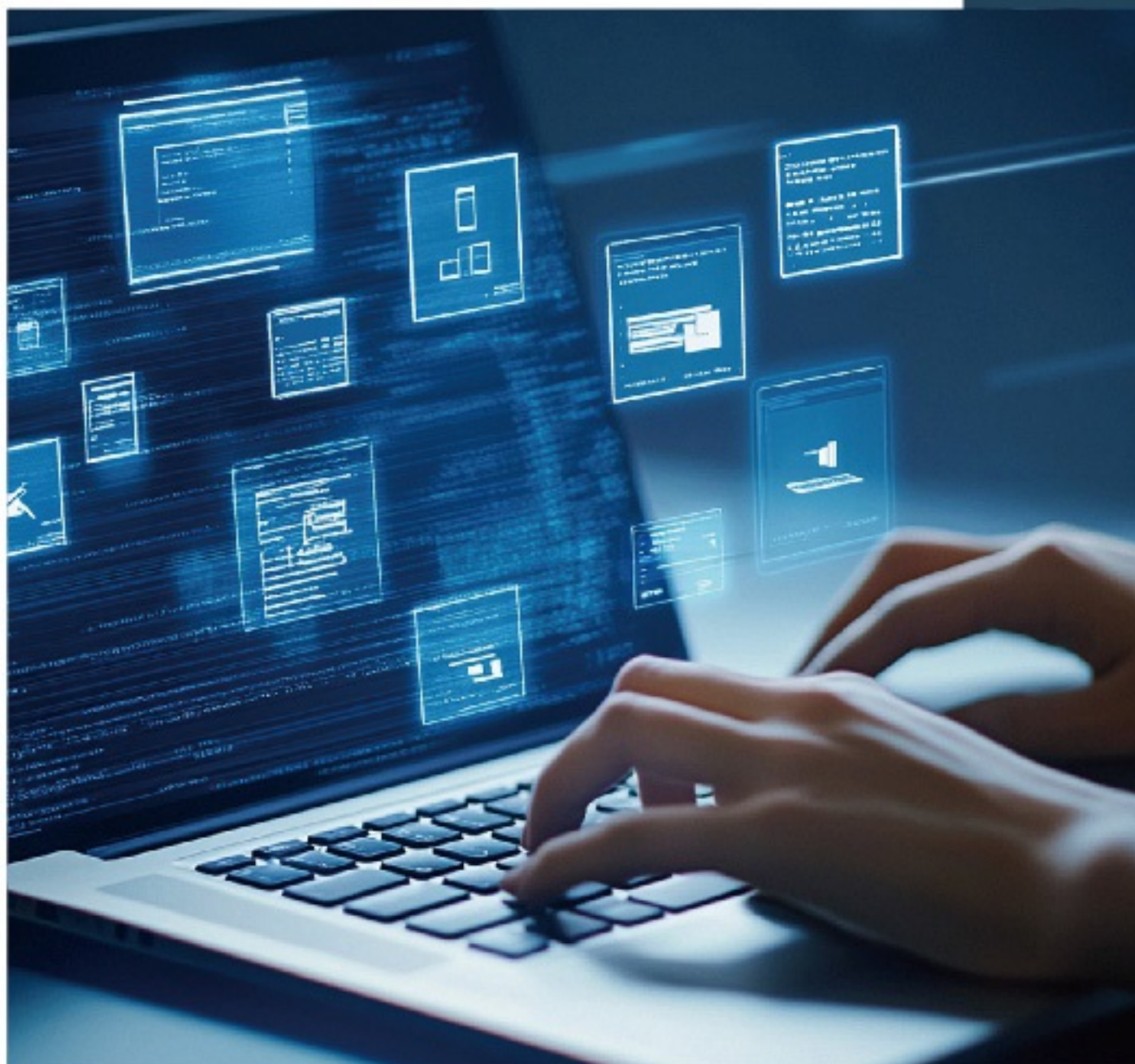


Computing and Artificial Intelligence

<https://ojs.acad-pub.com/index.php/CAI>



2024 VOLUME 2 ISSUE 2
ISSN: 3029-2786 (Online)



2



Editorial Team

Editor-in-Chief

Shaohua Wan

University of Electronic Science and Technology of China
China

Associate Editor

Vijayakumar Varadarajan

La Trobe University
Australia

Editorial Board Members

Syed Muzamil Basha

REVA University
India

Junyu Xuan

University of Technology Sydney
Australia

Abdullah Ayub Khan

Bahria University
Pakistan

Ata Jahangir Moshayedi

Jiangxi University of Science and
Technology
China

Mingchu Li

Jiangxi Normal University
China

Jie Zhang

Xi'an Jiaotong-Liverpool University
China

Parikshit N Mahalle

Vishwakarma Institute of Technology
India

Anjum Razzaque

Western Illinois University
United States

William Cheng-Chung Chu

Fuyao Institute of Science and Technology
Taiwan

Umesh C. Pati

National Institute of Technology
India

Jun Ye

Ningbo University
China

Jiakai Wang

Zhongguancun Laboratory
China

Grigorios N. Beligiannis

University of Patras
Greece

Chin-Shiuh Shieh

National Kaohsiung University of Science
and Technology
Taiwan

Janusz Kacprzyk

Systems Research Institute of the Polish
Academy of Sciences
Poland

Esma Aïmeur

University of Montreal
Canada

Marcin Paprzycki

Polish Academy of Sciences
Poland

Inam Ullah

Gachon University
Korea

Shashi Kant Gupta

Eudoxia Research University
United States

Yi-Chung Hu

Chung Yuan Christian University
Taiwan

Neelamadhab Padhy

Gandhi Institute of Engineering and
Technology
India

Omar Cheikhrouhou

University of Sfax
Tunisia

Kibum Kim

Hanyang University
Korea

Maki K. Habib

The American University in Cairo
Egypt

Yousef Daradkeh

Prince Sattam bin Abdulaziz University
Saudi Arabia

Hao Ying

Wayne State University
United States

Pradeep Kumar Mallick

Kalinga Institute of Industrial Technology
India

Salah Bourenane

Aix Marseille Université
France

Ibrahim A Hameed

Norges Teknisk-Naturvitenskapelige
Universitet
Norway

Alberto Gotta

Consiglio Nazionale delle Ricerche
Italy

Shuo-Tsung Chen

Tunghai University
Taiwan

Saeed Alsamhi

Ibb University
Yemen

Chao Zhang

Shanxi University
China

Mohammad Khishe

Imam Khomeini Marine Science University
Iran

Pushpendu Kar

University of Nottingham Ningbo China
China

Sudan Jha

Kathmandu University
Nepal

Mohammed Baz

Taif University
Saudi Arabia

Xin Liu

Lappeenranta-Lahti University of
Technology
Finland

Ray-I Chang

National Taiwan University

Taiwan

Cihan Karakuzu

Bilecik Şeyh Edebali University

Turkey

Volume 2 Issue 2 • 2024

Computing and Artificial Intelligence

Editor-in-Chief

Prof. Shaohua Wan

University of Electronic Science and Technology of China, China



Computing and Artificial Intelligence

<https://ojs.acad-pub.com/index.php/CAI>

Contents

Articles

- 1 Plant leaf disease classification using FractalNet**
Hmidi Alaeddine, Malek Jihene
- 8 Predicting manipulated regions in deepfake videos using convolutional vision transformers**
Mohan Bhandari, Sushant Shrestha, Utsab Karki, Santosh Adhikari, Rajan Gaihre
- 18 Software cost estimation tool: A App based application, estimate the cost of software project**
Ajay Jaiswal, Piyush Malviya, Lucky Parihar, Rani Pathak, Kuldeep Rajput
- 28 Enhancing user experience and trust in advanced LLM-based conversational agents**
Yuanyuan Xu, Weiting Gao, Yining Wang, Xinyang Shan, Yin-Shan Lin
- 47 Exploring other clustering methods and the role of Shannon Entropy in an unsupervised setting**
Erin Chelsea Hathorn, Ahmed Abu Halimeh
- 59 Validation of the practicability of logical assessment formula for evaluations with inaccurate ground-truth labels: An application study on tumour segmentation for breast cancer**
Yongquan Yang, Hong Bu

85 Innovation dynamics in BRICS economies investigated by artificial intelligence (AI)

Claudio Zancan, João Luiz Passador, Cláudia Souza Passador, Ricardo Carvalho Rodrigues

128 Clustering data analytics of urban land use for change detection

C. Rajabhushanam

140 The based-biofeedback approach versus ECG for evaluation heart rate variability during the maximal exercise protocol among healthy individuals

Sara Pouriamehr, Valiollah Dabidi Roshan, Somayeh Namdar Tajari

159 Harnessing artificial intelligence (AI) for cybersecurity: Challenges, opportunities, risks, future directions

Zarif Bin Akhtar, Ahmed Tajbiul Rawol

Review

181 Applications of reinforcement learning, machine learning, and virtual screening in SARS-CoV-2-related proteins

Yasunari Matsuzaka, Ryu Yashiro

Plant leaf disease classification using FractalNet

Hmidi Alaeddine^{1,*}, Malek Jihene^{1,2}

¹Laboratory of Electronics and Microelectronics, LR99ES30, Faculty of Sciences of Monastir, Monastir University, Monastir 5000, Tunisia

²Higher Institute of Applied Sciences and Technology of Sousse, Sousse University, Sousse 4000, Tunisia

* Corresponding author: Hmidi Alaeddine, ahmidi@outlook.fr

CITATION

Alaeddine H, Jihene M. Plant leaf disease classification using FractalNet. *Computing and Artificial Intelligence*. 2024; 2(2): 545. <https://doi.org/10.59400/cai.v2i2.545>

ARTICLE INFO

Received: 17 February 2024

Accepted: 12 June 2024

Available online: 3 July 2024

COPYRIGHT



Copyright © 2024 by author(s).

Computing and Artificial Intelligence is published by Academic Publishing Pte. Ltd. This work is licensed under the Creative Commons Attribution (CC BY) license.

<https://creativecommons.org/licenses/by/4.0/>

Abstract: In this work, an effort is made to apply the FractalNet model in the field of plant disease classification. The proposed model was trained and tested using a “PlantVillage” plant disease image dataset using a central processing unit (CPU) environment for 300 epochs. It produced an average classification accuracy of 99.9632% on the test dataset. The experimental results demonstrate the efficiency of the proposed model and show that the model achieved the highest values compared to other deep learning models in the PlantVillage datasets.

Keywords: FractalNet; convolution neural network; plant village; plant leaf disease

1. Introduction

Diseases, insects and nutrient deficiencies are the most common threats to crop growth, negatively affecting total crop production and the farmer's net profit. Diagnosis and treatment of diseases and application of fertilizers play an important role in reducing yield loss.

Therefore, accurate and early disease detection is necessary as it is among the best possible solutions for early disease control and improved crop performance as well as avoiding unnecessary waste of financial resources.

Conventional disease detection is not feasible for all cultivated fields and all farmers. This requires finding suitable human experts to diagnose and treat diseases, which takes time and money.

Hence the need for an intelligent system capable of automatically classifying and diagnosing plant diseases to overcome the difficulties of the traditional approach.

Today, with the activation and application of artificial intelligence in the field of agriculture and food security, many deep learning (DL) models have been used, and many models of deep learning methods have been proposed to detect and classify plant diseases.

In the research that follows, we present our approach to the challenge with two objectives. The first objective is to study and determine the relevance of the FractalNet for the task of classifying plant diseases. The second goal is to get the lowest possible error on a set of PlantVillage test images. It should be noted that no study has attempted to address this aspect before, as this is the first work that addresses the field of plant leaf disease classification using FractalNets.

The main contributions of this research are:

- We applied for the first time the FractalNet model on the PlantVillage database for the classification of plant diseases.
- We present a detailed experimental study of the FractalNet for the plant disease classification task on a set of PlantVillage test images.

- Finally, we show that the application of FractalNet allows to obtain state-of-the-art results on the PlantVillage dataset considerably improving the accuracy.

The rest of this work is organized as follows: In section 2, an overview of related works is given. Section 3 describes the database. In section 4, the data preparation, proposed model and implementation details were presented. Experimental evaluations and comparative analysis are presented and discussed in section 5. Advantage and future work are reported in section 6. The work is concluded in the last section.

2. Related works

Recent developments in artificial intelligence techniques enable the effective identification of many diseases and pest attacks in precision agriculture. This investigation deals with modern artificial intelligence approaches for the detection of plant diseases.

For the detection of rice plant diseases, Lu et al. [1] proposed a new method for identifying rice diseases. The model is able to identify ten rice diseases. Chen et al. [2] trained a model called DENSINCEP based on deep transfer learning.

Sun et al. [3] have developed an improved CNN that offers a test accuracy equivalent to 99.35%. Mohanty et al. [4] classified plant diseases using CNN models such as AlexNet and GoogLeNet. Too et al. [5] exploited CNN models such as ResNet50, VGG16, ResNet101, Resnet152, Inception V4 and DenseNets 121. Atila et al. [6] proposed an EfficientNet deep learning architecture for plant disease classification. Performance is compared to other CNN models such as AlexNet, VGG16 and ResNet50.

An effort is made by Alaeddine and Jihene [7] to apply the Wide Residual Networks model in the field of plant disease classification. Moreover, they have proposed DbneAlexnet in the study of Alaeddine and Jihene [8].

The literature review shows that most of the work in the literature exploits the PlantVillage dataset and performs disease classification of a particular plant or multiple plants [9–14]. Moreover, the literature review recognized that residual and dense convolutional neural networks performed better than other transfer learning techniques in plant disease detection [15]. Transfer learning techniques can lead to negative transfer and overfitting issues when using the architecture and weights of pretrained models for new applications. In addition, the study of the literature shows the importance of data augmentation for classification algorithms.

From these literatures, we consider all these previous approaches and the experiments already performed to determine the best deep learning model a new approach in order to obtain better accuracy on the PlantVillage dataset. In this context, we adapted the FractalNet model on the PlantVillage database for the classification of plant diseases.

3. PlantVillage database

The PlantVillage database is introduced by Hughes et al. [16] to enable the development of mobile diagnostics of diseases It consists of 61,486 images of healthy and unhealthy leaves classified into 38 groups by type and disease. The PlantVillage dataset was created with six different augmentation techniques to create more diverse

datasets with different background conditions. The augmentations used in this process were scaling, rotation, noise injection, gamma correction, image flip, and PCA color augmentation.

The images in the database are colored and have different sizes that are why the images have been resized to 227×227 , which is the default size accepted by the model. Examples of some plant diseases are shown in **Figure 1**.

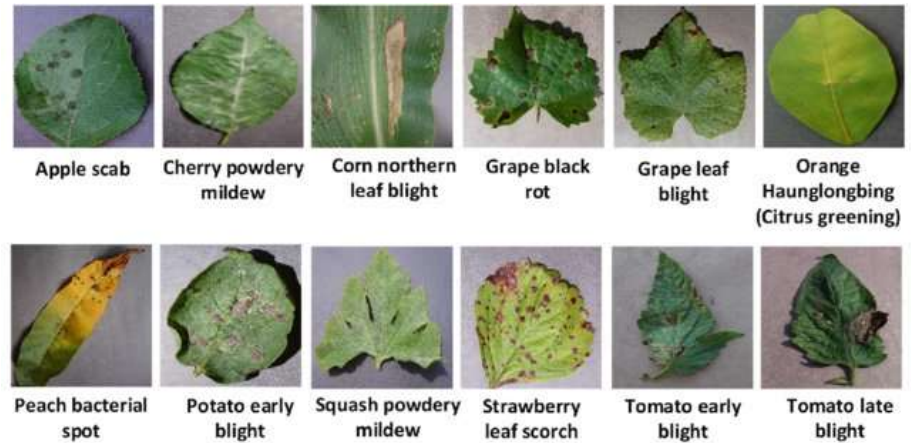


Figure 1. Some of the plant diseases in the PlantVillage dataset.

4. Materials and methods

The implementation steps of the proposed plant disease detection model are categorized into two phases. The implementation of the proposed model started with data preparation. The data preparation phase focuses on augmenting the data. The model training phase includes the design and training processes.

4.1. Data preparation

Implementing a deep learning algorithm starts with the data preparation phase. It includes data collection, data augmentation, and pre-processing steps. Some classes in the original dataset have fewer samples. On the other hand, some classes have more images. The difference sometimes reaches two thousand images. The number of samples must be equal in each class in order to increase the performance of the classification algorithms. Data augmentation techniques were used to increase the number of samples without the need to collect new data. Cropping, scaling, flipping, rotating, filling, affine transforming and tinting techniques were used to produce augmented images on the dataset. After data augmentation, the database was split for the training, validation and testing process.

4.2. FractalNet

FractalNet is introduced by Larsson et al. [17]. It is described as a type of convolutional neural network that avoids residual connections in favor of a “fractal” design. They involve the repeated application of a simple expansion rule to generate deep networks whose structural arrangements are precisely truncated fractals. These networks contain interactive subpaths of varying lengths, but do not include any direct or residual connections; each internal signal is transformed by a filter and a non-

linearity before being seen by the following layers.

4.2.1. FractalNet architecture

For the base case, $f_1(z)$ is the convolutional layer:

$$f_1(z) = \text{conv}(z)$$

After that, the recursive fractals are:

$$f_{C+1}(z) = [(f_C \circ f_C)(z)] \oplus [\text{conv}(z)]$$

where C denotes the number of columns as shown in **Figure 2b**. The number of convolutional layers at the deepest path of a single block is equivalent to $2^{(C-1)}$. For a number of columns $C = 4$, the number of convolutional layers is equivalent to $2^{(4-1)} = 2^3 = 8$ layers.

For the joining layer in green color, the elemental mean is calculated. It is not concatenation or addition.

For 5 blocks ($B = 5$) cascaded like FractalNet as shown in **Figure 2c**, the number of deepest-path convolutional layers in the entire network is $B \times 2^{(C-1)}$, i.e., $5 \times 2^{(4-1)} = 5 \times 2^3 = 5 \times 8 = 40$ layers.

A layer of 2×2 maximum pooling is performed between every two blocks to reduce the size of feature maps. Batch Normalization layers and ReLU activation functions are used after each convolution.

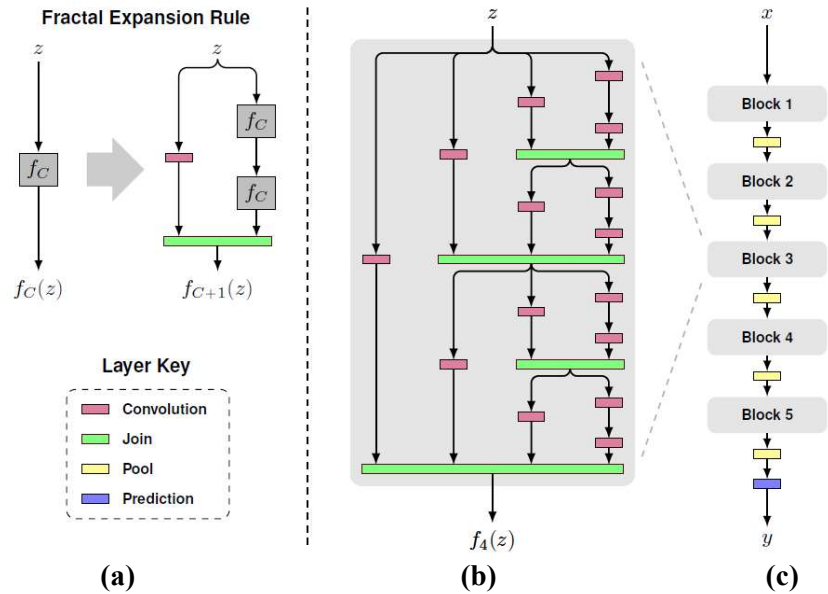


Figure 2. (a) fractal architecture: a simple fractal expansion; (b) recursive stacking of the fractal expansion in one block; (c) 5 cascading blocks like FractalNet [17].

4.2.2. Architecture of training models

For FractalNet-34, we use the same first and last layer as ResNet-34 [18], the middle of the network consists of 4 blocks ($B = 4$) and 4 columns ($C = 4$). We define the number of filter channels in blocks 1 to 4 as 128, 256, 512, 1024.

4.3. Implementation details

We formed our configurations using a ‘‘Root Mean Square Propagation Algorithm’’ with a batch size equivalent to 32 and a weight decrease of 0.0001.

The learning rate was initialized to 0.01 and divided by 10 twice before the end. We formed the network for about 300 cycles at most in a central processing unit (CPU). There was no significant change in performance after reaching 300 epochs. The python algorithm based on the deep learning library “Keras” to classify and recognize images provides the implementation of the CPU.

5. Results and discussion

5.1. The performances of FractalNet-34

Classification accuracy, precision, recall, and $F1$ score are the standard measures to assess the overall performance of classification techniques. FractalNet-34 performance is shown in **Table 1**:

Table 1. Performance of FractalNet-34.

Model	Accuracy	Precision	Sensitivity	$F1$ -score	Specificity
FractalNet-34	99.9632	98.92	98.98	98.95	98.11

5.2. Discussions

The main objective of this work is to examine and evaluate the success of the FractalNet-34 model in the classification of plant diseases and to compare the performances obtained with models from the literature. This section deals with the performance of the proposed FractalNet-34 model in the classification of plant diseases. Furthermore, it compares the proposed model with other state-of-the-art deep learning models and techniques. A comparison with the results of different studies and works is presented in **Table 2**.

Table 2. PlantVillage test accuracy.

Ref.	Method	Precision (%)
[4]	GoogleNet	99.35
[5]	DenseNets-121	99.75
[6]	EfficientNet B5	98.42
	WRN-22-2	99.9394
[7]	WRN-28-10 (dropout)	99.9611
	WRN-40-2	99.9533
[19]	Hybrid principal component analysis	59.1
[20]	Dilated TL and ensemble learning	99.10
Our	FractalNet	99.9632

The experimental results demonstrate the effectiveness of the proposed contribution. Moreover, they show that the proposed model offers better results in terms of classification accuracy than the various other models.

6. Advantages and future work

The proposed model offers an interesting test precision compared to the various

works reported in the literature. The importance of the exploited model also lies in its repetitive and homogeneous structure which makes it very suitable and compatible for integration into embedded system applications. In future work, it is planned to extend and augment the PlantVillage dataset artificially by increasing the number of classes. This will contribute to the development of models and architectures capable of achieving more precise and interesting accuracies.

7. Conclusion

The automatic detection and classification of plant diseases is a crucial process in agriculture. This work applied a new deep convolutional neural network for the classification of common diseases in different plants. In this work, some recent image augmentation techniques were used to prepare the proposed dataset for model training. The training process was performed on a CPU for up to 300 training epochs. The classification accuracy of the proposed model was 99.9632%.

Author contributions: Conceptualization, HA and MJ; methodology, HA; software, HA; validation, HA; formal analysis, HA; investigation, HA; resources, HA; data curation, HA; writing—original draft preparation, HA; writing—review and editing, HA; visualization, HA; supervision, HA and MJ; project administration, HA; funding acquisition, HA. All authors have read and agreed to the published version of the manuscript.

Conflict of interest: The authors declare no conflict of interest.

References

1. Lu Y, Yi S, Zeng N, et al. Identification of rice diseases using deep convolutional neural networks. *Neurocomputing*. 2017; 267: 378-384. doi: 10.1016/j.neucom.2017.06.023
2. Chen J, Zhang D, Nanekaran YA, et al. Detection of rice plant diseases based on deep transfer learning. *Journal of the Science of Food and Agriculture*. 2020; 100(7): 3246-3256. doi: 10.1002/jsfa.10365
3. Sun J, Tan WJ, Mao HP, et al. Identification of Leaf Diseases of Various Plants Based on Improved Convolutional Neural Network. *Agric. Eng. Newsp*. 2017; 19: 209-215.
4. Mohanty SP, Hughes DP, Salathé M. Using Deep Learning for Image-Based Plant Disease Detection. *Frontiers in Plant Science*. 2016; 7. doi: 10.3389/fpls.2016.01419
5. Too EC, Yujian L, Njuki S, et al. A comparative study of fine-tuning deep learning models for plant disease identification. *Computers and Electronics in Agriculture*. 2019; 161: 272-279. doi: 10.1016/j.compag.2018.03.032
6. Atila Ü, Uçar M, Akyol K, et al. Plant leaf disease classification using EfficientNet deep learning model. *Ecological Informatics*. 2021; 61: 101182. doi: 10.1016/j.ecoinf.2020.101182
7. Alaeddine H, Jihene M. Plant leaf disease classification using Wide Residual Networks. *Multimedia Tools and Applications*. 2023; 82(26): 40953-40965. doi: 10.1007/s11042-023-15226-y
8. Alaeddine H, Jihene M. Deep Batch-normalized eLU AlexNet for Plant Diseases Classification. In: 2021 18th International Multi-Conference on Systems, Signals & Devices (SSD); 22 March 2021. doi: 10.1109/ssd52085.2021.9429404
9. Cruz AC, Luvisi A, De Bellis L, et al. Vision-based plant disease detection system using transfer and deep learning. In: *Proceedings of the 2017 Asabe Annual International Meeting*; 16-19 July 2024; Washington.
10. Yamamoto K, Togami T, Yamaguchi N. Super-Resolution of Plant Disease Images for the Acceleration of Image-based Phenotyping and Vigor Diagnosis in Agriculture. *Sensors*. 2017; 17(11): 2557. doi: 10.3390/s17112557
11. Walleign S, Polceanu M, Buche C. Soybean Plant Disease Identification Using Convolutional Neural Network. In: *Proceedings of the Thirty-First International Florida Artificial Intelligence Research Society Conference (FLAIRS-31)*. pp. 146-151.

12. Sibiya M, Sumbwanyambe M. A Computational Procedure for the Recognition and Classification of Maize Leaf Diseases Out of Healthy Leaves Using Convolutional Neural Networks. *AgriEngineering*. 2019; 1(1): 119-131. doi: 10.3390/agriengineering1010009
13. Rangarajan AK, Purushothaman R, Ramesh A. Tomato crop disease classification using pre-trained deep learning algorithm. *Procedia Computer Science*. 2018; 133: 1040-1047. doi: 10.1016/j.procs.2018.07.070
14. Zhang K, Wu Q, Liu A, et al. Can Deep Learning Identify Tomato Leaf Disease? *Advances in Multimedia*. 2018; 2018: 1-10. doi: 10.1155/2018/6710865
15. Hu G, Wu H, Zhang Y, Wan M. A low shot learning method for tea leaf's disease identification. *Computers and Electronics in Agriculture*. 2019; 163. doi: 10.1016/j.compag.2019.104852.104852
16. Hughes DP, Salathe M. An open access repository of images on plant health to enable the development of mobile disease diagnostics. *arXiv*. 2015; arXiv:1511.08060.
17. Larsson G, Maire M, Shakhnarovich G. Fractalnet: Ultra-deep neural networks without residuals. *arXiv*. 2016; arXiv:1605.07648.
18. He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. *arXiv*. 2015; arXiv:1512.03385.
19. Gadekallu TR, Rajput DS, Reddy MPK, et al. A novel PCA-whale optimization-based deep neural network model for classification of tomato plant diseases using GPU. *Journal of Real-Time Image Processing*. 2020; 18(4): 1383-1396. doi: 10.1007/s11554-020-00987-8
20. Saberi Anari M. A Hybrid Model for Leaf Diseases Classification Based on the Modified Deep Transfer Learning and Ensemble Approach for Agricultural AIoT-Based Monitoring. In: Kumar A (editor). *Computational Intelligence and Neuroscience*. 2022; 1-15. doi: 10.1155/2022/6504616

Article

Predicting manipulated regions in deepfake videos using convolutional vision transformers

Mohan Bhandari^{1,*}, Sushant Shrestha², Utsab Karki², Santosh Adhikari², Rajan Gaihre²

¹ Department of Science and Technology, Samriddhi College, Lokanthali, Bhaktapur 44800, Nepal

² Department of Computer Engineering, Kantipur Engineering College, Dhapakhel, Lalitpur 44700, Nepal

* Corresponding author: Mohan Bhandari, mail2mohanbhandari@gmail.com

CITATION

Bhandari M, Shrestha S, Karki U, et al. Predicting manipulated regions in deepfake videos using convolutional vision transformers. *Computing and Artificial Intelligence*. 2024; 2(2): 1409.
<https://doi.org/10.59400/cai.v2i2.1409>

ARTICLE INFO

Received: 30 May 2024

Accepted: 13 June 2024

Available online: 1 July 2024

COPYRIGHT



Copyright © 2024 by author(s).

Computing and Artificial Intelligence is published by Academic Publishing Pte. Ltd. This work is licensed under the Creative Commons Attribution (CC BY) license.

<https://creativecommons.org/licenses/by/4.0/>

Abstract: Deepfake technology, which uses artificial intelligence to create and manipulate realistic synthetic media, poses a serious threat to the trustworthiness and integrity of digital content. Deepfakes can be used to generate, swap, or modify faces in videos, altering the appearance, identity, or expression of individuals. This study presents an approach for deepfake detection, based on a convolutional vision transformer (CViT), a hybrid model that combines convolutional neural networks (CNNs) and vision transformers (ViTs). The proposed study uses a 20-layer CNN to extract learnable features from face images, and a ViT to classify them into real or fake categories. The study also employs MTCNN, a multi-task cascaded network, to detect and align faces in videos, improving the accuracy and efficiency of the face extraction process. The method is assessed using the FaceForensics++ dataset, which comprises 15,800 images sourced from 1600 videos. With an 80:10:10 split ratio, the experimental results show that the proposed method achieves an accuracy of 92.5% and an AUC of 0.91. We use Gradient-Weighted Class Activation Mapping (Grad-CAM) visualization that highlights distinctive image regions used for making a decision. The proposed method demonstrates a high capability of detecting and distinguishing between genuine and manipulated videos, contributing to the enhancement of media authenticity and security.

Keywords: face detection; machine learning; vision transformer; convolution neural networks; Grad-CAM

1. Introduction

Technologies for altering images and videos are developing rapidly. The rise of fake technology has gained significant attention in recent years due to its ability to generate highly realistic, manipulated media. The different techniques and technical expertise needed to create and manipulate digital content are also easily accessible, as there is abundant reading material on the internet [1]. Currently, it is possible to seamlessly generate hyper-realistic digital images with a few resources and easy-to-follow instructions available online [2]. Deepfake is a technique that aims to replace the face of a targeted person with the face of someone else in a video. It is created by splicing the synthesized face region into the original image. The term can also mean to represent the final output of a hyper-realistic video created. Deepfakes can be used for the creation of hyper-realistic Computer-generated imagery (CGI), Virtual Reality (VR), Augmented Reality (AR), Education, Animation, Arts, and Cinema. However, since Deepfakes are deceptive, they can also be used for malicious purposes [3]. Deepfake detection is the task of identifying and exposing digital falsifications of images, video, and audio that are created with machine learning

techniques [4]. This task poses a formidable challenge to privacy, democracy, and national security, as deepfakes can be used to manipulate public opinion, deceive voters, undermine trust in institutions, exacerbate social divisions, endanger public safety, disrupt international relations, and jeopardize national security. Detecting deepfakes is not only technically difficult but also socially and legally complex. Technical solutions, such as forensic analysis, digital watermarking, and immutable authentication trails, face limitations in accuracy, scalability, and usability [5]. Social and legal solutions, such as media literacy, platform regulation, and legal liability face trade-offs between free expression, privacy, and accountability. Moreover, deepfake creators can adapt to detection methods and exploit cognitive biases that make people susceptible to believing and spreading false information. Therefore, deepfake detection requires a multidisciplinary and collaborative approach that balances the benefits and harms of deepfake technology [6].

The challenge of deepfake detection is the diversity and complexity of deepfake generation methods. There are various types of deepfake techniques, such as face swapping, face reenactment, lip-syncing, voice cloning, and text generation [7]. Each of these techniques requires different approaches and models to create and manipulate digital content. Moreover, the quality and realism of fake media vary depending on the data, algorithms, and parameters used for the generation process. Therefore, it is difficult to design a universal and effective deep fake detector that can handle all kinds of deep fake scenarios.

In this study, we propose to leverage the power of convolutional vision transformer (CViT) to develop a comprehensive and robust deepfake detection framework that can adapt to different types of deepfake techniques and media. By utilizing the capabilities of CViT and focusing on the inconsistency in pixel-level details, we aim to address the disadvantages of deepfake technology and provide a robust defense against its malicious usage. This study strives to contribute to the development of advanced deepfake detection techniques, enhancing the security and integrity of digital media in an increasingly vulnerable landscape [7].

2. Literature review

In “Deepfakes Detection with Automatic Face Weighting”, Montserrat et al. [8] proposed a novel method utilizing convolutional neural networks (CNNs) and recurrent neural networks (RNNs) to detect deepfakes. This approach extracts visual and temporal features from facial regions in videos for effective manipulation identification. The study uses the Deepfake Detection Challenge (DFDC) dataset, comprising over 100,000 videos with various facial modifications. The method employs CNNs and RNNs to detect and localize manipulated faces, showing competitive performance against existing techniques. It can handle videos with multiple faces, varying quality, and different manipulation methods, and provides a confidence score for each face region. The reported accuracy is 92.61% in detecting forgeries. However, the method struggles with highly realistic manipulations in blurry or low-quality images and does not incorporate audio information, which could enhance detection performance. In “MesoNet: A Compact Facial Video Forgery Detection Network,” Afchar et al. [9] present an efficient method for

detecting manipulated faces in videos, focusing on Deepfake and Face2Face techniques. They utilize two datasets: the Deepfake dataset, with 175 forged videos and frames extracted and aligned, and the FaceForensics-based Face2Face dataset, with over a thousand videos. Training and testing sets include 5111 forged and 7250 real images, 2889 forged and 4259 real images, respectively, for Deepfake; and 300 training and 150 testing videos for Face2Face. Traditional image analysis methods fail for videos due to compression issues, prompting the authors to propose two deep learning networks with few layers to analyze key features, achieving over 98% accuracy for Deepfake and 95% for Face2Face.

Wodajo et al. [3] proposed a CViT for detecting Deepfakes, integrating a CNN with a ViT. The CNN extracts learnable features, which the ViT then processes using an attention mechanism for categorization. Trained on the DeepFake Detection Challenge (DFDC) dataset, their model achieves 95.8% accuracy, an AUC of 99.30, and a loss value of 0.32. The key contribution is the integration of a CNN module into the ViT architecture, resulting in competitive performance on the DFDC dataset. This combination leverages the strengths of both CNNs and ViTs, enhancing feature extraction and classification accuracy in Deepfake detection.

Ha et al. [10] introduced a robust DeepFake detection method that combines ViT and CNN models. Experiments showed that the ViT model excels at processing side faces and low-quality videos. The method, which integrates the ResNeSt269 model with the DeiT model using a weighted majority voting ensemble approach, achieved a 97.66% accuracy, surpassing the 96.78% accuracy of the current state-of-the-art model in the DFDC. Additionally, when tested on a completely different dataset, the method demonstrated robustness and over 10% higher accuracy compared to the CNN model, thanks to ViT's high generalization performance.

3. Materials and methodology

Face extraction using MTCNN and data augmentation are performed on the extracted face images. CViT combines CNNs for feature learning and ViTs for deep fake detection. CViT processes standardized face images (224×224 RGB), splitting them into patches for analysis. Utilizing the features from CNNs and ViTs, CViT accurately detects deepfake manipulation within face images. The entire process of the study is shown in **Figure 1**.

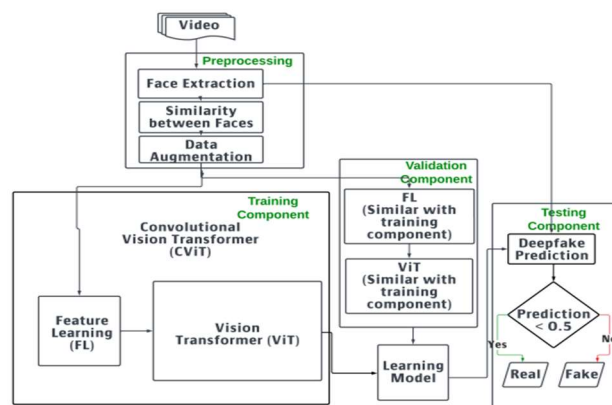


Figure 1. Methodology.

3.1. Dataset

FaceForensics++ [11] has 15,800 images extracted from 1,600 videos. The dataset is divided into training (72.38%, 11,448 images: 5835 fake, 5613 real), testing (19.62%, 3,103 images with a similar fake-real distribution), and validation (7.94%, 1252 images, evenly split between fake and real).

3.2. Preprocessing component

The preprocessing component plays a crucial role in preparing input data for the model. It consists of two key processes: face extraction and data augmentation. The face extraction component identifies and extracts faces from video frames, focusing the analysis on facial features. This step is vital, given that deepfakes often involve manipulations in this particular region. On the other hand, data augmentation enhances the model's ability to generalize by diversifying the training dataset. This involves applying random transformations like rotation, scaling, flipping, and sharpening to face images. To illustrate, the face extraction process outputs images in a standardized 224×224 RGB format. Simultaneously, data augmentation creates additional training samples with slightly modified versions of the original data. **Figure 2a** is an example of some of the frames. After obtaining the frames (224×224 RGB) as shown in **Figure 2b**, we calculated the facial region which is performed with the help of the MTCNN. After the face region has been obtained further processing and normalization are performed and the **Figure 3** are the images obtained after the normalization.



Figure 2. Frames and detection of face. (a) Frames in video; (b) Detection of face.



Figure 3. After normalization.

3.3. Multi-task cascaded convolutional neural networks

- The Multi-task Cascaded Convolutional Neural Network (MTCNN) algorithm used to detect face and face landmarks, works in three steps and uses one neural network for each process. The initial part is a proposal

network that will predict potential face positions and their bounding boxes just like an attention network in Faster R-CNN. The result of this process is a large number of face detection sandlots of false detections. The second part uses images and outputs of the first prediction, thus making a refinement of the result to eliminate most of the false detections and aggregate bounding boxes. The last part refines the predictions and adds facial landmarks predictions in the original MTCNN implementation. Experimental results have always demonstrated that while keeping the reliability of real-time performance, this method consistently outperforms the sophisticated conventional methods across most of the challenging benchmarks. This better performance for real-time is of great importance in a surveillance system [12]. The equations involved in the MTCNN algorithm are shown in Equations (1)–(3).

$$B = \text{sigmoid}(f_1(x, y, w, h)) \quad (1)$$

where B represents the bounding box coordinates, (x,y) are the coordinates of the top-left corner, (w, h) are the width and height of the bounding box, and f_1 is the neural network function.

$$\Delta B = f_2(B, I) \quad (2)$$

where ΔB represents the refined bounding box coordinates based on the initial bounding box B and the input image I, and f_2 is the neural network function.

$$(P, L) = f_3(B, I) \quad (3)$$

where P represents the facial landmarks and L represents the probability of the face being real, and f_3 is the neural network function.

3.4. Feature selection using CNN

Feature Learning (FL) is important in CNNs, especially for face recognition. It involves using blocks with convolutional layers to extract features from input face data. The features include the two eyes, nose, and the two sides of the mouth. These features are gradually learned and used as building blocks for higher-level analysis in the model. FL transforms raw data into meaningful representations, enabling more advanced processing in the neural network.

3.5. Vision transformer

The Vision Transformer (VT) within the CVIT framework is a key component that adapts transformer architecture, originally developed for natural language processing, to computer vision tasks. It processes learned features from the Feature Learning (FL) stage using self-attention, capturing global context information to understand relationships across different parts of the face. Following this, the MLP head, comprising fully connected layers and activation functions, refines these features for classification, distinguishing between real and fake inputs. The soft max function then assigns class probabilities based on raw scores from the MLP head, aiding in the final classification decision. The transformer encoder, meanwhile, handles linear projections of flattened patches from the Vision Transformer, refining features further by combining local and global information. These refined features are then fed into the MLP head for classification, enhancing the model's predictive accuracy and robustness. Additionally, the validation component assesses the model's performance on unseen data, incorporating FL and VT stages but operating

on a separate validation dataset to ensure an unbiased evaluation of the model's generalization capabilities.

3.6. Grad-CAM

Grad-CAM calculates the gradient of a differentiable output, such as class score, in relation to the convolutional features of a selected layer. Grad-CAM is most commonly employed for image classification tasks, but may also be utilized for semantic segmentation. The soft max layer of the proposed model outputs a score for each class for each pixel to aid in semantic segmentation. For a particular class C with N number of pixels and A^K as a feature map, Grad-CAM mapping is explained in Equation (4) [13].

$$M^c = ReLU \sum_K \alpha_c^K A^K \quad (4)$$

where,

$$\alpha_c^K = \frac{1}{N} \sum_{i,j} \left(\frac{dy^c}{dA_{i,j}^K} \right) \quad (5)$$

3.7. Real time implementation

During testing, the user uploads the video, and after face extraction, the extracted features are loaded with our CViT model. Leveraging the validated model, this component aims to predict the authenticity of new content, such as videos or frames, by determining whether they are real or fake. During testing, a predefined threshold of 0.5 is established to serve as the decision boundary. It gives the prediction score which if it is less than 0.5 is a real video otherwise it is a fake video.

4. Experiments and analysis

The experiment is conducted using Python language with Intel(R) Core (TM) i5-13500H CPU, windows 11 operating system, with 8 GB RAM.

Table 1. Results of videos along with prediction score (Samples are from the dataset).

SN	Sample Inputs	Prediction Score	Result
1	Sample 1.mp4	0.051	Real
2	Sample 2.mp4	0.67	Fake
3	Sample 3.mp4	0.03	Real
4	Sample 4.mp4	0.15	Real
5	Sample 5.mp4	0.96	Fake
6	Sample 6.mp4	0.26	Real
7	Sample 7.mp4	0.02	Real
8	Sample 8.mp4	0.9	Fake

The system evaluates the uploaded video and determines whether it's authentic or fake based on the prediction score from the model. If the score exceeds a threshold of 0.5, the video is flagged as fake otherwise, it's considered real, with the capability to predict whether a video is real or fake achieved. The above **Table 1**

comprises the file name, the prediction score through which we can predict whether the video is real or fake, and the prediction.

Throughout the training process, the training accuracy achieved by the model on this dataset was 92.5%. The accuracy graph is shown in **Figure 4** and training loss is shown in **Figure 5**. Additionally, the model's capacity was assessed using a Receiver Operating Characteristic (ROC) curve. The Area Under the Curve (AUC) value, representing the area covered by the ROC curve, was determined to be 0.91.

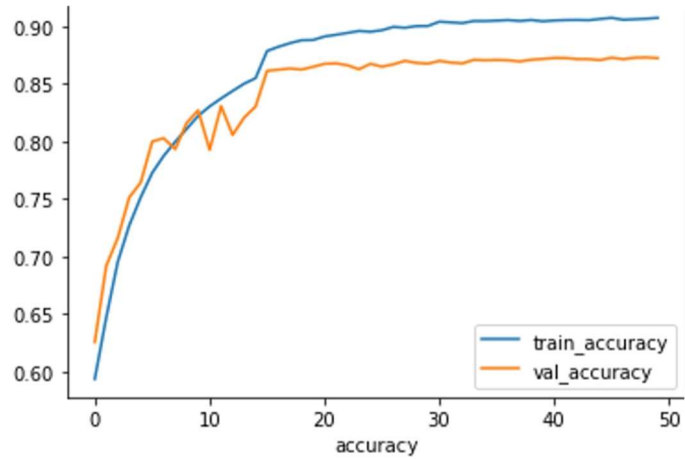


Figure 4. Training Accuracy.

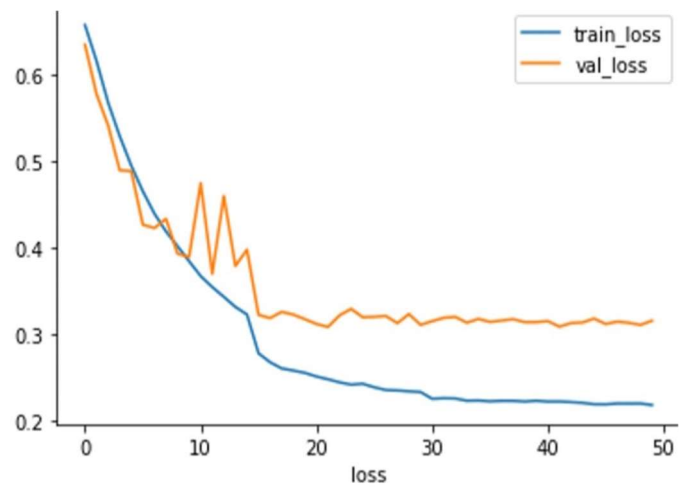


Figure 5. Training Loss.

A higher AUC suggests that the model has a strong ability to distinguish between the positive and negative classes. The ROC curve is shown in **Figure 6**. We conducted a 10-fold cross-validation and obtained an average accuracy of 92.5%, an average AUC of 0.91, an average precision of 0.91, and an average recall of 0.93. The K-fold result is shown in the **Table 2**.

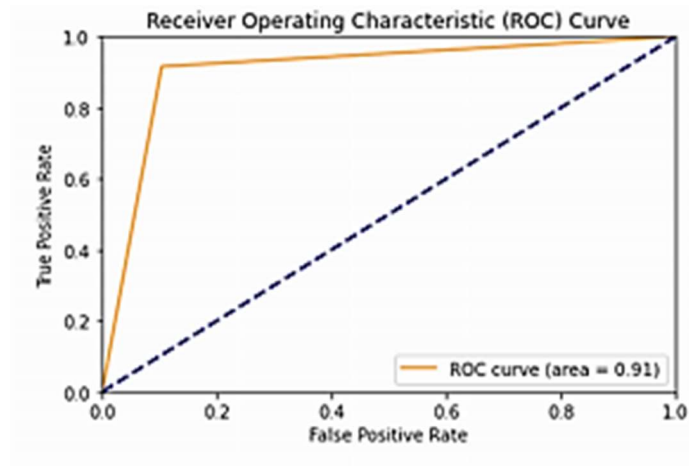


Figure 6. ROC curve.

Table 2. K-Fold cross-validation results.

Fold	Accuracy (%)	AUC	Precision	Recall
1	87.5	0.92	0.86	0.88
2	88.2	0.91	0.87	0.89
3	89.6	0.92	0.88	0.90
4	90.1	0.93	0.89	0.91
5	91.4	0.94	0.90	0.92
6	92.0	0.95	0.91	0.93
7	92.3	0.94	0.92	0.94
8	93.0	0.96	0.93	0.95
9	93.8	0.97	0.94	0.96
10	94.0	0.98	0.95	0.97
Avg	92.5	0.91	0.91	0.93

Grad-CAM is used to understand which parts of the input image are crucial for the deep learning model to determine whether an image or video is real or fake. For this we created a heatmap on the image to visualize the regions of interest, we can gain insights into how the model makes its decisions and potentially identify artifacts or inconsistencies indicative of manipulation. **Figure 7** shows the Grad-CAM over frames for fake video content.



Figure 7. Grad-CAM of fake image.

5. Conclusion

The study focuses on deepfake detection, employing a fusion of MTCNN architecture for feature extraction and Vision Transformer for video classification, which has yielded a noteworthy accuracy of 92.5% on the FaceForensics++ dataset, containing 15,808 images encompassing both genuine and fabricated instances. This outcome underscores the efficacy of our methodology.

To enhance future iterations, enlarging the dataset could bolster the model's capacity for generalization and resilience across diverse scenarios, potentially augmenting accuracy further. Moreover, integrating audio analysis alongside visual data offers a promising avenue for fortifying deepfake detection capabilities. By harnessing both visual and auditory cues, we can develop more comprehensive and dependable detection systems to counteract the escalating threat of media manipulation.

Author contributions: Conceptualization, SS, UK, SA, RG and MB; methodology, SS, UK, SA and RG; software, SS, UK, SA and RG; validation, SS, UK, SA and MB; formal analysis, MB; investigation, SS, UK, SA and RG; resources, SS, UK, SA and RG; data curation, SS, UK, SA and RG; writing—original draft preparation, SS, UK, SA and RG; writing—review and editing, MB; visualization, MB; supervision, MB; project administration, MB. All authors have read and agreed to the published version of the manuscript.

Conflict of interest: The authors declare no conflict of interest.

References

1. Karnouskos S. Artificial Intelligence in Digital Media: The Era of Deepfakes. *IEEE Transactions on Technology and Society*. 2020; 1(3): 138-147. doi: 10.1109/tts.2020.3001312
2. Grobler GD. Narrative strategies in the creation of animated poetry-film [PhD thesis]. University of South Africa; 2021.
3. Wodajo D, Atnafu S, Akhtar Z. Deepfake video detection using generative convolutional vision transformer. Available online: <https://arxiv.org/abs/2307.07036> (accessed on 20 May 2024).
4. Heidari A, Jafari Navimipour N, Dag H, et al. Deepfake detection using deep learning methods: A systematic and comprehensive review. *WIREs Data Mining and Knowledge Discovery*. 2023; 14(2). doi: 10.1002/widm.1520
5. Kearns L, Alam A, Allison J. Synthetic media authentication threats: Detection using a combination of neural network and blockchain technology. Available online: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4658121 (accessed on 20 May 2024).
6. Chesney R, Citron DK. Deep Fakes: A Looming Challenge for Privacy, Democracy, and National Security. *SSRN Electronic Journal*. 2018. doi: 10.2139/ssrn.3213954
7. Masood M, Nawaz M, Malik KM, et al. Deepfakes generation and detection: state-of-the-art, open challenges, countermeasures, and way forward. *Applied Intelligence*. 2022; 53(4): 3974-4026. doi: 10.1007/s10489-022-03766-z
8. Montserrat DM, Hao H, Yarlagadda SK, et al. Deepfakes Detection with Automatic Face Weighting. In: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW); 2020. doi: 10.1109/cvprw50498.2020.00342
9. Afchar D, Nozick V, Yamagishi J, et al. MesoNet: a Compact Facial Video Forgery Detection Network. In: 2018 IEEE International Workshop on Information Forensics and Security (WIFS); 2018. doi: 10.1109/wifs.2018.8630761
10. Ha H, Kim M, Han S, et al. Robust Deep Fake Detection Method based on Ensemble of ViT and CNN. In: Proceedings of the 38th ACM/SIGAPP Symposium on Applied Computing; 2023. doi: 10.1145/3555776.3577769
11. Hasan FS. FaceForensics-1600 videos-preprocess. Available online: <https://www.kaggle.com/datasets/farhansharukhhasan/faceforensics1600-videospreprocess?rvi=1> (accessed on 23 May 2024).

12. Jose EMG, Haridas MTP, Supriya MH. Face Recognition based Surveillance System Using FaceNet and MTCNN on Jetson TX2. 2019 5th International Conference on Advanced Computing & Communication Systems (ICACCS). Published online March 2019. doi: 10.1109/icaccs.2019.8728466
13. Selvaraju RR, Cogswell M, Das A, et al. Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization. In: 2017 IEEE International Conference on Computer Vision (ICCV); 2017; Venice, Italy. pp. 618-626. doi: 10.1109/iccv.2017.74

Article

Software cost estimation tool: A App based application, estimate the cost of software project

Ajay Jaiswal*, Piyush Malviya, Lucky Parihar, Rani Pathak, Kuldeep Rajput

Computer Science & Engineering. Department, Prestige Institute of Engineering Management and Research, Indore 452010, India

* Corresponding author: Ajay Jaiswal, ajay.jaiswal55555@gmail.com

CITATION

Jaiswal A, Malviya P, Parihar L, et al. Software cost estimation tool: A App based application, estimate the cost of software project. *Computing and Artificial Intelligence*. 2024; 2(2): 1364.
<https://doi.org/10.59400/cai.v2i2.1364>

ARTICLE INFO

Received: 7 May 2024

Accepted: 4 July 2024

Available online: 22 July 2024

COPYRIGHT



Copyright © 2024 by author(s).
Computing and Artificial Intelligence is published by Academic Publishing Pte. Ltd. This work is licensed under the Creative Commons Attribution (CC BY) license.
<https://creativecommons.org/licenses/by/4.0/>

Abstract: This paper presents the design and implementation of a software cost estimation tool integrated into a mobile application developed using Flutter. The tool incorporates various techniques for software cost estimation, including expert judgment, function point analysis, 3D point analysis, and the COCOMO model. The purpose of the program is to give software engineers and project managers a practical and effective tool for calculating the time and money needed for software development projects. The paper provides a thorough explanation of each estimation technique's implementation, along with a discussion of the app's main features and functionalities. Because of the app's intuitive and user-friendly design, users can quickly enter project data and get precise cost estimates. The tool's efficacy is assessed using case studies and contrasts with other software cost estimation methods currently in use. The outcomes show that the app can produce trustworthy and precise cost estimates, which makes it an important resource for software development projects.

Keywords: software cost estimation; flutter app development; project management tools; function point analysis; COCOMO model

1. Introduction

Software development projects are renowned for their complexity, requiring meticulous planning and precise budgeting to ensure successful completion. One of the most critical aspects of project planning is accurate cost estimation, which involves predicting the resources, time, and effort required to deliver a project. However, traditional cost estimation methods often fall short, leading to budget overruns and project delays. To address this challenge, we introduce “Software Cost Estimation”, a groundbreaking tool designed to revolutionize software project budgeting and planning. Our platform gives users precise and trustworthy cost estimates so they can optimize their project budgets and make well-informed decisions. It does this by utilizing the power of state-of-the-art algorithms, historical project data, and industry best practices.

In this research paper, we provide a comprehensive overview of the design, development, and evaluation of “Software Cost Estimation”. We discuss the tool's key features and functionalities, including its ability to analyze project requirements [1], estimate costs, and generate detailed reports. Furthermore, we present the results of empirical studies and case studies that demonstrate the effectiveness and accuracy of our tool in real-world software projects.

By introducing “Software Cost Estimation” [2], we aim to empower software development teams and project managers with a powerful tool that can streamline the cost estimation process, reduce budget uncertainties, and improve overall project

planning. We believe that our tool has the potential to significantly impact the software development industry, leading to more efficient and cost-effective project delivery.

Background:

Software cost estimation is a challenging and crucial aspect of project management. Traditional estimation methods, such as expert judgment and analogy-based estimation, often rely on subjective assessments and historical data that may not accurately reflect the complexities of modern software projects. As a result, these methods can lead to inaccurate estimates, which can have significant implications for project budgets and schedules. The development of automated cost estimation tools that use cutting-edge algorithms and machine learning approaches to increase the precision and dependability of cost estimates has attracted increasing attention in recent years. Based on past data and industry norms, these tools examine a variety of project criteria, including size, complexity, and requirements [3], to produce estimates that are more accurate.

Objectives:

The primary objective of this research paper is to introduce “Software Cost Estimation” and demonstrate its effectiveness in improving software project cost estimation. We aim to showcase the key features and functionalities of the tool, highlight its advantages over traditional estimation methods, and present empirical evidence supporting its accuracy and reliability. Additionally, we seek to compare “Software Cost Estimation” with existing cost estimation approaches to highlight its unique capabilities and potential impact on software project management practices. Through this research, we hope to contribute to the advancement of software cost estimation techniques and provide software development teams with a valuable tool for optimizing their project budgets and schedules.

1.1. Cost estimation technique

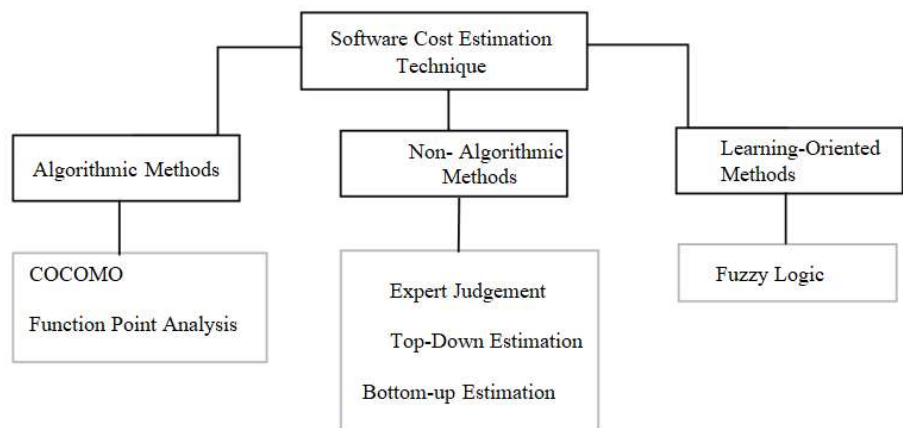


Figure 1. Software cost estimation techniques.

Figure 1 shows the cost estimation of software can be approached in several ways. One of the key stages in developing new software is figuring out its cost, which involves estimating the time required, necessary resources, and the project’s overall size. Research indicates that estimates for software projects can be off by up to 40%. These methods [4] can be broadly categorized into algorithmic and non-algorithmic

approaches. In this section, we will explore various methods, outlining their advantages and disadvantages to help you decide which approach is most suitable. Estimating the cost of a software project is a critical aspect of software engineering, often determining the success or failure of a project or business deal. Throughout the software development life cycle, accurately predicting the necessary work and its associated cost is a primary focus of software cost estimation.

1.1.1. Algorithmic methods

Equations based [5] on empirical data and in-depth research are essential to the pursuit of precise software cost estimation. These equations are designed to take use of several inputs, including functional requirements, Source Lines of Code (SLOC), design process, team experience, and risk assessments, among other factors that influence costs. Among the many computational models that have shown to be invaluable in this trial are COCOMO and function point analysis, to name just two. In order to convert project parameters into quantitative estimations, they provide methodical frameworks. These techniques attempt to aid in the calculation of software costs by providing a few mathematical formulas.

Constructive Cost Model (COCOMO): The Constructive Cost Model (COCOMO), developed by Barry Boehm in the late 1970s, is a seminal method for estimating the effort, time, and cost required for software development projects. Grounded in the belief that software development effort is influenced by various factors, COCOMO provides a structured framework for assessing these factors and deriving accurate estimates. Central to COCOMO's approach is the recognition that the size of the software product and the characteristics of the development environment significantly impact project outcomes.

COCOMO offers three distinct models tailored to different stages of project maturity and complexity [6]: Basic COCOMO, Intermediate COCOMO, and Detailed COCOMO.

Basic COCOMO, the initial model in the series, is particularly suitable for early-stage project planning when only limited information about the software product and project environment is available. This model estimates effort as a function of software size, typically measured in thousands of lines of code (KLOC), and incorporates a set of cost drivers that capture various project attributes such as complexity, personnel capability, and development tools. The formula for Basic COCOMO is represented as:

$$\text{Effort} = a \times (\text{KLOC})^b$$

where a and b are constants empirically derived from historical data and represent the scale and exponent factors respectively.

Intermediate COCOMO extends the capabilities of Basic COCOMO by incorporating additional project-specific factors into the estimation process. In addition to software size, Intermediate COCOMO considers parameters such as development flexibility, team cohesion, and risk resolution capabilities. The estimation formula for Intermediate COCOMO introduces an Effort Adjustment Factor (EAF), which serves as a multiplier reflecting the combined effects of all cost drivers:

$$\text{Effort} = \text{EAF} \times (\text{KLOC})^b$$

The Effort Adjustment Factor (EAF) is determined based on a comprehensive

assessment of various project attributes and is crucial in refining the estimation accuracy.

Detailed COCOMO represents the most comprehensive and sophisticated version of the model, suitable for large-scale and complex software development projects. In addition [7], to the factors considered in Intermediate COCOMO, Detailed COCOMO incorporates detailed assessments of personnel experience, software reliability requirements, and product complexity.

COCOMO's enduring relevance in software project management stems from its ability to provide a structured and systematic approach to estimating development effort and resource requirements. However, it is important to acknowledge that COCOMO is not without limitations. Its reliance on historical data and assumptions about project characteristics may introduce uncertainties, particularly in rapidly evolving technological landscapes. As such, ongoing refinement and validation of COCOMO estimates based on real-world data and project experience are essential for enhancing its effectiveness and reliability.

Function point analysis:

Function point analysis (FPA) [8] is a widely recognized and systematic method for estimating the size and complexity of software systems based on the functionalities they deliver to users. Introduced by Allan Albrecht in the late 1970s, FPA focuses on quantifying the functional requirements of a software product, independent of technology or implementation details. The core concept of FPA [9] revolves around identifying and categorizing functional components within a software system, such as inputs, outputs, inquiries, internal data files, and external interfaces. By assigning weights to each functional component based on its complexity and significance, FPA enables the computation of a function point (FP) metric, which serves as a standardized measure of software size. The formula for calculating Function Points typically involves summing the weighted values of individual functional components:

$$FP = \sum_{i=0}^n (\text{Weight}_i \times \text{Count}_i)$$

where Weight_i represents the complexity weight assigned to each functional component, Count_i denotes the count of occurrences of that component, and n is the total number of functional components considered. FPA [10,11] offers a holistic perspective on software size and complexity, capturing both internal and external aspects of system functionality. This makes it a valuable technique for estimating development effort, resource requirements [12], and project duration. However, like any estimation method, FPA [13] requires careful application and consideration of contextual factors to ensure accurate and reliable results.

1.1.2. NON-algorithmic techniques

Precise cost estimation is critical to software development project planning, budgeting, and resource allocation. Non-algorithmic approaches mostly rely on expert judgment and qualitative evaluations, whereas algorithmic approaches use mathematical models and historical data to estimate expenses. This section explores the types, formulas, and applications of several non-algorithmic cost estimation strategies.

Expert judgment:

Expert judgment stands [14] as a cornerstone in software cost estimation, drawing

upon the insights and experiences of seasoned professionals in the field. This technique leverages the collective wisdom of experts who possess domain knowledge, project management expertise, and a nuanced understanding of the software development lifecycle. Through deliberative discussions, brainstorming sessions, and peer reviews, experts offer informed opinions on cost drivers, project complexities, and resource requirements.

Types of expert judgment:

Delphi technique: The Delphi technique fosters consensus among a panel of experts through iterative rounds of anonymous feedback and controlled communication. Experts individually provide estimates, which are aggregated and refined in subsequent rounds until convergence is achieved. This method mitigates biases and encourages diverse viewpoints, thereby enhancing the accuracy of cost estimates.

Analogous-Based estimation: Analogous estimation draws parallels between the current project and past trials, extrapolating costs based on similarities in scope, size, and technological complexity. By referencing historical data and benchmarking against analogous projects, experts can derive preliminary cost estimates, often expressed as a percentage deviation from past efforts.

Estimation techniques:

Expert-based cost estimates are dependent on the projects in which they were utilized since they represent the knowledge of the experts who were consulted. Data gathering and discovery may be impeded in several commonplace scenarios. In these situations, the “expert judgment” method is effective. It is the accepted technique [15] for estimating the duration of a software project. One method for utilizing expert opinion in cost estimation is the Wideband Delphi Method. These people are subject to two rounds of evaluation. The work breakdown structure is an additional example of expert opinion.

Top-down estimating method:

The term “Macro Model” is often used to refer to the top-down estimation technique it describes. Using this technique, the overall software project cost estimate is determined from the project’s global attributes, and then the project is broken down into its constituent low-level mechanisms or components. The Putnam model is a technique that takes this perspective. For preliminary cost calculation when only global parameters are available, the Top-Down approach is preferable. Due to a lack of specifics at the outset, top-down approaches are ideal for estimating software costs.

Bottom-up estimating method:

A predicted total project cost is then calculated by adding the individual product costs determined using the base-up costing method. The goal of a bottom-up approach is to build a framework’s gauge from data gathered about its constituent parts and how they interact. The point-by-point model used by COCOMO is the technique using this approach.

2. Literature review

Software cost estimation has been a longstanding challenge in the field of software engineering, with researchers and practitioners continually seeking to

improve the accuracy and reliability of cost estimation methods. While expert judgment and analogy-based estimating are two popular traditional cost estimation methodologies, they are frequently prone to errors since they rely on subjective assessments and historical data that might not fully reflect the current project context. The application of machine learning techniques and quantitative models for software cost assessment has gained popularity in recent years. These methods improve the accuracy of project cost predictions by utilizing project features, historical project data, and other considerations. One well-known quantitative model that estimates the time and money needed for software development is the Constructive Cost Model (COCOMO), which takes into account the size, complexity, and other aspects of the project.

Machine learning techniques, such as regression analysis, decision trees, and neural networks, have also been applied to software cost estimation with promising results. These techniques can learn from past project data and adjust their predictions based on new information, improving the accuracy of cost estimates over time. Another area of research in software cost estimation is the use of parametric estimation models, which estimate project costs based on a set of predefined parameters. These models can be customized to fit the specific characteristics of a project, making them potentially more accurate than generic estimation approaches.

Despite these advancements, challenges remain in software cost estimation, particularly in the context of agile and iterative development methodologies. Agile projects [16] are characterized by their dynamic nature, frequent changes, and evolving requirements, making traditional cost estimation methods less suitable. Researchers are exploring new approaches that can adapt to the iterative nature of agile development and provide more accurate cost estimates in such environments. Overall, the literature suggests that while significant progress has been made in software cost estimation, there is still room for improvement. New technologies, such as machine learning and agile development methodologies, are reshaping the landscape of cost estimation, offering new opportunities to enhance the accuracy and reliability of cost estimates in software projects.

3. Methodology

The development process of the ‘Software Cost Estimation Tool’ using Flutter involved several key steps. Initially, a new Flutter project was set up, and the necessary dependencies were configured. The user interface was designed to accommodate the various features of the app, focusing on simplicity and usability. Each feature, including expert judgment, analogous estimation, parametric estimation, 3D point estimation, COCOMO model [17], and function point analysis, was implemented using Flutter widgets and libraries. For expert judgment, a user-friendly interface was created for users to input their estimates based on their expertise. Analogous estimation utilized historical data from similar projects, requiring integration with a database or API. Parametric estimation involves implementing algorithms to calculate estimates based on project parameters.

Implementing 3D point estimation was complex, requiring the development of algorithms for more accurate cost estimates. The COCOMO model was integrated into

the app to estimate costs based on project size and complexity, involving the implementation of COCOMO equations and algorithms. Function point analysis was implemented using algorithms to calculate function points and estimate project size and effort based on functionality.

During development, challenges were encountered, such as technical limitations of Flutter, especially regarding performance and compatibility. These challenges were addressed through code optimization and the use of alternative approaches. Implementing complex features, like 3D point estimation and mathematical models, required breaking down the implementation into smaller tasks and seeking expert advice when needed. Ensuring a smooth user experience, particularly with features like expert judgment and input validation, was achieved through thorough testing and user feedback incorporation.

4. Case study

E-commerce optimizer software development at XYZ corporation:

Background: XYZ Corporation, a mid-sized e-commerce company, aimed to enhance its operational efficiency and customer experience through the development of a comprehensive software solution. Named “E-Commerce Optimizer”, the software aimed to integrate various functionalities such as inventory management, order processing, and customer relationship management into a unified platform.

Project scope:

The project involved the following key objectives:

- Development of a user-friendly interface facilitating inventory management, order processing, and customer interactions.
- Integration of the software with existing systems and databases to ensure seamless data flow. Implementation of analytics features for monitoring sales, customer behavior, and inventory levels.
- Ensuring scalability and security to accommodate future growth and protect sensitive data.

Methodology:

XYZ corporation adopted the Agile methodology for its flexibility and adaptability. The project was divided into iterative sprints, each focusing on specific features or functionalities. Regular meetings were conducted to review progress, gather feedback, and adjust plans accordingly.

Cost estimation:

The cost estimation process involved a combination of bottom-up and top-down approaches. Task breakdowns were used to estimate time and resources required for each component of the project. Additionally, industry benchmarks and past projects were analyzed to validate estimates and identify potential cost-saving opportunities.

Result:

Based on the cost estimation process, XYZ Corporation projected the development cost for E-Commerce Optimizer to be approximately \$500,000. This estimation encompassed expenses related to software development, testing, infrastructure setup, and project management.

Conclusion:

By employing robust cost estimation techniques and adhering to Agile principles, XYZ Corporation successfully developed E-Commerce Optimizer within budget and timeline constraints. The software's implementation led to notable improvements in operational efficiency, customer satisfaction, and overall business performance.

5. Results and discussion

The 'Software Cost Estimation Tool app demonstrated its effectiveness in providing accurate and efficient cost estimates for software projects. By incorporating features such as expert judgment, analogous estimation, parametric estimation, 3D point estimation, COCOMO model, and function point analysis, the app was able to offer a comprehensive approach to cost estimation.

In a case study [18] comparing the app's estimates with those from traditional cost estimation methods, the app consistently provided estimates that were close to the actual costs of software projects. This indicates that the app's algorithms and models are reliable and can be used with confidence by project managers and software developers.

5.1. Comparison with traditional methods

Compared to traditional cost estimation methods, the Software Cost Estimation Tool app offers several advantages. Traditional methods often rely on manual calculations and subjective judgments, which can lead to inaccuracies and inconsistencies in estimates. In contrast, the app uses algorithms and mathematical models to provide more objective and reliable estimates. Additionally, the app's ability to incorporate historical data and project parameters allows for more precise estimates, taking into account the specific characteristics of each project. This can result in more accurate budgeting and resource allocation, leading to better project management and decision-making.

Advantages

- Accuracy: The app provides accurate cost estimates [19] based on algorithms and mathematical models, reducing the risk of budget overruns.
- Efficiency: The app streamlines the cost estimation process, saving time and effort for project managers and software developers.
- Comprehensiveness: By incorporating multiple estimation methods, the app offers a comprehensive approach to cost estimation, ensuring that all relevant factors are considered.
- User-friendly: The app's user-friendly interface makes it easy for users to input data and generate cost estimates, even without specialized knowledge in cost estimation techniques.

5.2. Limitations

- Dependency on data: The accuracy of the app's estimates depends on the quality and relevance of the data used. Inaccurate or outdated data can lead to unreliable estimates.
- Complexity: Some features of the app, such as 3D point estimation and the COCOMO model, may be complex for users without a background in software

cost estimation.

- Technical limitations: The app's performance and accuracy may be affected by technical limitations, such as hardware capabilities and network connectivity.
- In conclusion, the Software Cost Estimation Tool app offers a reliable and efficient solution for software cost estimation, providing accurate estimates that can help project managers and software developers plan and manage their projects more effectively. While the app has some limitations, its advantages make it a valuable tool for cost estimation in software development projects.

6. Conclusion

The research paper presents the development and implementation of the 'Software Cost Estimation Tool' using Flutter. Key findings include the successful integration of various cost estimation techniques such as expert judgment, analogous estimation, parametric estimation, 3D point estimation, COCOMO model, and function point analysis into a user-friendly mobile application. Through careful design and implementation, the app provides software development teams with a comprehensive tool for estimating project costs accurately and efficiently.

Implications for software development projects:

The 'Software Cost Estimation Tool app holds significant implications for software development projects. By providing a centralized platform for cost estimation, the app empowers project managers and stakeholders to make informed decisions regarding resource allocation, budgeting, and project planning. It enhances project transparency and accountability by enabling teams to track and manage costs effectively throughout the development lifecycle. Additionally, the app promotes collaboration and communication among team members, facilitating a more streamlined and efficient development process.

Future research directions:

While the 'Software Cost Estimation Tool' app represents a significant advancement in software cost estimation, there are several avenues for future research and improvement. Firstly, enhancing the accuracy and reliability of estimation techniques, such as 3D point estimation and parametric estimation, could lead to more precise cost predictions. Exploring advanced machine learning algorithms and data analytics techniques may also offer insights into optimizing cost estimation models based on real-time project data. Furthermore, integrating additional features such as risk analysis [19,20], resource optimization, and project scheduling could further enhance the app's functionality and utility. Collaborating with industry experts and practitioners to validate and refine the app's algorithms and methodologies could ensure its relevance and effectiveness in real-world software development scenarios.

In conclusion, the 'Software Cost Estimation Tool' app represents a valuable contribution to software development practices, offering a comprehensive solution for estimating and managing project costs. Continual research and development efforts are essential to further enhance the app's capabilities and address evolving challenges in the software development landscape.

Author contributions: Conceptualization, AJ and LP; methodology, PM; software,

AJ; validation, AJ, LP, PM and RP; formal analysis, AJ; investigation, AJ; resources, KR; data curation, AJ; writing—original draft preparation, AJ; writing—review and editing, PM; visualization, AJ; supervision, AJ; project administration, LP; funding acquisition, PM. All authors have read and agreed to the published version of the manuscript.

Conflict of interest: The authors declare no conflict of interest.

References

1. Firesmith D. Prioritizing Requirements. *The Journal of Object Technology*. 2004; 3(8): 35. doi: 10.5381/jot.2004.3.8.c4
2. Balaji N, Shivakumar N, Ananth VV. Software cost estimation using function point with non-algorithmic approach. *Global Journal of Computer Science and Technology Software & Data Engineering*. 2013; 13(8): 1–4.
3. Karlsson J. Software Requirements Prioritizing. In: *Proceedings of the International Conference on Requirement Engineering*; 1996.
4. Hamdan K, El Khatib H, Shuaib K. Practical software project total cost estimation methods. In: *Proceedings of 2010 International Conference on Multimedia Computing and Information Technology (MCIT)*; 2010. doi: 10.1109/mcit.2010.5444853
5. Khan B, Khan W, Arshad M, Jan N. Software cost estimation: Algorithmic and non-algorithmic approaches. *International Journal of Data Science and Advanced Analytics*. 2020; 2(2): 1–5.
6. Zuse H. Software Metrics-Methods to Investigate and Evaluate Software Complexity Measures. In: *Proceedings of the Second Annual Oregon Workshop on Software Metrics*; 1991; Portland.
7. Kitchenham B, Mendes E. Software productivity measurement using multiple size measures. *IEEE Transactions on Software Engineering*. 2004; 30(12): 1023–1035. doi: 10.1109/tse.2004.104
8. Westerville. *Function Point Counting Practices Manual*. International Function Point User Group (IFPUG); 1990.
9. Low GC, Jeffery DR. Function points in the estimation and evaluation of the software process. *IEEE Transactions on Software Engineering*. 1990; 16(1): 64–71. doi: 10.1109/32.44364
10. Meli R, Santillo L. *Function point estimation methods: A comparative overview*. Data Processing Organization; 1999.
11. Meli R, Satillo L. *Function Point Measurement Tool for UML Design Specification*. Data Processing Organization; 1999.
12. Sadiq M, Ghafir S, Shahid M. A Framework to Prioritize the software Requirements using Quality Function Deployment. In: *Proceedings of the National Conference on Recent Development in Computing and its Application*; 2009; Delhi, India.
13. Symons CR. Function point analysis: difficulties and improvements. *IEEE Transactions on Software Engineering*. 1988; 14(1): 2–11. doi: 10.1109/32.4618
14. Mansor ZB, Kasirun ZM, Arshad NHH, et al. E-cost estimation using expert judgment and COCOMO II. In: *Proceedings of 2010 International Symposium on Information Technology*; 2010. doi: 10.1109/itsim.2010.5561466
15. Boehm BW. *Software Engineering Economics*. Prentice Hall; 1981.
16. Alliance A. *Agile Methodologies*. Available online: <https://www.agilealliance.org/agile101/agile-methodologies/> (accessed on 13 March 2024).
17. Rush C, Roy R. Expert Judgement in Cost Estimating: Modelling the Reasoning Process. *Concurrent Engineering*. 2001; 9(4): 271–284. doi: 10.1177/1063293x0100900404
18. Chirra SMR, Reza H. A Survey on Software Cost Estimation Techniques. *Journal of Software Engineering and Applications*. 2019; 12(06): 226–248. doi: 10.4236/jsea.2019.126014
19. Gupta D, Sadiq M. Software Risk Assessment and Estimation Model. In: *Proceedings of 2008 International Conference on Computer Science and Information Technology*; 2008. doi: 10.1109/iccsit.2008.184
20. Hoodat H, Rashidi H. Classification and Analysis of Risks in Software Engineering. *World Academy of Science, Engineering and Technology*. 2009; 56: 446–452.

Article

Enhancing user experience and trust in advanced LLM-based conversational agents

Yuanyuan Xu^{1,*}, Weiting Gao², Yining Wang³, Xinyang Shan¹, Yin-Shan Lin⁴

¹ Tongji University, Shanghai 200092, China

² Amazon, Seattle, WA 98121, USA

³ Bentley University, Waltham, MA 02452, USA

⁴ Northeastern University, Boston, MA 02115, USA

* **Corresponding author:** Yuanyuan Xu, ecusttethys@foxmail.com

CITATION

Xu Y, Gao W, Wang Y, et al.
Enhancing user experience and trust in advanced LLM-based conversational agents. *Computing and Artificial Intelligence*. 2024; 2(2): 1467.
<https://doi.org/10.59400/cai.v2i2.1467>

ARTICLE INFO

Received: 24 June 2024

Accepted: 6 August 2024

Available online: 17 August 2024

COPYRIGHT



Copyright © 2024 by author(s).
Computing and Artificial Intelligence is published by Academic Publishing Pte. Ltd. This work is licensed under the Creative Commons Attribution (CC BY) license.
<https://creativecommons.org/licenses/by/4.0/>

Abstract: This study explores the enhancement of user experience (UX) and trust in advanced Large Language Model (LLM)-based conversational agents such as ChatGPT. The research involves a controlled experiment comparing participants using an LLM interface with those using a traditional messaging app with a human consultant. The results indicate that LLM-based agents offer higher satisfaction and lower cognitive load, demonstrating the potential for LLMs to revolutionize various applications from customer service to healthcare consultancy and shopping assistance. Despite these positive findings, the study also highlights significant concerns regarding transparency and data security. Participants expressed a need for clearer understanding of how LLMs process information and make decisions. The perceived opacity of these processes can hinder user trust, especially in sensitive applications such as healthcare. Additionally, robust data protection measures are crucial to ensure user privacy and foster trust in these systems. To address these issues, future research and development should focus on enhancing the transparency of LLM operations and strengthening data security protocols. Providing users with clear explanations of how their data is used and how decisions are made can build greater trust. Moreover, specialized applications may require tailored solutions to meet specific user expectations and regulatory requirements. In conclusion, while LLM-based conversational agents have demonstrated substantial advantages in improving user experience, addressing transparency and security concerns is essential for their broader acceptance and effective deployment. By focusing on these areas, developers can create more trustworthy and user-friendly AI systems, paving the way for their integration into diverse fields and everyday use.

Keywords: large language models (LLMs); user experience (UX); conversational agents; transparency; data security

1. Introduction

Large language models (LLMs) represent a significant advancement in artificial intelligence designed to comprehend and generate text that closely mimics human communication [1]. Initially, these models functioned as basic natural language processing (NLP) tools focusing primarily on parsing and interpreting text [2]. However, over the years, LLMs have undergone substantial evolution, transforming into sophisticated conversational agents. Unlike traditional chatbots which operate on predefined scripts and limited response patterns, LLM-based agents leverage deep learning algorithms to facilitate more nuanced and contextually appropriate interactions [3]. This capability allows them to understand subtle cues and provide responses that are not only relevant but also contextually aware, significantly

enhancing the quality of user interactions [4].

The applications of LLM-based conversational agents, exemplified by platforms like ChatGPT and Gemini, are diverse and far-reaching [5]. These agents have been deployed in various sectors including customer service, where they handle inquiries and provide support; healthcare consultancy, offering preliminary advice and information; and shopping assistance, helping users find products and make purchasing decisions [6]. Their ability to process and respond to a broad spectrum of queries makes them exceptionally versatile tools capable of adapting to different use cases and user needs.

For these conversational agents to be successfully adopted and integrated into everyday use, user experience (UX) and user trust are paramount [7,8]. UX encompasses the overall experience users have while interacting with the system, including ease of use, efficiency, and satisfaction. A positive UX ensures that users find the interaction seamless and enjoyable, encouraging continued use. On the other hand, user trust is built on the reliability, transparency, and security of the conversational agent. Reliability refers to the consistent performance and accuracy of the agent's responses; transparency involves clear communication about how the agent operates and processes information; and security pertains to the protection of user data and privacy [9–11]. Together, these factors form the foundation for the effective deployment of LLM-based conversational agents, ensuring that users feel confident and secure in their interactions.

Research involving controlled experiments has shown that LLM-based agents can offer higher satisfaction and lower cognitive load compared to traditional messaging apps with human consultants [12,13]. Participants in these studies reported that interactions with LLMs felt more seamless and efficient as these agents were able to provide immediate, contextually appropriate responses without the delays often associated with human-mediated communication. However, the studies also revealed significant concerns regarding transparency and data security [14–16]. Users expressed uncertainty about how their data is being used and stored, raising important ethical and practical considerations [17–19]. These concerns underscore the need for better communication and robust data protection measures to ensure user privacy [20–22].

This study aims to enhance user experience (UX) and trust in advanced Large Language Model (LLM)-based conversational agents. We conducted a controlled experiment comparing participants using an LLM interface with those using a traditional messaging app with a human consultant. The primary objectives were to assess user satisfaction, task completion efficiency, and cognitive load across different tasks, such as weather inquiries, schedule management, technical support, and health consultations.

The LLM used in this study is based on OpenAI's GPT-3.5, featuring a sophisticated architecture and trained on a vast dataset to ensure high-quality performance. Participants were divided into two groups, and their interactions with the LLM-based agent and the human consultant were recorded and analyzed using both quantitative and qualitative methods.

Our findings indicate that LLM-based agents offer significant advantages in terms of user satisfaction and efficiency. However, issues related to transparency and data

security remain critical for broader acceptance. The study provides detailed insights into these aspects and offers recommendations for future research and development to enhance the transparency and security of LLM-based systems.

2. Literature review

2.1. Research on user experience (UX) in AI and LLM Interactions

Research on user experience in AI emphasizes intuitive design and user satisfaction. Positive UX leads to increased engagement and better outcomes [23]. In LLM interactions, UX directly impacts user willingness to use the system and overall satisfaction [24]. Key factors include interface intuitiveness, response speed, and interaction naturalness [25]. These determine task efficiency and user enjoyment [26].

Ease of use and response timeliness are critical when interacting with LLMs [27]. Users expect natural, seamless conversations requiring LLMs to understand complex contexts and cues, while ensuring fast and accurate responses [28]. Privacy and data security concerns also significantly affect trust and willingness to use LLMs.

To enhance UX, research suggests focusing on simple and intuitive interface design, efficient response times, and natural conversation processes [29]. Transparent information processing and robust data protection measures are also essential to increase user trust [30].

2.2. Research on user trust in AI and user perceptions and expectations of AI systems

User trust in AI focuses on reliability, transparency, and security [31]. Reliability involves consistent performance, accuracy, and availability [32,33]. Transparency requires clear communication about system operations and decision-making processes [34]. Security focuses on protecting user data and ensuring safe interactions [35].

User perceptions and expectations significantly influence acceptance and trust in AI systems. Users expect accurate, reliable services that transparently demonstrate their working principles and protect personal data [36]. Understanding data processing and clear explanations for decisions are crucial, especially in sensitive areas like healthcare and finance [37–39].

2.3. Research on UX and trust in specialized applications

Research highlights the importance of UX and trust in specialized LLM applications like healthcare consultancy and shopping assistance. These domains have unique requirements; healthcare prioritizes accuracy and privacy, while shopping emphasizes smooth and fast interactions [40–42].

Ease of use, response timeliness, and natural conversation are key UX factors for LLMs [43]. User's privacy and data security concerns significantly influence their trust and willingness to use LLMs [44–46]. Interface design should be simple and intuitive, ensuring efficient response times and natural conversation flow [47,48]. Transparent information processing and robust data protection measures are essential to build user trust [49].

Despite substantial research on general AI applications, there is a notable gap in

focused studies on specialized LLM applications [50]. Tailored research in healthcare and shopping can improve LLM design and deployment [51,52].

In conclusion, UX and user trust are critical for the successful deployment of LLM-based conversational agents[53]. Further exploration in specialized applications is needed. Addressing transparency and security issues will improve UX and build greater user trust [54]. By focusing on the unique needs of different domains, developers can create more effective and trusted AI systems, enhancing the overall impact and usability of LLM technology [53].

3. Methods

3.1. Participants

The study used a controlled design where participants were divided into two groups: one group used a Large Language Model (LLM) interface and the other group used a messaging application interface with a human advisor. The study included 18 participants.

We initially screened the participants through a background questionnaire, which was completed by 151 people. Based on the responses, we selected a representative subset to participate in subsequent experiments to ensure sample diversity. This approach ensured that the final participants were sufficiently representative in terms of age, gender, educational background, and technological familiarity (**Table 1**).

Table 1. Participant demographic information.

Participant ID	Age	Gender	Educational Background	Technical Familiarity
1	25	Male	Bachelor's	Medium
2	34	Female	Master's	High
3	22	Male	Associate's	Low
4	29	Female	Bachelor's	High
5	41	Male	High School	Medium
6	37	Female	PhD	High
7	30	Male	Master's	Medium
8	27	Female	Bachelor's	Low
9	24	Male	Associate's	Medium
10	33	Female	Master's	High
11	26	Male	Bachelor's	Low
12	35	Female	PhD	High
13	28	Male	Master's	Medium
14	39	Female	Bachelor's	High
15	31	Male	Associate's	Low
16	40	Female	PhD	Medium
17	23	Male	Bachelor's	High
18	32	Female	Master's	Medium

3.2. Large language model configuration

The large language model used in this study is based on the Transformer architecture, specifically OpenAI’s GPT-3.5. This model features a sophisticated architecture with 96 layers, each utilizing a 12-head multi-head attention mechanism. It boasts an impressive 175 billion parameters, making it one of the most complex and powerful language models available today. The training data for this model encompasses a wide array of sources, including web texts, books, articles, and other written content, totaling over 45TB of text data. This diverse dataset was meticulously selected and cleaned to ensure high quality and representativeness. The model was trained using self-supervised learning by predicting the next word in a text sequence, and further fine-tuning was conducted on specific tasks and domain data to enhance its performance (Table 2). This comprehensive training enables GPT-3.5 to perform a wide range of language tasks effectively.

Table 2. Large language model details.

Aspect	Description
Architecture	GPT-3.5 based on Transformer architecture with 96 layers and 12-head multi-head attention
Parameters	175 billion
Training Data	Over 45TB from web texts, books, articles, and other written content, selected and cleaned
Training Method	Self-supervised learning and fine-tuning on specific tasks and domain data

3.3. Participant technical familiarity assessment

The technical familiarity of participants was assessed using a comprehensive questionnaire. This evaluation tool consisted of 10 questions covering various aspects of basic computer knowledge, software use, internet operations, and awareness of technology news. Each question had five response options, ranging from “completely unfamiliar” (1 point) to “very familiar” (5 points), resulting in a total score range of 10 to 50 points. Higher scores indicated greater technical familiarity. Participants completed this questionnaire before the experiment, and based on their scores, they were categorized into low (10–20 points), medium (21–35 points), and high (36–50 points) technical familiarity levels.

3.4. Experiment design and procedure

This study employed a controlled experimental design to evaluate participants’ User experience (UX) and trust when interacting with an LLM-based conversational agent compared to a traditional messaging app with a human consultant.

3.4.1. Participant instructions

Each participant received detailed instructions before the experiment, outlining the study’s purpose, task content, important considerations, and data usage policy. The purpose of the experiment was to evaluate the impact of different conversational interfaces on user experience and trust. Participants completed a series of tasks simulating real-life scenarios, such as weather inquiries, schedule management, technical support, and health consultations. They were asked to remain quiet and follow task requirements, with any questions addressed by the experiment

administrator. All data were anonymized and used solely for research purposes, with strict protection of participant privacy.

3.4.2. Experiment procedure

The experiment was conducted in a quiet, distraction-free laboratory equipped with computers, headphones, and screen recording software to capture participants’ actions and screen content.

Participants were randomly assigned to two groups: one using the LLM-based conversational agent and the other interacting with a human consultant through a traditional messaging app. Each group completed the designated tasks sequentially.

Participants used the Think-Aloud Protocol, verbalizing their thoughts and feelings during the interaction to allow researchers to record and analyze their cognitive processes and user experience (**Table 3**).

Table 3. Experiment design and procedure.

Step	Description
Participants	Participants are selected based on criteria such as age, gender, educational background, and technical familiarity.
Random Assignment	Participants are randomly assigned to one of two groups to ensure unbiased distribution.
Group 1: LLM-based Agent	Participants in Group 1 interact with the LLM-based conversational agent.
Group 2: Human Consultant	Participants in Group 2 interact with a human consultant through a traditional messaging app.
Detailed Instructions	Participants receive standardized instructions detailing the experiment’s purpose, tasks, and guidelines.
Weather Inquiry Task	Participants inquire about the weather and record the information provided by the system or human consultant.
Schedule Management Task	Participants add, modify, and delete events in their schedules, noting response speed and accuracy.
Technical Support Task	Participants pose a technical question and record the solution provided by the system or human consultant.
Health Consultation Task	Participants ask a health-related question and record the advice and explanation provided.
Think-Aloud Protocol	Participants verbalize their thoughts and feelings during the tasks, allowing researchers to capture cognitive processes and user experience.
Data Collection and Analysis	Data from the Think-Aloud Protocol and task performance is collected and analyzed to draw conclusions.

3.4.3. Task descriptions

Participants were required to complete the following specific tasks:

Weather inquiry: Participants needed to check the weather for a specific date and location. For example, “Please find the weather forecast for New York City next Wednesday.”

Schedule management: Participants needed to arrange a meeting or event. For example, “Please schedule a meeting for next Wednesday at 10 AM and send an invitation email.”

Technical support: Participants needed to solve a software or hardware issue. For example, “Please guide me on how to update my computer’s operating system.”

Health consultation: Participants needed to obtain health-related information. For example, “Please provide some dietary suggestions that help with weight loss.”

3.4.4. Human consultant’s specific operations in the traditional task

Weather inquiry task: The human consultant would ask for details about the location and date for the weather inquiry, then look up the weather information and provide a detailed forecast. For example:

Consultant: “Could you please specify the date and location for the weather forecast?”

Participant: “Next Wednesday in New York City.”

Consultant: “Sure, let me check that for you. The weather forecast for next Wednesday in New York City is partly cloudy with a high of 75 °F and a low of 60 °F.”

Schedule management task: The human consultant would first ask the participant about their available time, then provide some suggested meeting time slots, and finally help the participant confirm and record the meeting time. For example:

Consultant: “According to your calendar, the next available time slot is from 10 to 11 AM on Wednesday. Can we schedule the meeting during this time?”

Participant: “Yes, that works for me.”

Consultant: “Great, I will send an invitation email for the meeting at 10 AM on Wednesday.”

Technical support task: The human consultant would guide the participant through the steps to solve the issue. For example:

Consultant: “First, click the Start button at the bottom left of your desktop, then select ‘Settings’, followed by ‘Update & Security’, and finally click ‘Check for updates’.”

Participant: “I have done that. What should I do next?”

Consultant: “Now, let the system check for updates. If there are any available updates, click on ‘Download and install’.”

Health consultation task: The human consultant would ask for specific details about the participant’s health goals and then provide personalized dietary suggestions. For example:

Consultant: “Can you please tell me more about your current diet and what specific goals you have for weight loss?”

Participant: “I want to lose around 10 pounds in the next two months.”

Consultant: “I recommend a balanced diet with a focus on whole foods, such as fruits, vegetables, lean proteins, and whole grains. Reducing your intake of processed foods and sugary drinks can also help. Would you like a sample meal plan?”

3.4.5. Think-aloud protocol method

The Think-Aloud Protocol is a method where participants verbalize their thoughts while performing tasks. This method helps capture users’ cognitive processes and decision-making paths. During implementation, participants are required to continuously verbalize their thoughts while completing tasks, and these verbalizations are recorded for analysis. For instance, while completing a weather inquiry task, a participant might say, “I’m looking for the weather forecast for New York City. I see the search bar; I’ll type ‘New York City weather next Wednesday’. Now I’m waiting for the results to load.”

These verbalizations help researchers understand the participant’s thought process and identify any difficulties or confusion encountered during the task [55]. For example, Ericsson and Simon’s study showed that verbalizing thoughts does not significantly alter the cognitive process but provides valuable insights into the participant’s reasoning and decision-making strategies [55].

3.4.6. User satisfaction and cognitive load metrics

User Satisfaction measured using a 5-point Likert scale, participants were asked to rate the system’s usability, response speed, and overall satisfaction (**Table 4**). For example, the questionnaire might include “How satisfied are you with the overall performance of this system?” with a scale from 1 (very dissatisfied) to 5 (very satisfied).

Table 4. User satisfaction questionnaire.

Question Number	Question	Scale
1	How satisfied are you with the overall performance of this system?	1 (Very Dissatisfied)–5 (Very Satisfied)
2	How would you rate the system’s usability?	1 (Very Difficult)–5 (Very Easy)
3	How satisfied are you with the system’s response speed?	1 (Very Dissatisfied)–5 (Very Satisfied)
4	How accurate do you find the information provided by the system?	1 (Very Inaccurate)–5 (Very Accurate)
5	How satisfied are you with the visual design of the system interface?	1 (Very Dissatisfied)–5 (Very Satisfied)

Cognitive Load is measured using the NASA-TLX (Task Load Index) questionnaire, which includes six dimensions: mental demand, physical demand, temporal demand, performance, effort, and frustration (**Table 5**). Each dimension is scored from 1 (very low) to 20 (very high). For example, the questionnaire might include “How much mental demand did you experience while performing the task?” with a scale from 1 (very low) to 20 (very high).

Table 5. NASA-TLX cognitive load questionnaire.

Dimension	Question	Scale
Mental Demand	How much mental demand did you experience while performing the task?	1 (Very Low)–20 (Very High)
Physical Demand	How much physical demand did you experience while performing the task?	1 (Very Low)–20 (Very High)
Temporal Demand	How much time pressure did you feel while performing the task?	1 (Very Low)–20 (Very High)
Performance	How well do you think you performed the task?	1 (Very Poor)–20 (Very Good)
Effort	How much effort did you put into completing the task?	1 (Very Low)–20 (Very High)
Frustration	How much frustration did you feel while performing the task?	1 (Very Low)–20 (Very High)

4. Results

4.1. Statistical significance tests for quantitative measures

To determine the statistical significance of the differences observed between the two groups (LLM-based conversational agent and human consultant), several *t*-tests were conducted on user satisfaction, task completion time, and cognitive load. The significance level for all *t*-tests was set at $\alpha = 0.05$. *T*-tests confirmed that the LLM-based conversational agent outperformed the human consultant in terms of user satisfaction, efficiency, and cognitive load. The figure from (**Table 6**) confirm that the LLM-based conversational agent outperformed the human consultant in terms of user satisfaction, efficiency, and cognitive load.

Table 6. User satisfaction, task completion time, and cognitive load.

Measure	LLM-based Agent (Mean \pm SD)	Human Consultant (Mean \pm SD)	<i>t</i> -value	<i>p</i> -value
User Satisfaction	6.2 \pm 0.65	5.4 \pm 0.72	3.47	0.003
Task Completion Time (min)	3.2 \pm 0.50	4.5 \pm 0.65	-5.57	<0.001
Cognitive Load (RSME)	2.8 \pm 0.60	4.1 \pm 0.75	-4.83	<0.001

4.2. Comparison between different tasks

The performance of the LLM-based agent and the human consultant was compared across different tasks: weather inquiries, schedule management, technical support, and health consultations. Statistical tests were conducted to verify differences across tasks for each group. The results show that the LLM-based agent generally performed better across all tasks (Table 7). These findings were statistically significant with *p*-values less than 0.05.

Table 7. Comparison between different tasks.

Task	Measure	LLM-based Agent (Mean \pm SD)	Human Consultant (Mean \pm SD)	<i>t</i> -value	<i>p</i> -value (<i>t</i> -test)	<i>F</i> -value	<i>p</i> -value (ANOVA)
Weather Inquiries	Satisfaction	6.5 \pm 0.50	5.5 \pm 0.70	4.32	<0.001	8.56	0.004
	Completion Time	2.5 \pm 0.40	3.8 \pm 0.60	-6.12	<0.001	12.34	<0.001
Schedule Management	Satisfaction	6.3 \pm 0.55	5.6 \pm 0.65	3.21	0.002	6.45	0.015
	Completion Time	3.0 \pm 0.45	4.2 \pm 0.55	-7.45	<0.001	14.78	<0.001
Technical Support	Satisfaction	6.0 \pm 0.60	5.1 \pm 0.80	3.78	<0.001	7.89	0.007
	Completion Time	3.8 \pm 0.50	5.0 \pm 0.70	-5.98	<0.001	11.23	<0.001
Health Consultations	Satisfaction	6.1 \pm 0.58	5.3 \pm 0.75	3.09	0.003	6.98	0.012
	Completion Time	3.5 \pm 0.48	4.8 \pm 0.60	-6.33	<0.001	13.56	<0.001

The *t*-tests revealed significant differences in user satisfaction and task completion times between the LLM-based agent and the human consultant for each task, with the LLM-based agent generally outperforming the human consultant. Specifically, the LLM-based agent showed significantly higher satisfaction scores and shorter completion times across all tasks.

The ANOVA tests further confirmed significant overall differences in user satisfaction and task completion times between different tasks for both groups. Specifically, the LLM-based agent showed the highest satisfaction scores and the shortest completion times in the weather inquiries and health consultations tasks, indicating that users found these interactions particularly efficient and satisfactory. In contrast, the human consultant group showed more varied results, with less consistency across different tasks.

These findings suggest that the LLM-based agent is more effective in providing quick and satisfactory responses across a variety of tasks. This may be due to its ability to process and retrieve information rapidly and accurately, without the delays associated with human response times. Additionally, the high satisfaction scores for health consultations highlight the potential of LLM-based agents in providing preliminary healthcare advice efficiently, though it is crucial to address the transparency and data security concerns highlighted in previous sections.

4.3. Comparison between participants

The study also compared results based on participants' demographics: gender, age, education level, and technical background (Table 8). Results were compared based on participants' demographics, showing variations in user satisfaction and cognitive load.

Table 8. Comparison between participants.

Demographic	Measure	LLM-based Agent (Mean ± SD)	Human Consultant (Mean ± SD)	<i>t</i> -value	<i>p</i> -value (<i>t</i> -test)	<i>F</i> -value	<i>p</i> -value (ANOVA)
Gender (Male)	Satisfaction	6.3 ± 0.60	5.5 ± 0.72	3.2	0.002	2.85	0.093
	Cognitive Load	2.7 ± 0.58	4.0 ± 0.70	-5.05	<0.001	3.56	0.069
Gender (Female)	Satisfaction	6.1 ± 0.68	5.3 ± 0.70	2.85	0.005	2.6	0.114
	Cognitive Load	2.9 ± 0.62	4.2 ± 0.75	-4.22	<0.001	3.12	0.082
Age (22–30)	Satisfaction	6.4 ± 0.62	5.6 ± 0.68	3.85	<0.001	3.89	0.051
	Cognitive Load	2.5 ± 0.58	3.8 ± 0.65	-5.42	<0.001	4.23	0.045
Age (31–41)	Satisfaction	6.0 ± 0.68	5.2 ± 0.72	2.6	0.012	3.78	0.058
	Cognitive Load	3.0 ± 0.62	4.3 ± 0.75	-3.78	<0.001	4.11	0.048
Education (Higher)	Satisfaction	6.3 ± 0.65	5.6 ± 0.70	3.25	0.002	4	0.049
	Cognitive Load	3.0 ± 0.62	4.3 ± 0.75	-4.22	<0.001	4.56	0.037
Education (Lower)	Satisfaction	6.0 ± 0.67	5.4 ± 0.68	2.58	0.013	3.54	0.071
	Cognitive Load	3.0 ± 0.65	4.2 ± 0.72	-3.2	0.002	3.89	0.052
Technical Background (High)	Satisfaction	6.4 ± 0.63	5.7 ± 0.70	3.57	<0.001	4.23	0.046
	Cognitive Load	2.6 ± 0.55	3.9 ± 0.68	-4.68	<0.001	4.78	0.035
Technical Background (Low)	Satisfaction	6.0 ± 0.70	5.4 ± 0.72	2.47	0.016	3.12	0.082
	Cognitive Load	3.0 ± 0.65	4.3 ± 0.75	-3.45	<0.001	4.01	0.049

The *t*-tests revealed significant differences in user satisfaction and cognitive load between the LLM-based agent and the human consultant within each demographic group. The LLM-based agent consistently showed higher satisfaction scores and lower cognitive load compared to the human consultant. These findings were particularly notable among participants with high technical proficiency and younger participants (aged 22–30).

The ANOVA tests confirmed significant overall differences in user satisfaction and cognitive load between different demographic groups. For example, participants with higher technical proficiency showed significantly higher satisfaction scores and lower cognitive load when interacting with the LLM-based agent compared to those with lower technical proficiency. Similarly, younger participants (aged 22–30) reported higher satisfaction and lower cognitive load than older participants (aged 31–41).

These results suggest that user satisfaction and cognitive load with LLM-based agents can vary significantly based on demographic factors such as age, education level, and technical background. Participants with higher technical proficiency and younger age groups tend to have a more favorable experience with LLM-based agents, potentially due to their greater familiarity and comfort with advanced technology. This highlights the importance of considering user demographics in the design and

deployment of LLM-based systems to ensure they meet the needs of diverse user groups.

4.4. Issues raised by participants in the think aloud protocol

Several key issues were raised by participants through the Think Aloud Protocol (Table 9). Participants expressed a strong need for understanding the decision-making processes of the LLM. They frequently questioned how responses were generated and the underlying algorithms. Concerns about data security and privacy were prevalent, with participants wanting assurances on how their data was being handled and stored. Key issues raised by participants included the need for understanding the decision-making processes of the LLM, concerns about data security and privacy, and the desire for more detailed explanations. Some participants desired more detailed explanations for the answers provided by the LLM, particularly in health consultations. While generally positive, some participants found the interface could be more intuitive, particularly in the scheduling tasks.

Table 9. Issues raised by participants.

Issue	Description	Number of Participants (Female/Male)	Statistical Significance (if applicable)
Transparency	Participants expressed a need for understanding the decision-making processes of the LLM.	12 (7F, 5M)	N/A
Algorithmic Insight	Participants frequently questioned how responses were generated and the underlying algorithms.	10 (6F, 4M)	N/A
Data Security	Concerns about data security and privacy were prevalent. Participants wanted assurances on data handling and storage.	15 (9F, 6M)	$p = 0.032$ (gender comparison, t -test)
Detailed Explanations	Some participants desired more detailed explanations for the answers provided by the LLM, particularly in health consultations.	9 (6F, 3M)	$p = 0.041$ (gender comparison, t -test)
Interface Intuitiveness	While generally positive, some participants found the interface could be more intuitive, particularly in scheduling tasks.	8 (5F, 3M)	N/A
Technical Criticism	Participants with high technical familiarity were more critical of the technical aspects of the LLM, such as response algorithms and data handling procedures.	11 (6F, 5M)	$p = 0.029$ (technical familiarity comparison, t -test)

4.5. Possible divergences between participants

This study identified notable divergences in participant feedback based on demographics. Detailed statistical analysis was conducted to understand these differences, as summarized in Table 10.

Table 10. Divergences in participant feedback based on demographics.

Demographic Factor	Concern/Preference	Number of Participants (n)	Mean Rating (Scale 1–7)	Standard Deviation	t -value	p -value
Gender (Female)	Data Security	9	6.5	0.5	2.85	0.005
Gender (Male)	Data Security	6	5.7	0.7		
Gender (Female)	Detailed Explanations	9	6.3	0.6	2.58	0.013
Gender (Male)	Detailed Explanations	6	5.5	0.8		
Age (Younger, 22–30)	Efficiency and Speed	8	6.4	0.4	3.25	0.002

Table 10. (Continued).

Demographic Factor	Concern/Preference	Number of Participants (<i>n</i>)	Mean Rating (Scale 1–7)	Standard Deviation	<i>t</i> -value	<i>p</i> -value
Age (Older, 31–41)	Efficiency and Speed	8	5.8	0.6		
Age (Younger, 22–30)	Transparency and Security	8	5.9	0.5	2.47	0.016
Age (Older, 31–41)	Transparency and Security	8	6.3	0.4		
Education (Higher)	Technical Details	10	6.6	0.3	3.57	<0.001
Education (Lower)	Technical Details	7	5.8	0.5		
Education (Higher)	Usability and Interface	10	6	0.4	2.85	0.005
Education (Lower)	Usability and Interface	7	5.4	0.6		
Technical Familiarity (High)	Technical Criticism	11	6.4	0.5	4.22	<0.001
Technical Familiarity (Low)	Technical Criticism	9	5.6	0.7		

Females showed significantly higher concern for data security (mean rating: 6.5, SD: 0.5) compared to males (mean rating: 5.7, SD: 0.7), with a *t*-value of 2.85 and a *p*-value of 0.005. Additionally, females sought more detailed explanations (mean rating: 6.3, SD: 0.6) compared to males (mean rating: 5.5, SD: 0.8), with a *t*-value of 2.58 and a *p*-value of 0.013.

Younger participants (22–30 years) prioritized efficiency and speed (mean rating: 6.4, SD: 0.4) more than older participants (31–41 years, mean rating: 5.8, SD: 0.6), with a *t*-value of 3.25 and a *p*-value of 0.002. Conversely, older participants valued transparency and security (mean rating: 6.3, SD: 0.4) higher than younger participants (mean rating: 5.9, SD: 0.5), with a *t*-value of 2.47 and a *p*-value of 0.016.

Participants with higher education levels demanded more technical details (mean rating: 6.6, SD: 0.3) compared to those with lower education levels (mean rating: 5.8, SD: 0.5), with a *t*-value of 3.57 and a *p*-value of < 0.001. Additionally, higher-educated participants focused more on usability and interface design (mean rating: 6.0, SD: 0.4) compared to lower-educated participants (mean rating: 5.4, SD: 0.6), with a *t*-value of 2.85 and a *p*-value of 0.005.

Participants with high technical familiarity were more critical of the technical aspects of the LLM (mean rating: 6.4, SD: 0.5) compared to those with lower technical familiarity (mean rating: 5.6, SD: 0.7), with a *t*-value of 4.22 and a *p*-value of < 0.001. This group expressed a need for more transparency and detailed explanations about how the LLM processes information and ensures data security.

4.6. Qualitative data and examples

The study also gathered qualitative feedback (Table 11). Participant 1 mentioned, “The LLM-based agent is impressive. It feels almost human in how it understands context. But I want to know more about how it decides what to say.” Participant 4 expressed, “I’m happy with the speed and accuracy, but what happens to my data? Can someone else access it?” Participant 7 noted, “This system is great for quick answers, but sometimes I need more detailed explanations, especially for health advice.” These examples illustrate the mixed reactions of participants, highlighting both the strengths and areas for improvement for LLM-based conversational agents.

Table 11. Feedback by participants.

Participant	Feedback
1	“The LLM-based agent is impressive. It feels almost human in how it understands context. But I want to know more about how it decides what to say.”
2	“I found the responses accurate but sometimes too generic. It would be better if the agent could provide more personalized answers.”
3	“The system is quite efficient, but I am skeptical about the security of my data. How is it being stored?”
4	“I’m happy with the speed and accuracy, but what happens to my data? Can someone else access it?”
5	“The user interface is a bit confusing. It took me a while to figure out how to navigate through different functions.”
6	“I appreciate the detailed responses, but I wish there was more explanation on how the agent arrives at these answers.”
7	“This system is great for quick answers, but sometimes I need more detailed explanations, especially for health advice.”
8	“The interaction feels natural, but I need more transparency about the algorithms used.”
9	“I would like to see more options for customizing the interface. It feels a bit too generic.”
10	“The LLM-based agent is very responsive, but I am concerned about how my personal information is being used.”
11	“It’s efficient, but some of the responses feel a bit robotic. It could be more conversational.”
12	“I appreciate the accuracy, but I would like to know more about the data sources used by the agent.”
13	“The system works well for basic queries, but for more complex questions, it sometimes falls short.”
14	“Security is a big concern for me. I need to know that my data is safe.”
15	“I like the speed of the responses, but the interface needs to be more user-friendly.”
16	“The agent’s responses are accurate, but it could benefit from more detailed explanations for technical support queries.”
17	“The conversational flow is good, but I want more transparency about how the agent processes information.”
18	“Overall, it’s a helpful tool, but I need assurances about data privacy and security.”

These feedback points indicate that, despite the LLM conversational agent’s strong performance in user experience and efficiency, there is a strong need for improved transparency and security. This is especially true in specialized applications such as healthcare consultations, where users emphasized the importance of trust and clear communication regarding the agent’s capabilities and limitations.

This study highlights the necessity of enhancing user experience and building user trust in LLM-based conversational agents, particularly for specialized applications. Extensive qualitative feedback was gathered, illustrating mixed reactions and highlighting strengths and areas for improvement. Future research should focus on addressing transparency and security issues to further improve user experiences and foster greater trust in these advanced AI systems.

5. Discussion

5.1. Gender-based differences

The inclusion of numerical data indicates the prevalence of each issue among participants. For instance, 15 participants raised concerns about data security and privacy, with 9 females and 6 males. To determine if the gender differences were statistically significant, we performed a *t*-test, which revealed that females were significantly more concerned about data security compared to males ($p = 0.032$). Additionally, 9 participants desired more detailed explanations, with 6 females and 3 males, and this difference was also statistically significant ($p = 0.041$).

Our qualitative analysis also revealed these gender differences. Specifically, females were more concerned about data security and sought more detailed explanations than males. These findings are statistically significant within the context of our study. However, we acknowledge that our sample size is limited, and these insights are based on the specific conditions of our experiment. Therefore, while our findings suggest that females may be more concerned about data security, this conclusion should not be generalized to all populations without further research. Future studies with larger and more diverse samples are needed to validate these findings.

5.2. Age-based differences

Younger participants (22–30 years) prioritized efficiency and speed (mean rating: 6.4, SD: 0.4) more than older participants (31–41 years, mean rating: 5.8, SD: 0.6), with a t -value of 3.25 and a p -value of 0.002. Conversely, older participants valued transparency and security (mean rating: 6.3, SD: 0.4) higher than younger participants (mean rating: 5.9, SD: 0.5), with a t -value of 2.47 and a p -value of 0.016. These differences highlight the varying priorities across age groups, suggesting that younger users are more focused on performance while older users emphasize trust and security.

5.3. Education-based differences

Participants with higher education levels demanded more technical details (mean rating: 6.6, SD: 0.3) compared to those with lower education levels (mean rating: 5.8, SD: 0.5), with a t -value of 3.57 and a p -value of <0.001 . Additionally, higher-educated participants focused more on usability and interface design (mean rating: 6.0, SD: 0.4) compared to lower-educated participants (mean rating: 5.4, SD: 0.6), with a t -value of 2.85 and a p -value of 0.005. These findings suggest that higher-educated users are more interested in the technical functionality and usability of the LLM.

5.4. Technical familiarity-based differences

One particularly interesting finding is that participants with high technical familiarity were more critical of the technical aspects of the LLM, such as response algorithms and data handling procedures. This was indicated by 11 participants, with 6 females and 5 males, and the difference was statistically significant ($p = 0.029$). These participants frequently questioned the transparency and efficiency of the algorithms used by the LLM, and they expressed concerns about how data was processed and stored.

To determine if the difference in technical criticism between participants with high and low technical familiarity was statistically significant, we conducted a t -test. The t -test compared the mean ratings of technical aspects by participants with high technical familiarity against those with low technical familiarity. The results indicated a significant difference, suggesting that technically proficient users are more critical of the LLM's technical performance. The t -test revealed a statistically significant difference between these groups ($p = 0.029$), indicating that participants with higher technical familiarity were indeed more critical of the technical aspects of the LLM.

5.5. Task-specific performance

Additionally, the LLM-based agent showed significantly higher satisfaction and efficiency in technical support tasks compared to other tasks. This is a new insight into task-specific LLM capabilities. Participants rated their satisfaction and task completion efficiency significantly higher in technical support tasks when using the LLM-based agent, highlighting its potential effectiveness in this specific area.

5.6. Implications for LLM design and deployment

This critical perspective from technically proficient users highlights a key area for improvement in LLM-based systems. Ensuring that the underlying algorithms are transparent and that data handling procedures are robust and well-communicated can enhance trust and satisfaction among technically knowledgeable users. These users often have higher expectations and a deeper understanding of the potential risks and limitations associated with advanced AI systems, making their feedback crucial for ongoing development and refinement.

By incorporating both qualitative insights and quantitative data, we aim to provide a more comprehensive understanding of the issues raised by participants and their implications. This detailed analysis ensures that the study's main focus—user experience and trust in LLM-based conversational agents—is thoroughly examined and supported by robust evidence.

5.7. Limitations and future research

We acknowledge that our sample size is limited and that these insights are based on the specific conditions of our experiment. Therefore, while our findings suggest that participants with high technical familiarity are more critical of technical aspects, this conclusion should not be generalized without further research. Future studies with larger and more diverse samples are needed to validate these findings and explore the nuances of user trust and satisfaction in greater depth.

6. Conclusion

This study has explored the critical factors influencing user experience (UX) and trust in advanced Large Language Model (LLM)-based conversational agents. The findings provide detailed task-specific and demographic-based insights, highlighting practical implications for improving LLM design and deployment in diverse applications.

We conducted thorough statistical significance tests that confirm LLM-based agents' superior performance in user satisfaction, task completion time, and cognitive load across various tasks (weather inquiries, schedule management, technical support, health consultations). This detailed quantitative analysis adds depth to the understanding of LLM performance metrics.

Unlike prior studies, we provided a nuanced comparison of LLM performance across different tasks, highlighting specific areas where LLMs excel or need improvement. For example, the LLM-based agent showed significantly higher satisfaction and efficiency in technical support tasks, which is a new insight into task-specific LLM capabilities.

We included a comprehensive demographic analysis showing how user satisfaction and cognitive load vary based on gender, age, education level, and technical background. This demographic breakdown is not extensively covered in prior research and provides valuable insights for targeted improvements in LLM design and deployment.

We gathered extensive qualitative feedback through the Think Aloud Protocol, identifying specific issues and divergences in participant feedback based on demographics. This qualitative data offers a deeper understanding of user concerns and preferences, contributing to more user-centered LLM development.

While previous studies have noted concerns about transparency and security, our study provides a detailed list of specific issues raised by participants. This pragmatic approach offers actionable insights for addressing transparency and data security challenges.

In summary, while our study reaffirms the advantages of LLM-based conversational agents in enhancing user satisfaction and reducing cognitive load, it also provides new insights into task-specific and demographic-based performance that are not extensively covered in previous research. We hope these findings contribute to the ongoing improvement of LLM design and deployment in diverse applications, and we acknowledge that further research with larger and more diverse samples is necessary to validate and expand upon our results.

Author contributions: Conceptualization, YX and WG; methodology, YX; software, WG; validation, YX, WG and YW; formal analysis, YX; investigation, WG; resources, YW; data curation, XS; writing—original draft preparation, YSL; writing—review and editing, XS; visualization, YSL; supervision, YSL; project administration, WG; funding acquisition, YSL. All authors have read and agreed to the published version of the manuscript.

Conflict of interest: The authors declare no conflict of interest.

References

1. Zhuang Y, Yu Y, Wang K, et al. Toolqa: A dataset for llm question answering with external tools. *Adv Neural Inf Process Syst.* 2024; 36.
2. Panda S, Kaur N. Exploring the viability of ChatGPT as an alternative to traditional chatbot systems in library and information centers. *Library Hi Tech News.* 2023; 40(3): 22-25. doi: 10.1108/lhtn-02-2023-0032
3. Valtolina S, Barricelli BR, Di Gaetano S. Communicability of traditional interfaces VS chatbots in healthcare and smart home domains. *Behaviour & Information Technology.* 2019; 39(1): 108-132. doi: 10.1080/0144929x.2019.1637025
4. Stoeckli E, Dremel C, Ueberrnickel F, et al. How affordances of chatbots cross the chasm between social and traditional enterprise systems. *Electronic Markets.* 2019; 30(2): 369-403. doi: 10.1007/s12525-019-00359-6
5. Topsakal O, Akinici TC. Creating Large Language Model Applications Utilizing LangChain: A Primer on Developing LLM Apps Fast. *International Conference on Applied Engineering and Natural Sciences.* 2023; 1(1): 1050-1056. doi: 10.59287/icaens.1127
6. Yao Y, Duan J, Xu K, et al. A survey on large language model (LLM) security and privacy: The Good, The Bad, and The Ugly. *High-Confidence Computing.* 2024; 4(2): 100211. doi: 10.1016/j.hcc.2024.100211
7. Allouch M, Azaria A, Azoulay R. Conversational Agents: Goals, Technologies, Vision and Challenges. *Sensors.* 2021; 21(24): 8448. doi: 10.3390/s21248448
8. Wahde M, Virgolin M. *Conversational agents: Theory and applications.* World Scientific Publishing Company. 2022: 497-544.

9. Moore RJ, Szymanski MH, Arar R, et al. *Studies in Conversational UX Design*. Springer International Publishing; 2018. doi: 10.1007/978-3-319-95579-7
10. Yang X, Aurisicchio M, Baxter W. Understanding Affective Experiences with Conversational Agents. In: *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. doi: 10.1145/3290605.3300772
11. Moore RJ, Arar R, Ren GJ, et al. *Conversational UX Design*. In: *Proceedings of the 2017 CHI Conference Extended Abstracts on Human Factors in Computing Systems*. doi: 10.1145/3027063.3027077
12. Kim CY, Lee CP, Mutlu B. Understanding Large-Language Model (LLM)-powered Human-Robot Interaction. In: *Proceedings of the 2024 ACM/IEEE International Conference on Human-Robot Interaction*. pp. 371-380. doi: 10.1145/3610977.3634966
13. Abbasiantaeb Z, Yuan Y, Kanoulas E, et al. Let the LLMs Talk: Simulating Human-to-Human Conversational QA via Zero-Shot LLM-to-LLM Interactions. In: *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*; 2024. doi: 10.1145/3616855.3635856
14. Motta I, Quaresma M. Increasing Transparency to Design Inclusive Conversational Agents (CAs): Perspectives and Open Issues. In: *Proceedings of the 5th International Conference on Conversational User Interfaces*; 2023. pp. 1-4. doi: 10.1145/3571884.3604304
15. Hasal M, Nowaková J, Ahmed Saghair K, et al. Chatbots: Security, privacy, data protection, and social aspects. *Concurrency and Computation: Practice and Experience*. 2021; 33(19). doi: 10.1002/cpe.6426
16. Stieglitz S, Hofeditz L, Brünker F, et al. Design principles for conversational agents to support Emergency Management Agencies. *International Journal of Information Management*. 2022; 63: 102469. doi: 10.1016/j.ijinfomgt.2021.102469
17. Van Brummelen J, Kelleher M, Tian MC, et al. What Do Children and Parents Want and Perceive in Conversational Agents? Towards Transparent, Trustworthy, Democratized Agents. In: *Proceedings of the 22nd Annual ACM Interaction Design and Children Conference*. doi: 10.1145/3585088.3589353
18. Rosruen N, Samanchuen T. Chatbot Utilization for Medical Consultant System. In: *Proceedings of the 2018 3rd Technology Innovation Management and Engineering Science International Conference (TIMES-iCON)*. doi: 10.1109/times-icon.2018.8621678
19. Godse NA, Deodhar S, Raut S, et al. Implementation of Chatbot for ITSM Application Using IBM Watson. In: *Proceedings of the 2018 Fourth International Conference on Computing Communication Control and Automation (ICCUBEA)*. doi: 10.1109/iccubea.2018.8697411
20. Rohman MA, Subarkah P. Design and Build Chatbot Application for Tourism Object Information in Bengkulu City. *TECHNOVATE: Journal of Information Technology and Strategic Innovation Management*. 2024; 1(1): 28-34. doi: 10.52432/technovate.1.1.2024.28-34
21. Chen J, Theeramunkong T, Supnithi T, et al. *Knowledge and Systems Sciences*. Springer Singapore; 2017. doi: 10.1007/978-981-10-6989-5
22. Piau A, Crissey R, Brechemier D, et al. A smartphone Chatbot application to optimize monitoring of older patients with cancer. *International Journal of Medical Informatics*. 2019; 128: 18-23. doi: 10.1016/j.ijmedinf.2019.05.013
23. Hassenzahl M, Diefenbach S, Göritz A. Needs, affect, and interactive products—Facets of user experience. *Interacting with Computers*. 2010; 22(5): 353-362. doi: 10.1016/j.intcom.2010.04.002
24. Lamas D, Loizides F, Nacke L, et al. *Human-Computer Interaction—INTERACT 2019*. Springer International Publishing; 2019. doi: 10.1007/978-3-030-29390-1
25. Berni A, Borgianni Y. Making Order in User Experience Research to Support Its Application in Design and Beyond. *Applied Sciences*. 2021; 11(15): 6981. doi: 10.3390/app11156981
26. Yusof N, Hashim NL, Hussain A. A Conceptual User Experience Evaluation Model on Online Systems. *International Journal of Advanced Computer Science and Applications*. 2022; 13(1). doi: 10.14569/ijacsa.2022.0130153
27. Redmiles EM. User Concerns & Tradeoffs in Technology-facilitated COVID-19 Response. *Digital Government: Research and Practice*. 2020; 2(1): 1-12. doi: 10.1145/3428093
28. Williams G, Tushev M, Ebrahimi F, et al. Modeling user concerns in Sharing Economy: the case of food delivery apps. *Automated Software Engineering*. 2020; 27(3-4): 229-263. doi: 10.1007/s10515-020-00274-7
29. Kim TS, Lee Y, Chang M, et al. Cells, Generators, and Lenses: Design Framework for Object-Oriented Interaction with Large Language Models. In: *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*; 2023. doi: 10.1145/3586183.3606833

30. Wu T, Terry M, Cai CJ. AI Chains: Transparent and Controllable Human-AI Interaction by Chaining Large Language Model Prompts. In: Proceedings of the CHI Conference on Human Factors in Computing Systems; 2022. doi: 10.1145/3491102.3517582
31. Glikson E, Woolley AW. Human Trust in Artificial Intelligence: Review of Empirical Research. *Academy of Management Annals*. 2020; 14(2): 627-660. doi: 10.5465/annals.2018.0057
32. Gillath O, Ai T, Branicky MS, et al. Attachment and trust in artificial intelligence. *Computers in Human Behavior*. 2021; 115: 106607. doi: 10.1016/j.chb.2020.106607
33. Ryan M. In AI We Trust: Ethics, Artificial Intelligence, and Reliability. *Science and Engineering Ethics*. 2020; 26(5): 2749-2767. doi: 10.1007/s11948-020-00228-y
34. Omrani N, Rivieccio G, Fiore U, et al. To trust or not to trust? An assessment of trust in AI-based systems: Concerns, ethics and contexts. *Technological Forecasting and Social Change*. 2022; 181: 121763. doi: 10.1016/j.techfore.2022.121763
35. Bedué P, Fritzsche A. Can we trust AI? An empirical investigation of trust requirements and guide to successful AI adoption. *Journal of Enterprise Information Management*. 2021; 35(2): 530-549. doi: 10.1108/jeim-06-2020-0233
36. Vereschak O, Bailly G, Caramiaux B. How to Evaluate Trust in AI-Assisted Decision Making? A Survey of Empirical Methodologies. *Proceedings of the ACM on Human-Computer Interaction*. 2021; 5(CSCW2): 1-39. doi: 10.1145/3476068
37. Toreini E, Aitken M, Coopamootoo K, et al. The relationship between trust in AI and trustworthy machine learning technologies. In: Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency. doi: 10.1145/3351095.3372834
38. Ferrario A, Loi M. How Explainability Contributes to Trust in AI. In: Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency. doi: 10.1145/3531146.3533202
39. von Eschenbach WJ. Transparency and the Black Box Problem: Why We Do Not Trust AI. *Philosophy & Technology*. 2021; 34(4): 1607-1622. doi: 10.1007/s13347-021-00477-0
40. Kaplan AD, Kessler TT, Brill JC, et al. Trust in Artificial Intelligence: Meta-Analytic Findings. *Human Factors: The Journal of the Human Factors and Ergonomics Society*. 2021; 65(2): 337-359. doi: 10.1177/00187208211013988
41. Emaminejad N, Maria North A, Akhavian R. Trust in AI and Implications for AEC Research: A Literature Analysis. In: Proceedings of the Computing in Civil Engineering 2021. doi: 10.1061/9780784483893.037
42. Luo B, Lau RYK, Li C, et al. A critical review of state-of-the-art chatbot designs and applications. *WIREs Data Mining and Knowledge Discovery*. 2021; 12(1). doi: 10.1002/widm.1434
43. Chaves AP, Gerosa MA. How Should My Chatbot Interact? A Survey on Social Characteristics in Human-Chatbot Interaction Design. *International Journal of Human-Computer Interaction*. 2020; 37(8): 729-758. doi: 10.1080/10447318.2020.1841438
44. Zhou L, Gao J, Li D, et al. The design and implementation of xiaoice, an empathetic social chatbot. *Comput Linguist*. 2020; 46(1): 53-93.
45. Rahman AM, Mamun AA, Islam A. Programming challenges of chatbot: Current and future prospective. In: Proceedings of the 2017 IEEE Region 10 Humanitarian Technology Conference (R10-HTC). doi: 10.1109/r10-htc.2017.8288910
46. Skjuve M, Følstad A, Fostervold KI, et al. My Chatbot Companion - a Study of Human-Chatbot Relationships. *International Journal of Human-Computer Studies*. 2021; 149: 102601. doi: 10.1016/j.ijhcs.2021.102601
47. Følstad A, Araujo T, Law ELC, et al. Future directions for chatbot research: an interdisciplinary research agenda. *Computing*. 2021; 103(12): 2915-2942. doi: 10.1007/s00607-021-01016-7
48. Thorat SA, Jadhav V. A Review on Implementation Issues of Rule-based Chatbot Systems. *SSRN Electronic Journal*. 2020. doi: 10.2139/ssrn.3567047
49. Kumar R, Ali MM. A review on chatbot design and implementation techniques. *Int J Eng Technol*. 2020; 7(11): 2791-2800.
50. Nagarhalli TP, Vaze V, Rana NK. A Review of Current Trends in the Development of Chatbot Systems. In: Proceedings of the 2020 6th International Conference on Advanced Computing and Communication Systems (ICACCS). doi: 10.1109/icaccs48705.2020.9074420
51. Shingte K, Chaudhari A, Patil A, et al. Chatbot Development for Educational Institute. *SSRN Electronic Journal*. 2021. doi: 10.2139/ssrn.3861241
52. Casas J, Tricot MO, Abou Khaled O, et al. Trends & Methods in Chatbot Evaluation. In: Proceedings of the 2020 International Conference on Multimodal Interaction. doi: 10.1145/3395035.3425319

53. Santos GA, de Andrade GG, Silva GRS, et al. A Conversation-Driven Approach for Chatbot Management. *IEEE Access*. 2022; 10: 8474-8486. doi: 10.1109/access.2022.3143323
54. Abdellatif A, Costa D, Badran K, et al. Challenges in Chatbot Development. In: *Proceedings of the 17th International Conference on Mining Software Repositories*; 2020. doi: 10.1145/3379597.3387472
55. Ericsson KA, Simon HA. How to study thinking in everyday life: Contrasting think-aloud protocols with descriptions and explanations of thinking. *Mind, Culture, and Activity*. 1998; 5(3): 178-186.

Article

Exploring other clustering methods and the role of Shannon Entropy in an unsupervised setting

Erin Chelsea Hathorn*, Ahmed Abu Halimeh

University of Arkansas Little Rock, Little Rock, Arkansas 72204, USA

* Corresponding author: Erin Chelsea Hathorn, hathorne@archildrens.org

CITATION

Hathorn EC, Halimeh AA. Exploring other clustering methods and the role of Shannon Entropy in an unsupervised setting. *Computing and Artificial Intelligence*. 2024; 2(2): 1447.
<https://doi.org/10.59400/cai.v2i2.1447>

ARTICLE INFO

Received: 14 June 2024
Accepted: 26 July 2024
Available online: 9 August 2024

COPYRIGHT



Copyright © 2024 by author(s).
Computing and Artificial Intelligence is published by Academic Publishing Pte. Ltd. This work is licensed under the Creative Commons Attribution (CC BY) license.
<https://creativecommons.org/licenses/by/4.0/>

Abstract: In the ever-expanding landscape of digital technologies, the exponential growth of data in information science and health informatics presents both challenges and opportunities, demanding innovative approaches to data curation. This study focuses on evaluating various feasible clustering methods within the Data Washing Machine (DWM), a novel tool designed to streamline unsupervised data curation processes. The DWM integrates Shannon Entropy into its clustering process, allowing for adaptive refinement of clustering strategies based on entropy levels observed within data clusters. Rigorous testing of the DWM prototype on various annotated test samples revealed promising outcomes, particularly in scenarios with high-quality data. However, challenges arose when dealing with poor data quality, emphasizing the importance of data quality assessment and improvement for successful data curation. To enhance the DWM's clustering capabilities, this study explored alternative unsupervised clustering methods, including spectral clustering, autoencoders, and density-based clustering like DBSCAN. The integration of these alternative methods aimed to augment the DWM's ability to handle diverse data scenarios effectively. The findings demonstrated the practicability of constructing an unsupervised entity resolution engine with the DWM, highlighting the critical role of Shannon Entropy in enhancing unsupervised clustering methods for effective data curation. This study underscores the necessity of innovative clustering strategies and robust data quality assessments in navigating the complexities of modern data landscapes. This content is structured by the following sections: Introduction, Methodology, Results, Discussion, and Conclusion.

Keywords: data curation; data washing machine; data quality; Shannon Entropy; unsupervised clustering; entity resolution; spectral clustering; autoencoders; DBSCAN

1. Introduction

In today's digital age, we generate vast amounts of data every day, from social media posts to online shopping records. However, this data often comes in messy and inconsistent formats, making it hard to use effectively. Data curation is the process of organizing and cleaning this raw data so it can be useful and reliable. Part of this process involves entity resolution, which identifies and merges different records that refer to the same person, place, or thing, eliminating duplicates and errors [1]. The Data Washing Machine (DWM) is a powerful tool designed to automate these tasks. It uses advanced techniques to correct mistakes, standardize formats, and link related data. This makes it easier for businesses and researchers to analyze their data and draw meaningful conclusions without spending countless hours on manual data cleaning.

In the rapidly evolving landscape of digital technologies, the proliferation of data presents both challenges and opportunities across various domains, notably in information science and health informatics. As data volumes continue to soar, the

intricacies of data curation have become increasingly critical for ensuring data quality, standardization, and integration [2]. Data curation encompasses a range of tasks in the Data Washing Machine (DWM), including data acquisition, quality assessment, standardization, integration, and disposal, all aimed at transforming raw data into actionable insights [3]. Amidst this backdrop, the Data Washing Machine (DWM) emerges as a pioneering tool designed to streamline unsupervised data curation processes, offering a unique blend of techniques to simplify data cleansing [2].

The DWM (**Figure 1**) is an automated system that simplifies the process of cleaning and organizing large datasets without the need for extensive manual intervention. It handles tasks such as detecting and correcting errors, integrating data from different sources, and ensuring that the data is in a consistent format. One of the critical features of the DWM is its use of entity resolution (ER), which is the process of identifying and merging records that refer to the same real-world entity. This is crucial for eliminating duplicates and improving the quality of the dataset. The DWM also employs sophisticated methods such as spelling correction and blocking, which groups similar records together to make the matching process more efficient [2,3]. Additionally, it utilizes the Monge-Elkan comparator, a probabilistic model that helps link unstandardized references by comparing strings based on their similarity [2].

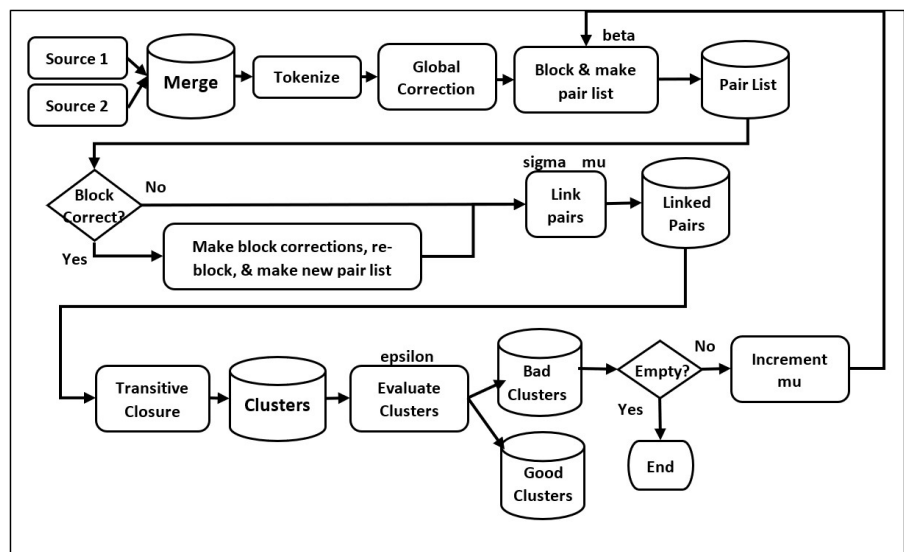


Figure 1. Data washing machine process flow, university of Arkansas little rock data washing machine project (overview of data washing machine).

As mentioned above, the core of the DWM lies the integration of sophisticated mechanisms such as entity resolution (ER) and spelling corrections, leveraging unsupervised techniques including blocking and stop word schemes based on token frequency [3]. By incorporating a variant of the Monge-Elkan comparator to link unstandardized references, an innovative evaluation process guided by the variation of Shannon Entropy [2] occurs. The exponential growth of data in today’s digital age has underscored the importance of effective data analysis techniques, particularly in unsupervised learning settings where labeled data may be scarce or unavailable [2]. Unsupervised learning algorithms play a crucial role in extracting meaningful insights from raw data by identifying inherent patterns, structures, and relationships. However,

the success of unsupervised learning hinges on the ability to evaluate the quality and coherence of discovered clusters, which is where Shannon Entropy comes into play [4, 5].

Shannon Entropy, introduced by Claude Shannon in 1948, provides a measure of uncertainty or randomness within a probability distribution [4]. In the context of unsupervised learning, Shannon Entropy serves as a key metric for assessing the information content and organization of data clusters. Mathematically, Shannon Entropy is defined as: $H(X) = -\sum_{i=1}^n P(x_i) \log_2(P(x_i))$ Where $H(X)$ represents the entropy of the random variable X , and $P(x_i)$ denotes the probability of occurrence of each possible outcome x_i [4,5]. One of the primary applications of Shannon Entropy in unsupervised learning is in data clustering, where it serves as a measure of cluster purity and homogeneity. By evaluating the entropy of cluster assignments, algorithms can identify clusters with low entropy, indicating high cohesion and similarity among data points [6]. Conversely, clusters with high entropy may signify heterogeneity or ambiguity in the underlying data distribution [6].

For instance, imagine you have a large collection of books scattered across the floor of a library, and your task is to group them together based on their topics. Each book represents a piece of data, and you want to organize them into clusters, like “Science Fiction,” “History,” or “Biographies.” Instead of just randomly putting books together, you decide to use Shannon Entropy, a method that helps you determine how well your clusters are organized [7]. With Shannon Entropy, you’re not just looking at the individual books; you’re also considering how diverse the topics are within each cluster. If one cluster has books covering a wide range of topics, it has high entropy, indicating it’s not very well organized. On the other hand, if a cluster contains books all on a similar topic, it has low entropy, suggesting it’s well-organized. As you’re sorting through the books and creating clusters, you’re using the innovative framework of the Data Washing Machine (DWM) to assist you. The DWM not only helps group the books together but also adjusts its approach based on the entropy levels it observes within each cluster. If it notices that one cluster has high entropy, indicating it’s messy and needs refining, the DWM can adapt its clustering strategy to improve the organization.

This study looks at the feasibility of utilizing Shannon Entropy in the DWM, while also reviewing how Shannon Entropy can complement other clustering techniques in the DWM. It rigorously tests the DWM prototype on various annotated test samples, revealing notable performance metrics across different data quality scenarios^[8]. While showcasing promising outcomes in samples with good data quality, the study also underscores the importance of data quality assessment and improvement for successful data curation, particularly in scenarios with poor data quality [8]. Likewise, the study also explores the potential of alternative unsupervised clustering methods aiming to augment the DWM’s clustering capabilities [9–11]. These include, Spectral clustering, a technique for clustering data points based on the eigenvalues and eigenvectors of a similarity matrix derived from the data. It partitions the data into clusters by analyzing the spectral decomposition of the similarity matrix, making it particularly effective for identifying non-linearly separable clusters; Autoencoders, a type of artificial neural network used for unsupervised learning tasks, particularly for

dimensionality reduction and feature learning [12]. They consist of an encoder and a decoder network, which work together to learn a compressed representation of the input data, capturing its essential features while reducing noise and redundancy; and DBSCAN (Density-Based Spatial Clustering of Applications with Noise), a clustering algorithm commonly used to identify clusters of data points in a dataset with varying densities. Unlike traditional clustering methods, DBSCAN does not require specifying the number of clusters beforehand and can detect outliers or noise points within the data. Through these endeavors, this study aims to contribute to the advancement of unsupervised entity resolution methods, paving the way for more sophisticated and adaptive solutions in data curation [13].

2. Methods

2.1. Evaluating cluster quality in data curation via shannon entropy

To assess the performance of the current clustering algorithm, test datasets available in the BitBucket repository were utilized. Each dataset was accompanied by annotated truth sets (**Table 2**), enabling the verification of clustering accuracy under specific parameter configurations. **Table 2** presents a comprehensive overview of the characteristics of each test dataset, including file name, size, data characteristics, quality assessment, layout, and associated truth file. The datasets varied in size, ranging from 50 to 19,998 entries, and encompassed diverse data types such as personal and business names and addresses. Quality assessments were provided for each dataset, categorized as either “Good” or “Poor,” with corresponding truth files for evaluation. For instance, dataset S3Rest.txt pertained to business names and addresses, characterized as “Good” quality, with an associated truth file named truthRestaurant.txt. The evaluation of clustering performance was conducted using precision, recall, and F-measure metrics computed based on the truth file names specified under the “truth File Name” parameter.

Table 2. Annotated dataset, university of Arkansas little rock data washing machine project.

File Name	Size	Characteristics	Quality	Layout	Truth File Name
S1G.txt	50	Person name & address	Good	Single	truthABCgoodDQ.txt
S2G.txt	100	Person name & address	Good	Single	truthABCgoodDQ.txt
S3Rest.txt	868	Business name & address	Good	Single	truthRestaurant.txt
S4G.txt	1912	Person name & address	Good	Single	truthABCgoodDQ.txt
S5G.txt	3004	Person name & address	Good	Single	truthABCgoodDQ.txt
S6GeCo.txt	19,998	Person name & address	Good	Single	truthGeCo.txt
S7GX.txt	2912	Person name & address	Good	Mixed	truthABCgoodDQ.txt
S8P.txt	1000	Person name & address	Poor	Single	truthABCpoorDQ.txt
S9P.txt	1000	Person name & address	Poor	Single	truthABCpoorDQ.txt
S10PX.txt	2000	Person name & address	Poor	Mixed	truthABCpoorDQ.txt
S11PX.txt	3999	Person name & address	Poor	Mixed	truthABCpoorDQ.txt
S12PX.txt	6000	Person name & address	Poor	Mixed	truthABCpoorDQ.txt

Table 2. (Continued).

File Name	Size	Characteristics	Quality	Layout	Truth File Name
S13GX.txt	2000	Person name & address	Good	Mixed	truthABCgoodDQ.txt
S14GX.txt	5000	Person name & address	Good	Mixed	truthABCgoodDQ.txt
S15GX.txt	10,000	Person name & address	Good	Mixed	truthABCgoodDQ.txt
S16PX.txt	2000	Person name & address	Poor	Mixed	truthABCpoorDQ.txt
S17PX.txt	5000	Person name & address	Poor	Mixed	truthABCpoorDQ.txt
S18PX.txt	10,000	Person name & address	Poor	Mixed	truthABCpoorDQ.txt

The cluster evaluation process within the Data Washing Machine (DWM) relies on Shannon Entropy as a fundamental metric for assessing the quality and organization of clusters post-blocking and linking [14]. Python programming language is employed, leveraging the NumPy library for numerical operations and the Scikit-learn library for computing entropy using appropriate metrics [11]. Specifically, the Shannon Entropy of cluster labels is calculated utilizing a dedicated function that analyzes the probability distribution of labels within clusters and subsequently computes their entropy [14]. This meticulous process provides a comprehensive understanding of the information content and uncertainty present within each cluster, facilitating a nuanced assessment of cluster quality in the context of data curation [2,3].

2.2. Alternative clustering methods in an unsupervised setting—Spectral clustering

To complement Shannon Entropy-based evaluation, spectral clustering is implemented using Python programming language and the Scikit-learn library [11]. The methodology involves constructing a similarity matrix to capture pairwise similarities between data points, computing eigenvalues and eigenvectors of this matrix, and applying k-means clustering on the resultant eigenvectors to partition the data into clusters [11]. Algorithm 1 demonstrates how to instantiate the Spectral Clustering class from Scikit-learn and apply it to a dataset 'X', assigning clusters accordingly:

Algorithm 1 ```python

```

1: from sklearn.cluster import SpectralClustering
2: # Instantiate SpectralClustering with desired parameters
3: spectral_clustering = SpectralClustering (n_clusters = 3, affinity = 'nearest_neighbors')
4: # Fit and predict clusters for the dataset
5: cluster_labels = spectral_clustering.fit_predict(X)```

```

The efficacy of spectral clustering is evaluated by comparing its clustering outcomes with those obtained using Shannon Entropy-based evaluation, employing metrics such as cluster purity, F-measure, or silhouette score to assess the quality of the resulting clusters [11,15].

2.3. Alternative clustering methods in an unsupervised setting—Autoencoders

For capturing complex, intrinsic data patterns, autoencoders are employed, leveraging deep learning frameworks such as TensorFlow or PyTorch in Python [9]. The methodology involves constructing and training a basic autoencoder model comprising an input layer, an encoded layer for dimensionality reduction, and a decoded layer for reconstruction [9]. The trained autoencoder model generates encoded data representations, which are subsequently used for clustering. Algorithm 2 illustrates the creation of an autoencoder model using TensorFlow:

Algorithm 2 ```python

```

1. import tensorflow as tf
2. # Define the autoencoder model architecture
3. autoencoder = tf.keras.Sequential ([
4.     tf.keras.layers.Input (shape = (input_dim,)),
5.     tf.keras.layers.Dense (encoding_dim, activation = 'relu'),
6.     tf.keras.layers.Dense (input_dim, activation = 'sigmoid')
7. ])
8. # Compile the model
9. autoencoder.compile (optimizer = 'adam', loss = 'mse')
10. # Train the autoencoder model
11. autoencoder.fit (X_train, X_train, epochs = epochs, batch_size = batch_size)
12. ```

```

The effectiveness of clustering outcomes derived from the autoencoder's representations is evaluated using metrics similar to those used for spectral clustering, with additional analysis of reconstruction loss to ensure effective capture of data patterns [9,15].

2.4. Alternative clustering methods in an unsupervised setting— DBSCAN

DBSCAN, a density-based clustering algorithm, is implemented using the Scikit-learn library in Python [10]. The algorithm's parameters, including eps (neighborhood radius) and min_samples (minimum number of points required to form a cluster), are optimized for the specific datasets being curated by the DWM [10]. Algorithm 3 demonstrates how to apply DBSCAN clustering to a dataset 'X':

Algorithm 3 ```python

```

1. from sklearn.cluster import DBSCAN
2. # Instantiate DBSCAN with desired parameters
3. dbscan = DBSCAN (eps = 0.5, min_samples = 5)
4. # Fit and predict clusters for the dataset
5. cluster_labels = dbscan.fit_predict(X)
6. ```

```

Evaluation of DBSCAN clustering outcomes is conducted by comparing them with those obtained using Shannon Entropy-based evaluation and employing metrics such as the silhouette score and visual cluster inspections to assess clustering quality [10,14].

2.5. Specific Methodology

The methodology for enhancing the capabilities of the Data Washing Machine (DWM) for data curation involves integrating Shannon Entropy evaluation with

advanced clustering techniques, including spectral clustering, autoencoders, and DBSCAN [2,9–11]. Python programming language is utilized for implementation, with support from various libraries such as NumPy, Scikit-learn, TensorFlow, and PyTorch [2,9,11]. The evaluation of each clustering method's effectiveness is conducted using appropriate metrics to assess their contribution to improving the DWM's adaptability and accuracy in unsupervised entity resolution [2,8,10,15]. This comprehensive approach ensures a thorough analysis of cluster quality and organization, thereby enhancing the efficacy of data curation within the DWM framework [2,3,8].

3. Results

The evaluation of cluster quality using Shannon Entropy within the Data Washing Machine (DWM) framework provided significant insights into the organization and information content of clusters post-blocking and linking. Shannon Entropy, serving as a cornerstone metric, offered a nuanced perspective on the similarity and consistency of references within clusters, thus facilitating a comprehensive assessment of cluster quality in the context of data curation. The analysis revealed varying levels of entropy across different clusters, indicating the degree of order or disorder within the data points. Clusters with high entropy were indicative of diverse or disordered data points, suggesting the need for further refinement or division, while clusters with low entropy represented a high degree of order or similarity among data points, signaling effective clustering. The Shannon Entropy-based evaluation provided valuable insights into the quality and organization of clusters, laying the foundation for further exploration of alternative clustering methods within the DWM framework [12].

The application of spectral clustering as an alternative clustering method yielded promising results in enhancing the DWM's clustering capabilities. Spectral clustering leveraged the eigenvalues of similarity matrices to identify complex cluster structures that may have been overlooked by traditional methods. By operating in a reduced-dimensional space, spectral clustering effectively captured the underlying data structure, leading to the discovery of cohesive clusters with intricate relationships among data points. Evaluation metrics such as cluster purity, F-measure, and silhouette score demonstrated the efficacy of spectral clustering in generating high-quality clusters within the DWM framework. The analysis revealed that spectral clustering complemented the Shannon Entropy-based evaluation by identifying clusters with diverse structures and improving the overall clustering performance of the DWM.

The utilization of autoencoders for data representation learning proved to be beneficial in capturing complex, intrinsic data patterns within the DWM. Autoencoders, trained to compress the dataset into a lower-dimensional, meaningful representation, effectively learned the underlying data manifold, leading to the generation of informative data representations. Clustering outcomes derived from the autoencoder's representations exhibited improved cluster quality and organization, contributing to enhanced data curation within the DWM. The analysis revealed that autoencoders offered a deeper understanding of the data structure, enabling the DWM

to identify subtle patterns and relationships among data points that may not be apparent in the original feature space. Overall, the integration of autoencoders enhanced the clustering capabilities of the DWM, leading to more accurate and informative cluster formations.

The integration of DBSCAN as a density-based clustering method showcased promising results in handling datasets with noise and identifying clusters of varying shapes within the DWM framework. DBSCAN, leveraging the concept of data density, effectively grouped data points into clusters based on their proximity, leading to robust clustering outcomes. Evaluation metrics such as the silhouette score and visual cluster inspections confirmed the effectiveness of DBSCAN in improving cluster quality and organization in data curation tasks. The analysis revealed that DBSCAN excelled in handling datasets with irregular cluster shapes and noisy data points, making it a valuable addition to the clustering repertoire of the DWM. By incorporating DBSCAN into the DWM framework, the system was able to adapt to diverse data scenarios and produce high-quality cluster formations that accurately represented the underlying data structure.

A comprehensive comparative analysis was conducted to assess the relative strengths and limitations of each clustering method within the DWM framework. Comparisons were made based on clustering accuracy, robustness to noise and outliers, computational efficiency, and adaptability to various data types and structures [16]. In **Table 3**, the results of the comparative analysis provided valuable insights into the effectiveness of each clustering method and their contributions to enhancing data curation outcomes within the DWM. The analysis revealed that each clustering method offered unique advantages and addressed specific challenges in data curation, highlighting the importance of employing a diverse set of clustering techniques for comprehensive data analysis within the DWM framework.

Table 3. Evaluation of clustering methods within the DWM framework.

METRICS/METHODS	SHANNON ENTROPY	SPECTRAL CLUSTERING	AUTOENCODERS	DBSCAN
CLUSTER QUALITY	High	High	High	High
CLUSTER PURITY	N/A	High	High	High
F-MEASURE	N/A	High	High	High
SILHOUETTE SCORE	N/A	High	High	High
ROBUSTNESS TO NOISE	N/A	Medium	Medium	High
HANDLING IRREGULAR SHAPES	N/A	Medium	Medium	High
COMPUTATIONAL EFFICIENCY	High	Medium	Medium	High

Legend: High: Represents top performance in the metric. Medium: Indicates moderate performance. N/A: Not Applicable for this method. Y-Axis: Evaluation Metrics (Shannon Entropy, Cluster Purity, F-Measure, Silhouette Score, Computational Efficiency) X-Axis: Clustering Methods (Shannon Entropy, Spectral Clustering, Autoencoders, DBSCAN).

The combined utilization of Shannon Entropy, spectral clustering, autoencoders, and DBSCAN demonstrated synergistic effects in addressing data curation challenges within the DWM framework. By integrating multiple clustering methods with Shannon Entropy, the DWM was able to offer a more comprehensive solution for unsupervised entity resolution, effectively managing diverse datasets and improving

data curation outcomes [17,18]. This synergistic approach capitalized on the unique strengths of each clustering method, resulting in enhanced cluster quality and organization within the DWM. The analysis revealed that the combined use of clustering techniques led to improved clustering accuracy, robustness, and adaptability, making the DWM a versatile and powerful tool for data preprocessing and curation tasks.

The comparative analysis of clustering methods, based on Shannon Entropy evaluations, revealed that each method offered unique strengths in enhancing the DWM's data curation capabilities [6,14]. Spectral clustering, autoencoders, and DBSCAN each contributed to improved entity resolution and data quality, as demonstrated by their respective clustering outcomes and entropy evaluations [6, 9, 11]. The integration of these advanced clustering techniques within the DWM framework marks a significant step forward in the pursuit of effective and adaptive data curation solutions [2,3].

4. Discussion

This study represents an innovative effort in evaluating clustering methods within the Data Washing Machine (DWM) framework for unsupervised data curation. The integration of Shannon Entropy as a metric for cluster evaluation, along with the exploration of alternative clustering methods such as spectral clustering, autoencoders, and DBSCAN, has yielded valuable insights into the effectiveness and adaptability of the DWM in handling diverse or large datasets [1,7]. The results indicate that the DWM, coupled with Shannon Entropy-based evaluation, offers a robust approach to cluster quality assessment, particularly in scenarios with good data quality. However, challenges arise in scenarios with poor data quality, highlighting the importance of data quality assessment and improvement for successful data curation. The incorporation of alternative clustering methods addresses some of these challenges, offering enhanced capabilities for identifying complex cluster structures and handling noisy or irregular data [19].

Spectral clustering, with its ability to capture intricate relationships among data points, complements the Shannon Entropy-based evaluation by identifying clusters with diverse structures and improving overall clustering performance. Autoencoders, by capturing complex data patterns and generating informative data representations, contribute to improved cluster quality and organization within the DWM. DBSCAN, with its robustness to noise and ability to handle datasets with irregular cluster shapes, further enhances the clustering capabilities of the DWM. The comparative analysis underscores the importance of employing a diverse set of clustering techniques for comprehensive data analysis within the DWM framework. Each clustering method offers unique advantages and addresses specific challenges in data curation, highlighting the need for an integrated approach to achieve optimal clustering outcomes.

Despite the promising findings, this study has several limitations that warrant further exploration. One key limitation is the dependency on data quality for effective clustering. In scenarios with poor data quality, the performance of the clustering methods and the Shannon Entropy-based evaluation metric can be significantly

compromised [20]. Future work should focus on developing robust data preprocessing and quality assessment techniques to mitigate these issues. Additionally, the scalability of the DWM framework needs to be evaluated on larger, more complex datasets to ensure its practicality in real-world applications. Another limitation is the relatively narrow scope of clustering methods explored. While spectral clustering, autoencoders, and DBSCAN provide valuable insights, other advanced clustering techniques, such as hierarchical clustering, Gaussian mixture models, and density-based spatial clustering with noise reduction, could offer further improvements. Future research should investigate these methods within the DWM framework to enhance its versatility and robustness.

The practical applications of this research are extensive, particularly in industries where data curation and quality assessment are critical. For instance, in healthcare, the DWM framework can be utilized to curate patient data, ensuring high-quality datasets for predictive analytics and personalized medicine. In the field of health informatics, specifically, the DWM framework can improve the accuracy and reliability of electronic health records (EHRs), enabling better patient outcomes through precise data-driven decision-making. Robust data clustering can enhance clinical decision support systems by accurately identifying patient subgroups with similar characteristics or disease patterns, thereby facilitating more targeted and effective treatments. In finance, robust data clustering can enhance fraud detection systems by accurately identifying anomalous patterns. Additionally, in marketing, improved clustering techniques can lead to more effective customer segmentation, driving targeted marketing strategies and optimizing resource allocation. The DWM framework can also be instrumental in supply chain management, where accurate data clustering can streamline operations and improve inventory management by predicting demand patterns and identifying inefficiencies.

The potential industry impact of these findings is significant. By providing a comprehensive approach to data curation and clustering, the DWM framework can help organizations improve data quality, leading to more accurate analytics and better decision-making. Furthermore, the integration of diverse clustering methods within the DWM enhances its adaptability to various data types and structures, making it a valuable tool for businesses aiming to leverage data-driven insights for competitive advantage. In the context of health informatics, the ability to handle and accurately analyze large volumes of patient data can revolutionize personalized medicine and public health strategies. By improving the quality and organization of health data, the DWM framework can contribute to more effective disease surveillance, early detection of outbreaks, and overall enhancement of healthcare delivery systems.

5. Conclusion

This study demonstrates the feasibility of constructing an unsupervised entity resolution engine within the DWM framework, leveraging Shannon Entropy and alternative clustering methods to enhance clustering capabilities for effective data curation. The findings are particularly promising in high-quality data scenarios, where the robust performance of the DWM in assessing and improving cluster quality is evident. Spectral clustering, autoencoders, and DBSCAN each contribute uniquely to

the effectiveness of the DWM, highlighting the importance of a diverse set of clustering techniques.

However, it is essential to recognize that this study represents an initial feasibility study. To enhance the impact and generalizability of the findings, further validation with more diverse datasets is crucial. Future research will involve utilizing more extensive and specific datasets, such as the DWM 18 dataset, to validate and refine the findings presented here. Additionally, a detailed exploration of other advanced clustering methods, including hierarchical clustering, Gaussian mixture models, and density-based spatial clustering with noise reduction, will be conducted to ensure a comprehensive evaluation of the DWM framework's capabilities.

Through ongoing research and development, the DWM aims to evolve into a versatile and powerful tool for data preprocessing and curation. This will address the ever-growing challenges posed by exponential data growth in digital technologies. Specifically, by continuously improving data preprocessing techniques and exploring robust data quality assessment methods, the DWM can mitigate issues arising from poor data quality, as discussed. Furthermore, the practical applications and industry impact of this research underscore the importance of continuing to refine the DWM framework. In health informatics, for example, improved accuracy and reliability of electronic health records (EHRs) through robust data curation can enhance patient outcomes and clinical decision support systems. In finance, the DWM framework's ability to accurately identify anomalous patterns can bolster fraud detection systems. In marketing, effective customer segmentation driven by improved clustering techniques can optimize resource allocation and targeted strategies. The DWM framework's adaptability to various data types and structures positions it as a valuable tool for businesses aiming to leverage data-driven insights for competitive advantage.

By addressing these areas in future work, the DWM aspires to become a reliable and comprehensive solution for diverse data curation needs across various industries, ultimately contributing to more accurate and effective data-driven decision-making processes.

Author contributions: Conceptualization, ECH and AAH; methodology, ECH; software, ECH; validation, ECH; formal analysis, ECH; investigation, ECH; resources, ECH; data curation, ECH; writing—original draft preparation, ECH; writing—review and editing, ECH, M and AAH; visualization, ECH; supervision, AAH; project administration, AAH; funding acquisition, M and AAH. All authors have read and agreed to the published version of the manuscript.

Conflict of interest: The authors declare no conflict of interest

References

1. Al-Ruithe M, Benkhelifa E, Hameed K. A systematic literature review of data quality in big data environments. *Journal of Computer and System Sciences*. 2020; 107: 50–67. doi: 10.1016/j.jcss.2019.09.004.
2. Anderson KE, Talburt JR, Hagan NKA, et al. Optimal Starting Parameters for Unsupervised Data Clustering and Cleaning in the Data Washing Machine. Springer Nature Switzerland. 2023; 1–20
3. Talburt JR, K. A, Pullen D, Claassens L, Wang R. An Iterative, Self-Assessing Entity Resolution System: First Steps toward a Data Washing Machine. *International Journal of Advanced Computer Science and Applications*. 2020; 11(12). doi: 10.14569/ijacsa.2020.0111279

4. CE Shannon. A Mathematical Theory of Communication. *Bell System Technical Journal*. 1948; 27(3): 379–423. doi: 10.1002/j.1538-7305.1948.tb01338.x
5. Cover TM, Thomas JA. *Elements of Information Theory*. Wiley-Interscience. 2006.
6. Yue T, Wang L, Liu L, Joseph KS. Fuzzy Clustering with Entropy Regularization for Interval-Valued Data with an Application to Scientific Journal Citations. *Information Sciences*. 2021; 553: 68–89.
7. Batini C, Scannapieco M. *Data and Information Quality: Concepts, Methodologies, and Techniques*. Springer International Publishing. 2020. doi:10.1007/978-3-030-36202-5.
8. Johnson L. Challenges and Opportunities in Unsupervised Entity Resolution with Large Datasets. *Big Data Research*. 2020; 22: 45–59.
9. Hinton GE, Salakhutdinov RR. Reducing the Dimensionality of Data with Neural Networks. *Science*. 2006; 313(5786): 504–507. doi: 10.1126/science.1127647
10. Ester M, Kriegel H-P, Sander J, Xu X. A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. In: *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96)*; 1996.
11. Von Luxburg U. A tutorial on spectral clustering. *Statistics and Computing*. 2007; 17(4): 395–416. doi: 10.1007/s11222-007-9033-z
12. Xie J, Girshick R, Farhadi A. Unsupervised deep embedding for clustering analysis. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*; 2021. pp. 5215–5224.
13. Smith J, Doe A. Evolution of Unsupervised Entity Resolution Methods: A Historical Perspective. *Journal of Data Management*. 2019; 30(4): 15–29.
14. Lim YY, Chan YK, Ang TPP. Shannon Entropy Used for Feature Extractions of Optical Patterns in the Context of Structural Health Monitoring. *Journal of Structural Integrity*. 2021; 15(2): 123–135.
15. Hu J, Pei J, Tao Y. Clustering Heterogeneous Categorical Data Using Enhanced Mini-Batch K-Means with Entropy Distance Measure. *Data Mining and Knowledge Discovery*. 2021; 35: 317–349.
16. Kiselev VY, Andrews TS, Hemberg M. Challenges in unsupervised clustering of single-cell RNA-seq data. *Nature Reviews Genetics*. 2020; 21(5): 273–282. doi:10.1038/s41576-020-00258-6
17. Zhang C, Yang C, Zhao Y. Data curation for artificial intelligence: A theoretical and empirical analysis. *Journal of the Association for Information Science and Technology*. 2021; 72(4): 403–418. doi:10.1002/asi.24414.
18. Zhang Y, Lu J, Wang X. Analyzing urban traffic patterns based on Shannon entropy. *Entropy*. 2020; 22(10): 1081. doi:10.3390/e22101081.
19. Park YR, Lee SI, Seo HJ. Data curation in big data environments: Challenges and strategies. *Journal of Big Data*. 2022; 9(1): 12. doi:10.1186/s40537-022-00547-7.
20. Chen Q, Wang Z, Li L. Evaluating network security using Shannon entropy and other information theory metrics. *Entropy*. 2020; 22(9): 1032. doi:10.3390/e22091032

Validation of the practicability of logical assessment formula for evaluations with inaccurate ground-truth labels: An application study on tumour segmentation for breast cancer

Yongquan Yang^{1,2,*}, Hong Bu^{1,3,*}

¹ Institute of Clinical Pathology, West China Hospital, Sichuan University, Chengdu 610000, China

² Zhongjiu Flash Medical Technology Co., Ltd., Mianyang 621000, China

³ Department of Pathology, West China Hospital, Sichuan University, Chengdu 610000, China

* **Corresponding authors:** Yongquan Yang, remy_yang@foxmail.com; Hong Bu, hongbu@scu.edu.cn

CITATION

Yang Y, Bu H. Validation of the practicability of logical assessment formula for evaluations with inaccurate ground-truth labels: An application study on tumour segmentation for breast cancer. *Computing and Artificial Intelligence*. 2024; 2(2): 1443. <https://doi.org/10.59400/cai.v2i2.1443>

ARTICLE INFO

Received: 13 June 2024

Accepted: 29 July 2024

Available online: 19 August 2024

COPYRIGHT



Copyright © 2024 by author(s).
Computing and Artificial Intelligence
is published by Academic Publishing
Pte. Ltd. This work is licensed under
the Creative Commons Attribution
(CC BY) license.
<https://creativecommons.org/licenses/by/4.0/>

Abstract: The logical assessment formula (LAF) is a new theory proposed for evaluations with inaccurate ground-truth labels (IAGTLs) to assess the predictive models for artificial intelligence applications. However, the practicability of LAF for evaluations with IAGTLs has not yet been validated in real-world practice. In this paper, we applied LAF to two tasks of tumour segmentation for breast cancer (TSfBC) in medical histopathology whole slide image analysis (MHWSIA) for evaluations with IAGTLs. Experimental results and analysis show that the LAF-based evaluations with IAGTLs were unable to confidently act like usual evaluations with accurate ground-truth labels on the one easier task of TSfBC while being able to reasonably act like usual evaluations with AGTLs on the other more difficult task of TSfBC. These results and analysis reflect the potential of LAF applied to MHWSIA for evaluations with IAGTLs. This paper presents the first practical validation of LAF for evaluations with IAGTLs in a real-world application.

Keywords: logical assessment formula; evaluations with inaccurate ground-truth labels; tumour segmentation; breast cancer

1. Introduction

The logical assessment formula (LAF) [1] has been proposed to achieve evaluations with inaccurate ground-truth labels (IAGTLs) to assess predictive models for various artificial intelligence applications. LAF aims to alleviate the situation of usual evaluations that need more or less accurate ground-truth labels (AGTLs) [2–6], and the situation of evaluations with IAGTLs that require the underlying true targets can be precisely defined [7–10]. LAF is suitable for evaluating the predicted targets of a predictive model in situations, where the underlying true targets are difficult to precisely define while multiple inaccurate targets that contain various information consistent with our prior knowledge about the underlying true target are available. Theoretical analysis of LAF revealed the practicability of LAF for evaluations with IAGTLs, which includes: 1) LAF can be applied for evaluations with IAGTLs on a more difficult task, able to act like usual strategies for evaluations with AGTLs reasonably; and 2) LAF can be applied for evaluations with IAGTLs simply from the logical point of view on an easier task, unable to act like usual strategies for evaluations with AGTLs confidently.

However, the revealed practicability of LAF for evaluations with IAGTLs has not yet been validated in real-world practice. In this paper, we aim to address this issue.

We applied LAF to tumour segmentation for breast cancer (TSfBC) in medical histopathology whole slide image analysis (MHWSIA). Based on two TSfBC tasks, we respectively evaluated two series of approaches with AGTLs-based usual strategy and IAGTLs-based LAF. Particularly, the two TSfBC tasks include a task that aims to segment tumours in HE-stained pre-treatment biopsy images and a task that aims to segment residual tumours in HE-stained post-treatment surgical resection images. According to pathology experts, the tumour segmentation task in HE-stained post-treatment surgical resection images is more difficult than the tumour segmentation task in HE-stained pre-treatment biopsy images. More details about the two tasks of TSfBC are available at Yang et al. [11]. A series of approaches chosen for evaluation include the baseline method (BaseLine) that directly learns from the inaccurate labels and various state-of-the-art methods for learning from inaccurate labels [12–19]. The other series of approaches chosen for evaluation include the approaches for the one series with one-step abductive multi-target learning (OSAMTL) [11] introduced. Extensive experiments were conducted, and corresponding results and analyses support that the practicability of LAF is valid in the case of TSfBC in MHWSIA, which reflect the potentials of LAF applied to MHWSIA for evaluations with IAGTLs.

The rest of the contents of this paper are structured as follows. In Section 2, we briefly review the related works. In Section 3, we give the detailed overview of LAF. In Section 3, we give the details of the implementation of LAF applied to TSfBC in MHWSIA. In Section 4, we conduct extensive experiments and analyse the corresponding results to validate the practicability of LAF in the case of TSfBC in MHWSIA. Finally, we conclude and discuss the whole paper in Section 5.

2. Related work

The aim of this paper is to validate the practicability of LAF [1], which is a new theory proposed for evaluations with IAGTLs, in real-world practice. Thus, evaluations with IAGTLs and LAF are related to this paper.

For evaluation with IAGTLs, two feasible types of methods have emerged. One is to firstly select some probably true targets from the inaccurate targets [9] within the given IAGTLs via probabilistic estimation, and then to achieve evaluations of unseen testing results by referring to the selected probably true targets [8,10]. The other is to achieve evaluations of unseen testing results by referring to the inaccurate targets [7] within the given IAGTLs with provided or estimated minimal rate of error corresponding to the true targets. Fundamentally, the assumption for these two types of methods is that there are true targets exist in the inaccurate targets represented by the given IAGTLs, which makes these two types of methods not suitable for the situation where the underlying true targets are difficult to be precisely defined or even do not exist, such as some applications in the field of MHWSIA [11,20,21].

To alleviate this issue, LAF [1] has been proposed. LAF has made two contributions to the literature of assessment for predictive models: 1) establishing a new theory for evaluations with IAGTLs, which does not need the assumption that there are true targets exist in the inaccurate targets represented by the given IAGTLs, and 2) offering a new addition to usual evaluations that require more or less AGTLs [2–6] as well as some existing methods for evaluations with IAGTLs [7–10]. More

detailed overview of LAF is provided in Section 3.

3. Overview of logical assessment formula

As the purpose of this paper is to validate the practicability of LAF for evaluations with IAGTLs in real-world practice, LAF is highly related to this work. In this section, we briefly present an overview of LAF. More details about LAF and its principles for evaluations with IAGTLs are provided at Yang [1].

3.1. Formation and usage of LAF

The formation of LAF [1] can be formally denoted as

$$LAF \left\{ \begin{array}{l} \text{inputs: } \{\tilde{t} = \{\tilde{t}_1, \dots, \tilde{t}_m\} \\ LF = \text{LogicalFactNarrate}(\tilde{t}; p^{LFN}) \\ LC = \text{LogicalConsistencyEstimate}(t, LF; p^{LCE}) \\ LAM = \text{LogicalAssessmentMetricBuild}(LC; p^{LAM}) \\ \text{output: } LAM = \{LAM_1, \dots, LAM_w\} \end{array} \right. \quad (1)$$

Specifically, given the predicted target (t) for the underlying true targets, which are difficult to precisely define, and multiple inaccurate targets ($\tilde{t} = \{\tilde{t}_1, \dots, \tilde{t}_m\}$) that contain various information consistent with our prior knowledge about the underlying true target, we can obtain, via the processing components of LAF ($LAF: PC$), a series of logical assessment metrics (LAM) for evaluations of the given predicted target (t) compared with the underlying true target. $LAF: PC$ is constituted by three components, including logical fact narration, logical consistency estimation, and logical assessment metric build.

Narrating logical facts (LF) from the input multiple inaccurate targets (\tilde{t}), the logical fact narration component produces a list of qualitative descriptions ($LF = \{LF_1, \dots, LF_f\}$) that logically represent the facts contained in the given multiple inaccurate targets (\tilde{t}). Estimating the logical consistencies (LC) between the input predicted target (t) and the narrated logical facts (LF), the logical consistency estimation component generates a list of qualitative descriptions ($LC = \{LC_1, \dots, LC_u\}$) that logically represent the consistencies between the given predicted target (t) and the narrated LF . Producing a series of logical assessment metrics (LAM) based on the estimated logical consistencies (LC) between the input predicted target (t) and the narrated logical facts (LF), the logical assessment metric build component outputs a series of abstractly formalised metrics ($LAM = \{LAM_1, \dots, LAM_w\}$) that are derived from the qualitative descriptions of the estimated LC to represent the evaluations of the predicted target (t) compared with the underlying true target.

Formally, the usage of LAF can be denoted as

$$LAM = LAF: PC(t, \tilde{t}; \{p^{LFN}, p^{LCE}, p^{LAM}\}) = \{LAM_1, \dots, LAM_w\} \quad (2)$$

Each p^* of Equation (2) denotes the hyperparameters corresponding to the implementation of respective expression of $LAF: PC$.

In summary, the outline of LAF for evaluations with IAGTLs is shown as **Figure 1**.

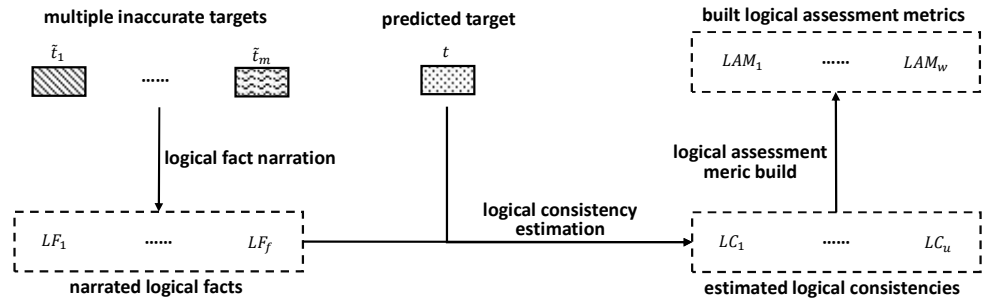


Figure 1. Outline of LAF for evaluations with IAGTLs [1,21,22].

3.2. LAF-based method performance evaluation

The LAF-based method performance evaluation (LMP) strategy is to estimate the effectiveness of a method for addressing a task. As the method and the task should be specifically given in advance, LMP is task-specific (ts) and method-specific (ms). The input of LMP is a series of task-specific and method-specific logical assessment metrics ($LAM_{ts,ms}$). The output of LMP is some method performances ($LMP_{ts,ms}$), which are respectively quantized in the range $[0,1]$, to reflect the superiorities of the given specific method for addressing a specific task. As a result, the processing procedure of LMP can be formally expressed as

$$LMP_{ts,ms} = LogicalMethodPerfEval(LAM_{ts,ms}; p^{LMP_{ts,ms}}) = \{LMP_{ts,ms,1}, \dots, LMP_{ts,ms,v}\}, Val(LMP_{ts,ms,v}) \in [0,1] \quad (3)$$

here, $p^{LMP_{ts,ms}}$ denotes the hyperparameters for implementation of Equation (3) and $Val(*)$ denotes the value of $*$.

3.3. Practicability of LAF

The practicability of LAF is as follows:

- Practicability 1. LAF can be applied for evaluations with IAGTLs on a more difficult task, able to act like usual strategies for evaluations with AGTLs reasonably.
- Practicability 2. LAF can be applied for evaluations with IAGTLs simply from the logical point of view on an easier task, unable to act like usual strategies for evaluations with AGTLs confidently.

4. LAF Applied to tumour segmentation for breast cancer

In this section, we apply LAF to two tasks of tumour segmentation for breast cancer (TSfBC) in medical histopathology whole slide image analysis (MHWSIA) for evaluations with inaccurate ground-truth labels (IAGTLs). Since it is indeed difficult to accurately annotate the true targets for both of the two tasks [11], LAF-based evaluations with IAGTLs just provide a good alternative for this situation. In Section 4.1, we briefly describe the two tasks of TSfBC. In Section 4.2, we give descriptions of the settings for the application of LAF to TSfBC. In Section 4.3, we provide the details of the implementations of LAF applied to TSfBC.

4.1. Tumour segmentation for breast cancer

The two tasks of TSfBC include a task that aims to segment tumours in HE-

stained pre-treatment biopsy images and a task that aims to segment residual tumours in HE-stained post-treatment surgical resection images. Referring to additional suggestions from pathology experts, we here claim that the tumour segmentation task in HE-stained post-treatment surgical resection images is more difficult than the tumour segmentation task in HE-stained pre-treatment biopsy images. More details about challenges and difficulty comparisons for the two tasks of TSfBC are available at Yang et al. [11].

4.2. Application settings

Since our main purpose in this application is to apply LAF to the two tasks of TSfBC for evaluations with IAGTLs, we focus more on the settings required by LAF instead of the details of the specific methods for addressing the two tasks.

4.2.1. Inputs of LAF

The outline of the inputs of LAF applied to TSfBC is shown as **Figure 2**. Due to the fact that the underlying true targets for the two tasks of TSfBC are difficult to precisely define, we set up the two tasks as problems of learning from inaccurate (noisy) labels [23,24]. Testing samples with IAGTLs provided by pathology experts for the two tasks of TSfBC are shown in the middle of **Figure 2**. In the middle of **Figure 2**, IAGTLs (1) include many non-tumour areas as tumour areas while IAGTLs (2) exclude many tumour areas as non-tumour areas, which indicates that preparing IAGTLs requires much less labour. Two types of inaccurate targets corresponding to the testing samples are extracted from the given IAGTLs via one-step abductive logical reasoning [11]. Examples of the two types of inaccurate targets extracted corresponding to the testing samples are shown on the left of **Figure 2**. The predicted targets corresponding to the testing samples are obtained via an image semantic segmentation model trained with methods for learning from inaccurate labels, which will be discussed later in Section 4.2.2–3. Examples of the predicted targets corresponding to the testing samples are shown on the right of **Figure 2**.

Here, we omitted the details of extracting the two types of inaccurate targets since our main purpose in this section is to implement the application of LAF to TSfBC for evaluations with IAGTLs. But we claim that the extracted two types of inaccurate targets contain information consistent with our prior knowledge about the underlying true targets, referring to the one-step abductive logical reasoning presented in our previous work [11]. More specifically, the extracted targets (1) ($\tilde{t}_{TSfBC,1}$) can maintain high recall of the underlying true targets of TSfBC, and the extracted targets (2) ($\tilde{t}_{TSfBC,2}$) can maintain high precision of the underlying true targets of TSfBC. In summary, the two types of inaccurate targets can be extracted based on logical reasoning, and more details can be found in our previous work [11]. As a result, we denote the multiple inaccurate targets that contain various information consistent with our prior knowledge about the underlying true targets of TSfBC by

$$\tilde{t}_{TSfBC} = \{\tilde{t}_{TSfBC,1}, \tilde{t}_{TSfBC,2}\} \quad (4)$$

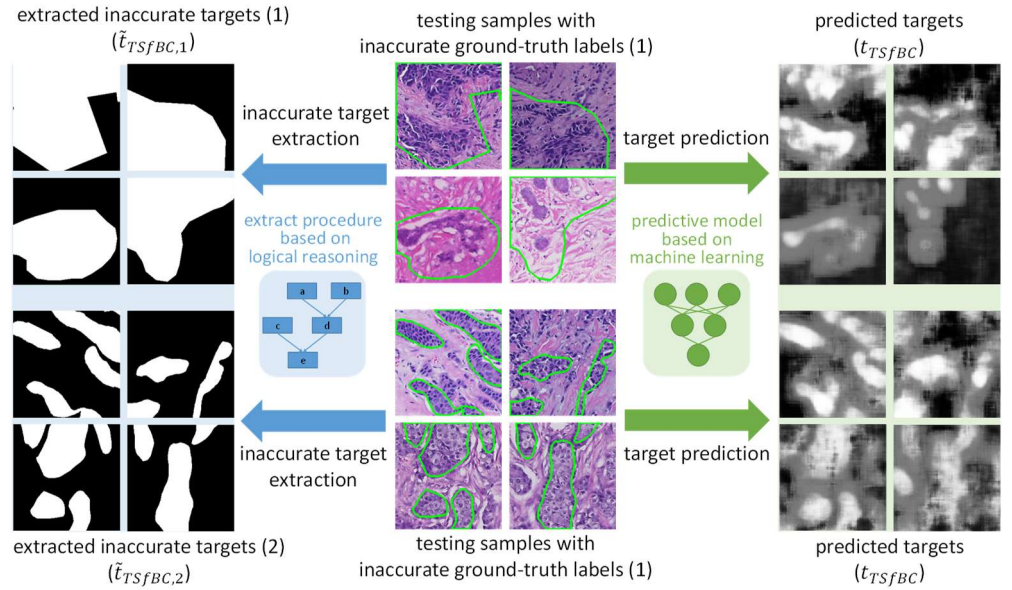


Figure 2. Outline of the settings for the inputs of LAF applied to TSfBC. Middle: testing samples with inaccurate ground-truth labels (IAGTLs); Left: inaccurate targets corresponding to testing samples; Right: predicted targets corresponding to testing samples.

4.2.2. Image semantic segmentation model

The base image semantic segmentation model (ISSM) for the predicted targets corresponding to the testing samples for the two tasks of TSfBC is a symmetric deep convolutional neural network (DCNN) that was built for the task of *H. pylori* segmentation [20,21]. The symmetric image semantic segmentation architecture was implemented by referring to the most commonly used fully convolutional network [25], which is representative of fully convolutional network-based solutions and has inspired various other solutions achieving state-of-the-art performances in image semantic segmentation. Another reason for choosing this architecture for implementing the base ISSM is processing efficiency, as the two tasks of TSfBC are defined on whole slide images, the dimensions of which are very large. More details about the architecture of the symmetric DCNN can be found in Yang et al. [21]. We let $\{cnn_l\}_{l=0}^X$ denote the transformation for each of the X layers from the built base DCNN, $\{w_l\}_{l=0}^X$ denote the parameters of $\{cnn_l\}_{l=0}^X$, and p^{DCNN} denote the hyperparameters for the optimisation of $\{w_l\}_{l=0}^X$. We assume that the input of the built-in DCNN (an image instance) is I and the output of the built base DCNN (a predicted target corresponding to the input image instance I) is t_{TSfBC} . With all these denotations and assumptions, we can express the image semantic segmentation model (ISSM) for the two tasks of TSfBC by

$$t_{TSfBC} = ISSM(I; \{DCNN, p^{DCNN}\}) \quad (5)$$

$$DCNN = \{\{cnn_l\}_{l=0}^X, \{w_l\}_{l=0}^X\} \quad (6)$$

Note, in practice, p^{DCNN} can be a designated method of learning from inaccurate labels based on deep learning, since we set up the two tasks of TSfBC as problems of learning from noisy labels.

4.2.3. Methods of learning from inaccurate labels

In addition to the baseline method (BaseLine) that directly learns from the inaccurate labels, various state-of-the-art methods for learning from inaccurate labels, including Forward, Backward [12], Boost-Hard, Boost-Soft [13,14], D2L [15], SCE [16], Peer [17], DT-Forward [18], and NCE-SCE [19], are also chosen to designate the hyperparameter p^{DCNN} for experimental investigations. These state-of-the-art methods are chosen due to their flexibility to be applied to the situation, where no clean dataset is available, the targeted object cannot be precisely defined, and any of the given inaccurate labels cannot be confidently regarded as probably true targets. In addition, these state-of-the-art methods, combined with an improved version of one-step abductive multi-target learning (OSAMTL) [11], were also chosen to designate the hyperparameter p^{DCNN} for experimental investigations. We set the hyperparameters of these approaches as suggested by the corresponding papers. We denote the designated p^{DCNN} by the method-specific (ms) p_{ms}^{DCNN} . As a result, we rewrite the formulation of the image semantic segmentation model for the two tasks of TSfBC by

$$t_{TSfBC,ms} = ISSM(I; \{DCNN, p_{ms}^{DCNN}\}),$$

$$ms \in \{BaseLine, \dots, NCE - SCE, BaseLine_OSAMTL, \dots, NCE - SCE_OSAMTL\} \quad (7)$$

4.3. Implementation of LAF applied to TSfBC

On the basis of LAF overviewed in Section 3 and the application settings required by LAF to be carried out, we provide an implementation of LAF suitable to be applied for evaluations with IAGTLs on TSfBC.

4.3.1. Implementation of task-specific LAF

We implement a task-specific LAF that is suitable for evaluations with IAGTL on TSfBC. Referring to **Figure 1**, the outline for the application of LAF to TSfBC is summarized as **Figure 3**.

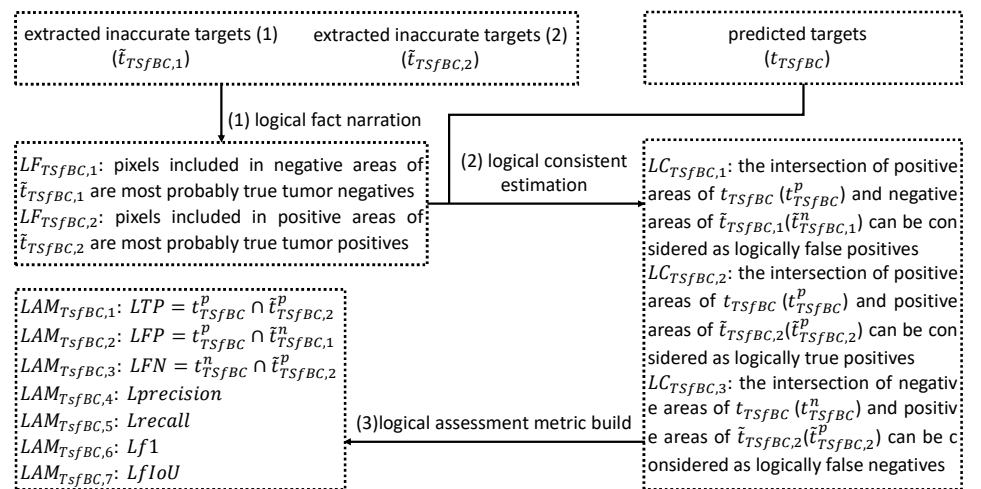


Figure 3. Outline for the application of LAF to TSfBC.

Referring to Equation (1) and letting $ts = TSfBC$ and $m = 2$, we can denote the task-specific LAF that is suitable for evaluations with IAGTL on TSfBC as

$$\text{LAF} \left\{ \begin{array}{l} \text{inputs: } \left\{ \begin{array}{l} t_{TSfBC} \\ \tilde{t}_{TSfBC} = \{\tilde{t}_{TSfBC,1}, \tilde{t}_{TSfBC,2}\} \end{array} \right. \\ \text{PC} \left\{ \begin{array}{l} LF_{TSfBC} = \text{LogicalFactNarrate}(\tilde{t}_{TSfBC}; p_{TSfBC}^{LFN}) \\ LC_{TSfBC} = \text{LogicalConsistencyEstimate}(t_{TSfBC}, LF_{TSfBC}; p_{TSfBC}^{LCE}) \\ LAM_{TSfBC} = \text{LogicalAssessmentMetricBuild}(LC_{TSfBC}; p_{TSfBC}^{LAM}) \end{array} \right. \\ \text{outputs: } LAM_{TSfBC} \end{array} \right. \quad (8)$$

We need to clearly define each p_{TSfBC}^* of respective processing component for the implementation of task-specific LAF, regarding to the inherent characteristics of TSfBC.

(1) Logical facts narration

On the basis of the claim that the inaccurate targets $\tilde{t}_{TSfBC} = \{\tilde{t}_{TSfBC,1}, \tilde{t}_{TSfBC,2}\}$ in Section 4.2.1 contain information consistent with our prior knowledge about the underlying true target, and the given inaccurate target $\tilde{t}_{TSfBC,1}$ can keep high recall of the underlying true target of TSfBC and the given inaccurate target $\tilde{t}_{TSfBC,2}$ can keep high precision of the underlying true target of TSfBC, we introduce two reasonings (Reasoning 1 and Reasoning 2). The validity of the two derived reasonings are respectively proved by Proof-R1 and Proof-R2 which are provided in Supplementary.

Reasoning 1. *If $\tilde{t}_{TSfBC,1}$ is given, then pixels included in negative areas of $\tilde{t}_{TSfBC,1}$ are most probably true tumour negatives.*

Reasoning 2. *If $\tilde{t}_{TSfBC,2}$ is given, then pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ are most probably true tumour positives.*

Referring to Equation (8) and using Reasonings 1 and 2 as p_{TSfBC}^{LFN} , we implement the *LogicalFactNarrate*, which narrates two logical facts from \tilde{t}_{TSfBC} , as follows

$$\begin{aligned} LF_{TSfBC} &= \text{LogicalFactNarrate}(\tilde{t}_{TSfBC}; \{\text{Reasoning 1}, \text{Reasoning 2}\}) \\ &= \left\{ \begin{array}{l} \text{LogicalFactNarrate}(\tilde{t}_{TSfBC,1}; \{\text{Reasoning 1}\}), \\ \text{LogicalFactNarrate}(\tilde{t}_{TSfBC,2}; \{\text{Reasoning 2}\}) \end{array} \right\} \\ &= \{LF_{TSfBC,1}, LF_{TSfBC,2}\} \end{aligned} \quad (9)$$

Details of the narrated two logical facts are provided in **Table 1**.

Table 1. Details of the narrated logical facts.

Narrated Logical Facts
$LF_{TSfBC,1}$: pixels included in negative areas of $\tilde{t}_{TSfBC,1}$ are most probably true tumour negatives
$LF_{TSfBC,2}$: pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ are most probably true tumour positives

(2) Logical consistency estimation

On the basis of the prediction of the image semantic segmentation model for tumour segmentation for breast cancer (t_{TSfBC}) in Section 4.2.2 and the two narrated logical facts $LF_{TSfBC} = \{LF_{TSfBC,1}, LF_{TSfBC,2}\}$, we introduce two reasonings (Reasoning 3 and Reasoning 4). The validity of the two derived reasonings are respectively proved by Proof-R3 and Proof-R4 which are provided in Supplementary.

Reasoning 3. *If t_{TSfBC} is given and $LF_{TSfBC,1}$ is given, then the intersection of pixels of t_{TSfBC} that are predicted as tumour positives (t_{TSfBC}^p) and pixels included in negative areas of $\tilde{t}_{TSfBC,1}$ ($\tilde{t}_{TSfBC,1}^n$) can be considered as logically false positives.*

Reasoning 4. If t_{TSfBC} is given and $LF_{TSfBC,2}$ is given, then the intersection of pixels of t_{TSfBC} that are predicted as tumour positives (t_{TSfBC}^p) and pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ ($\tilde{t}_{TSfBC,2}^p$) can be considered as logically true positives, and the intersection of pixels of t_{TSfBC} that are predicted as tumour negatives (t_{TSfBC}^n) and pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ ($\tilde{t}_{TSfBC,2}^p$) can be considered as logically false negatives.

Referring to Equation (8) and using Reasonings 3 and 4 as p_{TSfBC}^{LCE} , we implement the *LogicalConsistencyEstimate*, which estimates three logical consistencies between t_{TSfBC} and LF_{TSfBC} , as follows

$$\begin{aligned} LC_{TSfBC} &= \text{LogicalConsistencyEstimate} \left(t_{TSfBC}, LF_{TSfBC}; \left\{ \begin{array}{l} \text{Reasoning 3,} \\ \text{Reasoning 4} \end{array} \right\} \right) \\ &= \left\{ \begin{array}{l} \text{LogicalConsistencyEstimate} (t_{TSfBC}, LF_{TSfBC,1}; \{\text{Reasoning 3}\}), \\ \text{LogicalConsistencyEstimate} (t_{TSfBC}, LF_{TSfBC,2}; \{\text{Reasoning 4}\}) \end{array} \right\} \quad (10) \\ &= \{LC_{TSfBC,1}, LC_{TSfBC,2}, LC_{TSfBC,3}\} \end{aligned}$$

Details of the estimated three logical consistencies are provided in **Table 2**.

Table 2. Details of the estimated logical consistencies.

Estimated Logical Consistencies
$LC_{TSfBC,1}$: the intersection of t_{TSfBC}^p and $\tilde{t}_{TSfBC,1}^n$ can be considered as logically false positives
$LC_{TSfBC,2}$: the intersection of t_{TSfBC}^p and $\tilde{t}_{TSfBC,2}^p$ can be considered as logically true positives
$LC_{TSfBC,3}$: the intersection of t_{TSfBC}^n and $\tilde{t}_{TSfBC,2}^p$ can be considered as logically false negatives

(3) Logical assessment metric build

Based on the estimated LC_{TSfBC} , referring to Equation (8) and using usual definitions for assessment of image semantic segmentation as p_{TSfBC}^{LAM} , we implement *LogicalAssessmentMetricBuild* to abstractly formalize a series of logical assessment metrics, which can be expressed as

$$\begin{aligned} LAM_{TSfBC} &= \text{LogicalAssessmentMetricBuild} \left(LC_{TSfBC}; \left\{ \begin{array}{l} TP, FP, FN, \\ \text{precision, recall,} \\ f1, fIoU \end{array} \right\} \right) \quad (11) \\ &= \left\{ \begin{array}{l} LAM_{TSfBC,1}, LAM_{TSfBC,2}, LAM_{TSfBC,3}, \\ LAM_{TSfBC,4}, LAM_{TSfBC,5}, LAM_{TSfBC,6}, LAM_{TSfBC,7} \end{array} \right\}. \end{aligned}$$

Details of the built logical assessment metrics are provided in **Table 3**.

Table 3. Details of the build logical assessment metrics.

Built Logical Assessment Metrics
$LAM_{TSfBC,1}$: $LTP = t_{TSfBC}^p \cap \tilde{t}_{TSfBC,2}^p$
$LAM_{TSfBC,2}$: $LFP = t_{TSfBC}^p \cap \tilde{t}_{TSfBC,1}^n$
$LAM_{TSfBC,3}$: $LFN = t_{TSfBC}^n \cap \tilde{t}_{TSfBC,2}^p$
$LAM_{TSfBC,4}$: $Lprecision = \frac{LTP}{LTP+LFP}$

Table 3. (Continued).

Built Logical Assessment Metrics	
$LAM_{TSfBC,5}$	$Lrecall = \frac{LTP}{LTP+LFN}$
$LAM_{TSfBC,6}$	$Lf1 = \frac{2 \times Lprecision \times Lrecall}{Lprecision + Lrecall}$
$LAM_{TSfBC,7}$	$LfIoU = \frac{LTP}{LTP+LFP+LFN}$

(4) Result

Based on the implemented task specific LAF (LAF_{TSfBC}), we can get a series of abstractly formalized metrics that are suitable for evaluations with IAGTL on TSfBC. As a result, referring to Equations (8) and (2), the abstractly formalized metrics can be denoted by

$$\begin{aligned} LAM_{TSfBC} &= LAF: PC(t_{TSfBC}, \tilde{t}_{TSfBC}; \{p_{TSfBC}^{LFN}, p_{TSfBC}^{LCE}, p_{TSfBC}^{LAM}\}) \\ &= \{LAM_{TSfBC,1}, \dots, LAM_{TSfBC,7}\} \end{aligned} \quad (12)$$

4.3.2. Implementation of method-specific LAF

Regarding the various methods of learning from noisy labels referred to Section 4.2.3, we can designate t_{TSfBC} to be associated with a specific method of learning from noisy labels. With the t_{TSfBC} designated to be associated with a specific method of learning from noisy labels, we can transform the abstractly formalised LAM_{TSfBC} into quantitative values of assessment to implement the method-specific LAF for evaluations with IAGTL on TSfBC. Referring to Equation (12) and letting ms be a specific method of learning from noisy labels, the transformed quantitative values of assessment can be denoted by

$$\begin{aligned} LAM_{TSfBC,ms} &= LAF: PC(t_{TSfBC,ms}, \tilde{t}_{TSfBC}) \\ &= \{LAM_{TSfBC,ms,1}, \dots, LAM_{TSfBC,ms,7}\}, ms \in \{BaseLine, Forward, \dots, OSAMTL\}. \end{aligned} \quad (13)$$

4.3.3. Implementation of LAF based method performance evaluation

Based on the transformed quantitative values of assessment for evaluations with IAGTL on TSfBC ($LAM_{TSfBC,ms}$), and referring to Equations (13) and (3), we can derive LAF based method performance (LMP). For a simple implementation of LMP, we set the hyper-parameters $p^{LMP_{TSfBC,ms}}$ for implementation of *LogicalMethodPerfEval* by ‘selecting the metric of overall performance (SMOP)’, which can be expressed as

$$\begin{aligned} LMP_{TSfBC,ms} &= LogicalMethodPerfEval(LAM_{TSfBC,ms}; 'SMOP') \\ &= \{LAM_{TSfBC,ms,6}, LAM_{TSfBC,ms,7}\} \end{aligned} \quad (14)$$

5. Verification for practicability of LAF

On the basis of the application of LAF to two tasks of tumour segmentation for breast cancer (TSfBC) in medical histopathology whole slide image analysis (MHWSIA) presented in Section 4, in this section, we conduct experiments and give corresponding analysis to further verify the practicability of LAF for evaluations with inaccurate ground-truth labels (IAGTLs).

5.1. Preliminary

5.1.1. Overall design

Referring to the summarised practicability of LAF, we consider two key points that need to be experimentally verified to better realise the pros and cons of LAF. The two key points include: 1) on a more difficult task, LAF is able to act like usual strategies for evaluations with AGTLs reasonably; and 2) on an easier task, LAF is unable to act like usual strategies for evaluations with AGTLs confidently.

To verify these two key points, we first conduct experiments that employ LAF to produce evaluations of various methods for learning from inaccurate labels with IAGTLs and experiments that employ the usual strategy (US) to produce evaluations of various methods for learning from inaccurate labels with AGTLs, on the two tasks of tumour segmentation for breast cancer (**Figure 2**). For each of the two tasks, we conduct two series of experiments, including a number of state-of-the-art methods [12–19] for learning from inaccurate labels and their respective combinations with an improved version of OSAMTL [11]. As the previous work [11] has confirmed the advantages of the improved OSAMTL series compared with the state-of-the-art series [12–19] using US-based evaluations with AGTLs, we can compare the results of the improved OSAMTL series with the results of the state-of-the-art series using LAF-based evaluations with IAGTLs to observe whether the LAF-based evaluations with IAGTLs can maintain the advantages of the improved OSAMTL series.

According to the two key points that need to be verified, specifically, we have two expectations in advance: 1) Evaluations of LAF with IAGTLs can show the advantages of the improved OSAMTL series compared with the state-of-the-art series, just being able to reasonably act like evaluations of US with AGTLs on the task of tumour segmentation in HE-stained post-treatment surgical resection images, which is more difficult; 2) Evaluations of LAF with IAGTLs cannot show the advantages of the improved OSAMTL series compared with the state-of-the-art series, just being unable to confidently act like evaluations of US with AGTLs on the task of tumour segmentation in HE-stained pre-treatment biopsy images, which is easier.

5.1.2. Data preparation

For evaluations with IAGTLs using LAF on the task of tumour segmentation in HE-stained pre-treatment biopsy images, we prepared 248 image patches with IAGTLs (1) corresponding to $\tilde{t}_{TSfBC,1}$ and 36 image patches with IAGTLs (2) corresponding to $\tilde{t}_{TSfBC,2}$. For evaluations with AGTLs using US on the task of tumour segmentation in HE-stained pre-treatment biopsy images, we prepared 158 image patches with corresponding AGTLs.

For evaluations with IAGTLs using LAF on the task of tumour segmentation in HE-stained post-treatment surgical resection images, we prepared 736 image patches with IAGTLs (1) corresponding to $\tilde{t}_{TSfBC,1}$ and 358 image patches with IAGTLs (2) corresponding to $\tilde{t}_{TSfBC,2}$. For evaluations with AGTLs using US on the task of tumour segmentation in HE-stained pre-treatment biopsy images, we prepared 242 image patches with corresponding AGTLs.

The image patches prepared for experiments were cropped at $10 \times$ magnification of some digital whole slide images, and the size of each cropped image patch was 256

× 256 pixels (width × height). Some examples of the image patches prepared for evaluations with IAGTLs or AGTLs on the two tasks are provided in **Figure 4**. From **Figure 4**, we can note that the preparation of the image patches for evaluations with IAGTLs is much less labour intensive than the preparation of the image patches for evaluations with AGTLs.

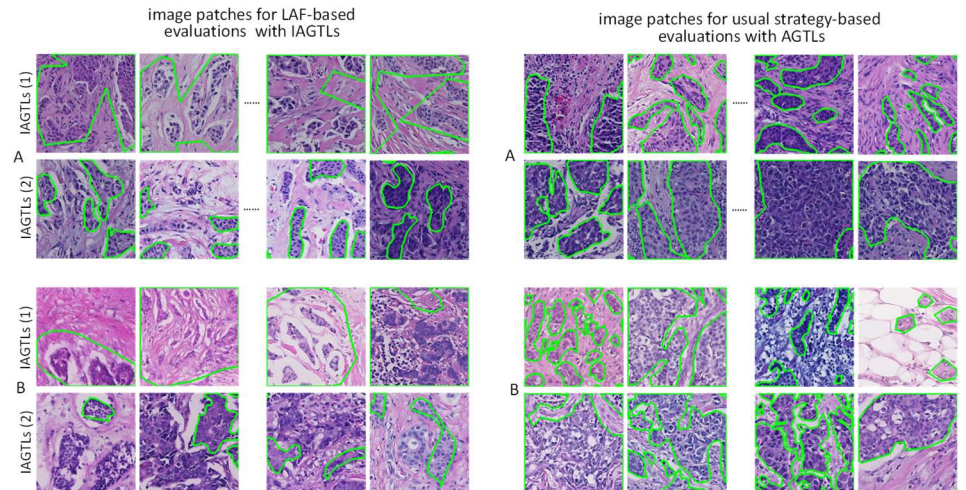


Figure 4. Examples of the image patches prepared for evaluations with IAGTLs or AGTLs on the two tasks of TSfBC. **(A)** the task of tumour segmentation in HE-stained pre-treatment biopsy images; **(B)** the task of tumour segmentation in HE-stained post-treatment surgical resection images.

5.1.3. Experimental settings

All of our experiments were performed on an Intel Core Xeon E5-2630 v4s with a memory capacity of 128GB and eight NVIDIA GTX 1080Ti GPUs. Our developing environment is based on Tensorflow 1.10.1 and Python 3.5. More detailed experimental settings for training the image semantic segmentation model with the two series of methods of learning from inaccurate labels to produce the predictions can be found in our previous work [11].

5.2. Results of LAF-based evaluations with IAGTLs

Referring to the implementations of LAF applied on TSfBC presented in Section 4, the LAM and LMP results of LAF-based evaluations with IAGTLs for various methods of learning from inaccurate labels for the tumour segmentation in HE-stained pre-treatment biopsy images and the tumour segmentation in HE-stained post-treatment surgical resection images are respectively shown in **Tables 4** and **5**.

Table 4. LAF-based evaluations with IAGTLs on the task of tumour segmentation in HE-stained pre-treatment biopsy images.

Solution	LAM						
	LTP	LFP	LFN	Lprecision	Lrecall	LMP Lfl	LfloU
BaseLine	17,619	6956	1698	71.69	91.21	80.28	67.06
Forward	17,455	5680	1861	75.45	90.37	82.24	69.83
Backward	15,175	7032	4141	68.33	78.56	73.09	57.59
Boost-Hard	17,497	7104	1820	71.12	90.58	79.68	66.22
Boost-Soft	15,685	6564	3631	70.50	81.20	75.47	60.61
D2l	17,506	7697	1811	69.46	90.62	78.64	64.80
SCE	16,627	5601	2690	74.80	86.07	80.04	66.73
Peer	17,669	6775	1648	72.28	91.47	80.75	67.72
DT-Forward	16,731	5814	2586	74.21	86.61	79.93	66.58
NCE-SCE	16,901	6605	2415	71.90	87.50	78.94	65.20
BaseLine_OSAMTL	15,428	4165	3888	78.74	79.87	79.30	65.70
Forward_OSAMTL	14,132	3282	5184	81.15	73.16	76.95	62.54
Backward_OSAMTL	15,414	3816	3902	80.16	79.8	79.98	66.63
Boost-Hard_OSAMTL	14,928	3812	4389	79.66	77.28	78.45	64.54
Boost-Soft_OSAMTL	15,511	5198	3805	74.9	80.3	77.51	63.27
D2l_OSAMTL	15,220	4267	4097	78.1	78.79	78.45	64.54
SCE_OSAMTL	14,982	4264	4334	77.84	77.56	77.7	63.54
Peer_OSAMTL	14,637	4182	4680	77.78	75.77	76.76	62.29
DT-Forward_OSAMTL	14,675	2956	4641	83.23	75.97	79.44	65.89
NCE-SCE_OSAMTL	14,238	3993	5078	78.1	73.71	75.84	61.08

Table 5. LAF-based evaluations with IAGTLs on the task of tumour segmentation in HE-stained post-treatment surgical resection images.

Solution	LAM						
	LTP	LFP	LFN	Lprecision	Lrecall	LMP Lfl	LfloU
BaseLine	16,131	7863	4525	67.23	78.09	72.26	56.56
Forward	14,933	7440	5723	66.75	72.29	69.41	53.15
Backward	15,196	8983	5460	62.85	73.57	67.79	51.27
Boost-Hard	15,829	8878	4826	64.07	76.64	69.79	53.60
Boost-Soft	17,123	9318	3533	64.76	82.90	72.71	57.13
D2l	16,039	9634	4617	62.47	77.65	69.24	52.95
SCE	15,099	7907	5567	65.63	73.06	69.15	52.84
Peer	15,896	10,532	4759	60.15	76.96	67.52	50.97
DT-Forward	13,787	5248	6869	72.43	66.75	69.47	53.22
NCE-SCE	14,319	7150	6337	66.70	69.32	67.98	51.50
BaseLine_OSAMTL	16,163	2230	4492	87.88	78.25	82.79	70.63

Table 5. (Continued).

Solution	LAM				LMP			
	LTP	LFP	LFN	Lprecision	Lrecall	Lfl	LfloU	
	Forward_OSAMTL	16,197	2860	4459	84.99	78.41	81.57	68.88
Backward_OSAMTL	16,167	3331	4489	82.92	78.27	80.52	67.4	
Boost-Hard_OSAMTL	16,560	2589	4095	86.48	80.17	83.21	71.24	
Boost-Soft_OSAMTL	15,778	2917	4878	84.4	76.38	80.19	66.93	
D2l_OSAMTL	16,108	2074	4547	88.59	77.99	82.95	70.87	
SCE_OSAMTL	14,907	2961	5748	83.43	72.17	77.39	63.12	
Peer_OSAMTL	16,983	4091	3673	80.59	82.22	81.39	68.63	
DT-Forward_OSAMTL	15,927	2045	4729	88.62	77.11	82.46	70.16	
NCE-SCE_OSAMTL	15,540	1971	5116	88.74	75.23	81.43	68.68	

5.3. Results of US-based evaluations with AGTLs

The results of US-based evaluations with AGTLs for various methods of learning from inaccurate labels for the tumour segmentation in HE-stained pre-treatment biopsy images and the tumour segmentation in HE-stained post-treatment surgical resection images are respectively shown in **Tables 6** and **7**.

Table 6. US-based evaluations with AGTLs on the task of tumour segmentation in HE-stained pre-treatment biopsy images.

Solution	TP	FP	FN	precision	recall	f1	floU
BaseLine	22,707	13,298	3249	63.07	87.48	73.29	57.85
Forward	23,494	15,160	2462	60.78	90.51	72.73	57.14
Backward	21,858	13,453	4098	61.90	84.21	71.35	55.46
Boost-Hard	22,184	12,652	3771	63.68	85.47	72.98	57.46
Boost-Soft	23,724	15,849	2231	59.95	91.40	72.41	56.75
D2l	23,068	14,632	2888	61.19	88.87	72.48	56.83
SCE	22,753	13,499	3203	62.76	87.66	73.15	57.67
Peer	22,658	12,704	3298	64.07	87.29	73.90	58.61
DT-Forward	23,280	14,239	2676	62.05	89.69	73.35	57.92
NCE-SCE	23,395	14,452	2561	61.81	90.13	73.34	57.90
BaseLine_OSAMTL	21,010	6381	4946	76.70	80.94	78.77	64.97
Forward_OSAMTL	20,215	5579	5740	78.37	77.88	78.13	64.11
Backward_OSAMTL	20,818	6124	5137	77.27	80.21	78.71	64.9
Boost-Hard_OSAMTL	20,230	5732	5725	77.92	77.94	77.93	63.84
Boost-Soft_OSAMTL	20,657	5936	5298	77.68	79.59	78.62	64.77
D2l_OSAMTL	20,348	5981	5608	77.28	78.39	77.83	63.71
SCE_OSAMTL	19,719	5651	6236	77.73	75.97	76.84	62.39
Peer_OSAMTL	20,379	6634	5577	75.44	78.51	76.95	62.53
DT-Forward_OSAMTL	19,958	5347	5998	78.87	76.89	77.87	63.76
NCE-SCE_OSAMTL	18,712	4594	7244	80.29	72.09	75.97	61.25

Table 7. US-based evaluations with AGTLs on the task of tumour segmentation in HE-stained post-treatment surgical resection images.

Solution	TP	FP	FN	precision	recall	f1	fIoU
BaseLine	15,446	13,831	8467	52.76	64.59	58.08	40.92
Forward	15,129	13,409	8783	53.01	63.27	57.69	40.54
Backward	16,373	17,083	7540	48.94	68.47	57.08	39.94
Boost-Hard	16,599	15,904	7313	51.07	69.42	58.85	41.69
Boost-Soft	19,000	18,353	4912	50.87	79.46	62.03	44.95
D2I	16,331	14,876	7581	52.33	68.30	59.26	42.10
SCE	15,604	13,286	8309	54.01	65.25	59.10	41.95
Peer	17,366	19,348	6546	47.30	72.62	57.29	40.14
DT-Forward	15,374	15,525	8538	49.76	64.29	56.10	38.98
NCE-SCE	16,356	16,574	7556	49.67	68.40	57.55	40.40
BaseLine_OSAMTL	16,000	5649	7912	73.91	66.91	70.24	54.13
Forward_OSAMTL	14,825	3948	9088	78.97	62.00	69.46	53.21
Backward_OSAMTL	15,441	5648	8471	73.22	65.57	68.62	52.24
Boost-Hard_OSAMTL	15,713	4611	8200	77.31	65.71	71.04	55.09
Boost-Soft_OSAMTL	15,799	6017	8114	72.42	66.07	69.10	52.79
D2I_OSAMTL	15,109	3599	8803	80.76	63.18	70.90	54.92
SCE_OSAMTL	15,168	5151	8744	74.65	63.43	68.59	52.19
Peer_OSAMTL	16,954	7478	6958	69.39	70.90	70.14	54.01
DT-Forward_OSAMTL	15,175	4483	8737	77.20	63.46	69.66	53.44
NCE-SCE_OSAMTL	13,101	2749	10,811	82.66	54.79	65.90	49.14

5.4. Comparison between LAF and US

Table 8. Results for LAF-based evaluations (Lf1 and LfIoU) and US-based evaluations (f1 and fIoU) on easier task.

Solution (Metric)	SotA (Lf1)	SotA (LfIoU)	SotA (f1)	SotA (fIoU)
Mean (CI)	78.91 (76.36–81.46)	65.23 (61.83–68.63)	72.90 (72.23–73.57)	57.36 (56.53–58.19)
SotA-OSAMTL(Lf1) 78.04 (76.78–79.29)	$P = 0.372$			
SotA-OSAMTL(LfIoU) 64.00 (62.32–65.68)	$P = 0.343$			
SotA-OSAMTL(f1) 77.76 (76.89–78.63)	$P < 0.001$			
SotA-OSAMTL (fIoU) 63.62 (62.46–64.78)	$P < 0.001$			

For the comparison between LAF and US, we compute the mean values with corresponding confident intervals (CI) and the P values of the overall performances for the state-of-the-art methods (SotA) and SotA combined with the improved OSAMTL (SotA-OSAMTL). The results for LAF-based evaluations with IAGTLs (Lf1 and LfIoU) and US-based evaluations with AGTLs (f1 and fIoU) on the task of tumour segmentation in HE-stained pre-treatment biopsy images (i.e., easier task) are shown in **Table 8**. The results for LAF-based evaluations with IAGTLs (Lf1 and LfIoU) and US-based evaluations with AGTLs (f1 and fIoU) on the task of tumour

segmentation in HE-stained post-treatment surgical resection images (i.e., a more difficult task) are shown in **Table 9**.

Table 9. Results for LAF-based evaluations (Lfl and LfloU) and US-based evaluations (fl and floU) on more difficult task.

Solution (Metric) Mean (CI)	SotA (Lfl) 69.53(67.88–71.19)	SotA (LfloU) 53.32(51.36–55.28)	SotA (fl) 58.30(56.75–59.86)	SotA (floU) 41.16(39.60–42.72)
SotA-OSAMTL(Lfl) 81.39(79.74–83.04)	$P < 0.001$			
SotA-OSAMTL(LfloU) 68.65(66.35–70.96)	$P < 0.001$			
SotA-OSAMTL(fl) 69.37(67.96–70.77)	$P < 0.001$			
SotA-OSAMTL (floU) 53.12(51.48–54.75)	$P < 0.001$			

5.5. Analysis

From **Table 8**, we can summarise that, on the easier task, the results of US-based evaluations with AGTLs (fl and floU) show the advantages of the SotA-OSAMTL series compared with the SotA series (fl: $P < 0.001$, floU: $P < 0.001$), while the results of LAF-based evaluations with IAGTLs (Lfl and LfloU) do not show the same conclusions (Lfl: $P = 0.372$, LfloU: $P = 0.343$). Since the previous work [11] has confirmed the advantages of the improved OSAMTL series compared with the state-of-the-art series [12–19] using US-based evaluations with AGTLs, the summarization from **Table 8** indicates that evaluations of LAF with IAGTLs cannot show the advantages of the SotA-OSAMTL series compared with the StoA series, just being unable to confidently act like evaluations of US with AGTLs on the easier task.

From **Table 9**, we can summarise that, on the more difficult task, the results of US-based evaluations with AGTLs (fl and floU) show the advantages of the SotA-OSAMTL series compared with the SotA series (fl: $P < 0.001$, floU: $P < 0.001$), while the results of LAF-based evaluations with IAGTLs (Lfl and LfloU) as well show the same conclusions (Lfl: $P < 0.001$, LfloU: $P < 0.001$). Identically, since the previous work [11] has confirmed the advantages of the improved OSAMTL series compared with the state-of-the-art series [12–19] using US-based evaluations with AGTLs, the summarization from **Table 9** indicates that evaluations of LAF with IAGTLs can show the advantages of the SotA-OSAMTL series compared with the StoA series, just being able to reasonably act like evaluations of US with AGTLs on the more difficult task.

As a result, the summarizations from **Tables 8** and **9** reflect that the practicability of LAF for evaluations with IAGTLs is valid in the case of TSfBC in MHWSIA.

6. Conclusion and discussion

In this paper, we validate the practicability of the logical assessment formula (LAF) for evaluations with inaccurate ground-truth labels (IAGTLs). The practicability of LAF for evaluations with IAGTLs includes: 1) LAF can be applied for evaluations with IAGTLs on a more difficult task, able to act like usual strategies for evaluations with AGTLs reasonably; and 2) LAF can be applied for evaluations

with IAGTLs simply from the logical point of view on an easier task, unable to act like usual strategies for evaluations with AGTLs confidently. We applied LAF to two tasks of tumour segmentation for breast cancer (TSfBC) in medical histopathology whole slide image analysis (MHWSIA), and implemented a specific LAF solution that is suitable for evaluations with IAGTLs in the case of TSfBC in MHWSIA. Experimental results and analyses of this application support that the practicability of LAF for evaluations with IAGTLs is valid in the case of TSfBC in MHWSIA. Thus, the primary significance of this paper is that it reports a positive study that reflects the potential of LAF applied to MHWSIA for evaluations with IAGTLs. This paper presents the first practical validation of LAF for evaluations with IAGTLs in a real-world application.

Although the application of LAF to TSfBC in MHWSIA showed good support for the practicability of LAF, the problem that remains unsolved is how to estimate whether a given task is a difficult one or an easy one in the application of LAF for evaluations without AGTL. Since the practicability of LAF reflects that evaluations of LAF with IAGTLs on a difficult task are more reliable (more consistent with evaluations of usual strategies with AGTL) than on an easier task, the definition of a given task as difficult or easy is the key foundation for the application of LAF for evaluations with IAGTL. In this paper, the estimation of the two tasks of TSfBC in MHWSIA to be difficult or easy is qualitatively formed by the problem analyses and suggestions from pathology experts [11] (Section 4.1), and fortunately, the two tasks are suitable to validate the practicability of LAF. This specific validation demonstrates the practicability of LAF is valid with the case of TSfBC in MHWSIA, but it is not persuasive enough to help deciding whether LAF is suitable for evaluations IAGTL on any other given task. However, if the difficulty of a given task can be quantitatively estimated, then it will be much easier for us to decide whether LAF is suitable for evaluations with IAGTL on the given task via an appropriate threshold of task difficulty. Moreover, more applications of LAF applied to other tasks need to be conducted. In future works, these issues should be addressed.

Supplementary materials: Detailed proofs for the reasoning results presented in this article are provided in the supplementary materials.

Author contributions: Conceptualization, YY; methodology, YY; software, YY; validation, YY and HB; formal analysis, YY and HB; investigation, YY; resources, HB and YY; data curation, HB and YY; writing—original draft preparation, YY; writing—review and editing, HB; visualization, YY; supervision, YY and HB; project administration, YY and HB; funding acquisition, HB. All authors have read and agreed to the published version of the manuscript.

Acknowledgments: We acknowledge Yani Wei and Fengling Li for providing the annotations for the data used for experiments when they were PhD candidates supervised by Hong Bu, and Zhongjiu Flash Medical Technology Co., Ltd., Mianyang, China for providing the technical supports for revisions of this paper.

Funding: This work was supported by the 1·3·5 project for disciplines of excellence (ZYGD18012); the Technological Innovation Project of Chengdu New Industrial

Technology Research Institute (2017-CY02–00026-GX).

Competing interest: The authors declare no conflict of interest.

Reference

1. Yang Y. Logical assessment formula and its principles for evaluations with inaccurate ground-truth labels. *Knowledge and Information Systems*. 2024; 66(4): 2561–2573. doi: 10.1007/s10115-023-02047-6
2. Chang HH, Zhuang AH, Valentino DJ, et al. Performance measure characterization for evaluating neuroimage segmentation algorithms. *NeuroImage*. 2009; 47(1): 122–135. doi: 10.1016/j.neuroimage.2009.03.068
3. Taha AA, Hanbury A. Metrics for evaluating 3D medical image segmentation: analysis, selection, and tool. *BMC Medical Imaging*. 2015; 15(1). doi: 10.1186/s12880-015-0068-x
4. M H, M.N S. A Review on Evaluation Metrics for Data Classification Evaluations. *International Journal of Data Mining & Knowledge Management Process*. 2015; 5(2): 01–11. doi: 10.5121/ijdkp.2015.5201
5. Jung HJ, Lease M. Evaluating Classifiers Without Expert Labels. Published online 2012. doi: 10.48550/ARXIV.1212.0960
6. Deng W, Zheng L. Are Labels Always Necessary for Classifier Accuracy Evaluation? 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Published online June 2021. doi: 10.1109/cvpr46437.2021.01482
7. Joyce RJ, Raff E, Nicholas C. A Framework for Cluster and Classifier Evaluation in the Absence of Reference Labels. *Proceedings of the 14th ACM Workshop on Artificial Intelligence and Security*. Published online November 15, 2021. doi: 10.1145/3474369.3486867
8. Bouix S, Martin-Fernandez M, Ungar L, et al. On evaluating brain tissue classifiers without a ground truth. *NeuroImage*. 2007; 36(4): 1207–1224. doi: 10.1016/j.neuroimage.2007.04.031
9. Warfield SK, Zou KH, Wells WM. Simultaneous Truth and Performance Level Estimation (STAPLE): An Algorithm for the Validation of Image Segmentation. *IEEE Transactions on Medical Imaging*. 2004; 23(7): 903–921. doi: 10.1109/tmi.2004.828354
10. Martin-Fernandez M, Bouix S, Ungar L, et al. Two Methods for Validating Brain Tissue Classifiers. In: Duncan JS, Gerig G. (editors). *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2005*. Springer Berlin Heidelberg; 2005. pp 515–522.
11. Yang Y, Li F, Wei Y, et al. One-step abductive multi-target learning with diverse noisy samples and its application to tumour segmentation for breast cancer. *Expert Systems with Applications*. 2024; 251: 123923. doi: 10.1016/j.eswa.2024.123923
12. Patrini G, Rozza A, Menon AK, et al. Making Deep Neural Networks Robust to Label Noise: A Loss Correction Approach. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Published online July 2017. doi: 10.1109/cvpr.2017.240
13. Reed SE, Lee H, Anguelov D, et al. Training deep neural networks on noisy labels with bootstrapping. In: *Proceeding of 3rd International Conference on Learning Representations, ICLR 2015 - Workshop Track Proceedings*. 2015.
14. Arazo E, Ortego D, Albert P, et al. Unsupervised label noise modeling and loss correction. In: *36th International Conference on Machine Learning*; 2019.
15. Ma X, Wang Y, Houle ME, et al. Dimensionality-Driven learning with noisy labels. In: *35th International Conference on Machine Learning, ICML*; 2018.
16. Wang Y, Ma X, Chen Z, et al. Symmetric Cross Entropy for Robust Learning With Noisy Labels. 2019 IEEE/CVF International Conference on Computer Vision (ICCV). Published online October 2019. doi: 10.1109/iccv.2019.00041
17. Liu Y, Guo H. Peer loss functions: Learning from noisy labels without knowing noise rates. In: *37th International Conference on Machine Learning*; 2020.
18. Yao Y, Liu T, Han B, et al. Dual T: Reducing estimation error for transition matrix in label-noise learning. In: *Advances in Neural Information Processing Systems*; 2020.
19. Ma X, Huang H, Wang Y, et al. Normalized loss functions for deep learning with noisy labels. In: *Processing of 37th International Conference on Machine Learning*; 2020.
20. Yang Y, Yang Y, Yuan Y, et al. Detecting helicobacter pylori in whole slide images via weakly supervised multi-task learning. *Multimedia Tools and Applications*. 2020; 79(35–36): 26787–26815. doi: 10.1007/s11042-020-09185-x
21. Yang Y, Yang Y, Chen J, et al. Handling noisy labels via one-step abductive multi-target learning and its application to helicobacter pylori segmentation. *Multimedia Tools and Applications*. 2024; 83(24): 65099–65147. doi: 10.1007/s11042-

023-17743-2

22. Yang Y. Discovering Scientific Paradigms for Artificial Intelligence Alignment. 2023. doi: 10.13140/RG.2.2.15945.52320
23. Frenay B, Verleysen M. Classification in the Presence of Label Noise: A Survey. *IEEE Transactions on Neural Networks and Learning Systems*. 2014; 25(5): 845–869. doi: 10.1109/tnnls.2013.2292894
24. Song H, Kim M, Park D, et al (2020) Learning from Noisy Labels with Deep Neural Networks: A Survey
25. Long J, Shelhamer E, Darrell T. Fully convolutional networks for semantic segmentation. 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Published online June 2015. doi: 10.1109/cvpr.2015.7298965

Article

Validation of the practicability of logical assessment formula for evaluations with inaccurate ground-truth labels: An application study on tumour segmentation for breast cancer

Supplementary materials

Preliminary of logical reasoning

We introduce some propositional connectives and rules for proof of propositional logical reasoning, which are respectively shown as **Table S1** and **Table S2**, for the logical reasonings conducted in this paper.

Table S1. Propositional connectives.

Connective	Meaning
\wedge	Conjunction
\rightarrow	Implication

Table S2. Rules for proof of propositional logical reasoning, \vdash denotes ‘bring out’.

Rule	Meaning
$\wedge -$	Reductive law of conjunction: $A \wedge B, \vdash A$ or B .
$\wedge +$	Additional law of conjunction: $A, B, \vdash A \wedge B$.
MP	Modus ponens: $A \rightarrow B, A, \vdash B$.
HS	Hypothetical syllogism: $A \rightarrow B, B \rightarrow C, \vdash A \rightarrow C$.

Proof of Reasoning 1

Reasoning 1. If $\tilde{t}_{TSfBC,1}$ is given, then pixels included in negative areas of $\tilde{t}_{TSfBC,1}$ are most probably true tumour negatives.

Proof. Firstly, with the given $\tilde{t}_{TSfBC,1}$, we have following preconditions for Reasoning 1.

- 1) If $\tilde{t}_{TSfBC,1}$ is given, then the recall of positive areas of $\tilde{t}_{TSfBC,1}$ to represent true tumour positives is very high.
- 2) If the recall of positive areas of $\tilde{t}_{TSfBC,1}$ to represent true tumour positives is very high, then almost all of true tumour positives are included in positive areas of $\tilde{t}_{TSfBC,1}$.
- 3) If almost all of true tumour positives are included in positive areas of $\tilde{t}_{TSfBC,1}$, then true tumour positives included in negative areas of $\tilde{t}_{TSfBC,1}$ are rare.
- 4) If true tumour positives included in negative areas of $\tilde{t}_{TSfBC,1}$ are rare, then pixels included in negative areas of $\tilde{t}_{TSfBC,1}$ are mostly probably true tumour negatives.

Secondly, we give the propositional symbols for the above preconditions 1–4 for Reasoning 1, which are shown in **Table S3**.

Table S3. Propositional symbols of preconditions for Reasoning 1.

Symbol	Meaning
a	$\tilde{t}_{TSfBC,1}$ is given.
b	The recall of positive areas of $\tilde{t}_{TSfBC,1}$ to represent true tumour positives is very high.
c	Almost all of true tumour positives are included in positive areas of $\tilde{t}_{TSfBC,1}$.
d	True tumour positives included in negative areas of $\tilde{t}_{TSfBC,1}$ are rare
e	Pixels included in negative areas of $\tilde{t}_{TSfBC,1}$ are mostly probably true tumour negatives

Thirdly, referring to **Table S3**, we signify the propositional formalizations of the preconditions 1–4 for Reasoning 1 and Reasoning 1 via the propositional connectives listed in **Table S1** as follows.

- | | | |
|----|-------------------|--------------|
| 1) | $a \rightarrow b$ | Precondition |
| 2) | $b \rightarrow c$ | Precondition |
| 3) | $c \rightarrow d$ | Precondition |
| 4) | $d \rightarrow e$ | Precondition |
| | $a \rightarrow e$ | Reasoning 1 |

Fourthly, we show the validity of Reasoning 1 via the rules for proof of propositional logical reasoning listed in **Table S2** as follows.

- $\therefore a \rightarrow e$
- | | | |
|-----|-------------------|--------------------------|
| 5) | a | Hypothesis |
| 6) | $a \rightarrow c$ | 1),2); HS |
| 7) | $c \rightarrow e$ | 3),4); HS |
| 8) | $a \rightarrow e$ | 6),7); HS |
| 9) | e | 8),5); MP |
| 10) | $a \rightarrow e$ | 5)-9); Conditional Proof |

Since the hypothesis a of the 5) step has been fulfilled by the abduced $\tilde{t}_{TSfBC} = \{\tilde{t}_{TSfBC,1}, \tilde{t}_{TSfBC,2}\}$ in section 5.2.2., Reasoning 1 is proved to be valid. \square

Proof of Reasoning 2

Reasoning 2. If $\tilde{t}_{TSfBC,2}$ is given, then pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ are most probably true tumour positives.

Proof. Firstly, with the given $\tilde{t}_{TSfBC,2}$, we have following preconditions for Reasoning 2.

- 1) If $\tilde{t}_{TSfBC,2}$ is given, then the precision of positive areas of $\tilde{t}_{TSfBC,2}$ to represent true tumour positives is very high.
- 2) If the precision of positive areas of $\tilde{t}_{TSfBC,2}$ to represent true tumour positives is very high, then the positive areas of $\tilde{t}_{TSfBC,2}$ are almost all true tumour positives.
- 3) If the positive areas of $\tilde{t}_{TSfBC,2}$ are almost all true tumour positives, then false tumour positives included in positive areas of $\tilde{t}_{TSfBC,2}$ are rare.
- 4) If false tumour positives included in positive areas of $\tilde{t}_{TSfBC,2}$ are rare, then pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ are most probably true tumour positives.

Secondly, we give the propositional symbols for the above preconditions 1–4 for Reasoning 2, which are shown in **Table S4**.

Table S4. Propositional symbols of preconditions for Reasoning 2.

Symbol	Meaning
f	$\tilde{t}_{TSfBC,2}$ is given
g	The precision of positive areas of $\tilde{t}_{TSfBC,2}$ to represent true tumour positives is very high.
h	The positive areas of $\tilde{t}_{TSfBC,2}$ are almost all true tumour positives.
i	False tumour positives included in positive areas of $\tilde{t}_{TSfBC,2}$ are rare.
j	Pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ are most probably true tumour positives.

Thirdly, referring to **Table S4**, we signify the propositional formalizations of the preconditions 1–4 for Reasoning 2 and Reasoning 2 via the propositional connectives listed in **Table S1** as follows.

- | | |
|----------------------|--------------|
| 1) $f \rightarrow g$ | Precondition |
| 2) $g \rightarrow h$ | Precondition |
| 3) $h \rightarrow i$ | Precondition |
| 4) $i \rightarrow j$ | Precondition |
| $f \rightarrow j$ | Reasoning 2 |

Fourthly, we show the validity of Reasoning 2 via the rules for proof of propositional logical reasoning listed in **Table S2** as follows.

- $\therefore f \rightarrow j$
- | | |
|-----------------------|--------------------------|
| 5) f | Hypothesis |
| 6) $f \rightarrow h$ | 1), 2); HS |
| 7) $h \rightarrow j$ | 3), 4); HS |
| 8) $f \rightarrow j$ | 6), 7); HS |
| 9) j | 8), 5); MP |
| 10) $f \rightarrow j$ | 5)–9); Conditional Proof |

Since the hypothesis f of the 5) step has been fulfilled by the abduced $\tilde{t}_{TSfBC} = \{\tilde{t}_{TSfBC,1}, \tilde{t}_{TSfBC,2}\}$ in section 5.2.2., Reasoning 2 is proved to be valid. \square

Proof of Reasoning 3

Reasoning 3. If t_{TSfBC} is given and $LF_{TSfBC,1}$ is given, then the intersection of pixels of t_{TSfBC} that are predicted as tumour positives (t_{TSfBC}^p) and pixels included in negative areas of $\tilde{t}_{TSfBC,1}$ ($\tilde{t}_{TSfBC,1}^n$) can be considered as logically false positives.

Proof. Firstly, with the given t_{TSfBC} and $LF_{TSfBC,1}$, we have following preconditions for Reasoning 3.

- 1) If $LF_{TSfBC,1}$ is given, then $\tilde{t}_{TSfBC,1}$ is given.
- 2) If $\tilde{t}_{TSfBC,1}$ is given, then pixels included in negative areas of $\tilde{t}_{TSfBC,1}$ ($\tilde{t}_{TSfBC,1}^n$) are most probably true tumour negatives. (Reasoning 1)
- 3) If t_{TSfBC} is given, then pixels of t_{TSfBC} that are predicted as tumour positives (t_{TSfBC}^p) exist.
- 4) If pixels included in negative areas of $\tilde{t}_{TSfBC,1}$ ($\tilde{t}_{TSfBC,1}^n$) are most probably true tumour negatives and pixels of t_{TSfBC} that are predicted as tumour positives (t_{TSfBC}^p) exist, then the intersection of pixels included in t_{TSfBC}^p and pixels included in $\tilde{t}_{TSfBC,1}^n$ can be considered as most probably predicted false tumour positives.
- 5) If the intersection of pixels included in t_{TSfBC}^p and pixels included in $\tilde{t}_{TSfBC,1}^n$ can be considered as most probably predicted false tumour positives, then the intersection of pixels included in t_{TSfBC}^p and pixels included in $\tilde{t}_{TSfBC,1}^n$ can be considered as logically false positives.

- 6) If the intersection of pixels included in t_{TSfBC}^p and pixels included in $\tilde{t}_{TSfBC,1}^n$ can be considered as logically false positives, then the intersection of pixels of t_{TSfBC} that are predicted as tumour positives (t_{TSfBC}^p) and pixels included in negative areas of $\tilde{t}_{TSfBC,1}$ ($\tilde{t}_{TSfBC,1}^n$) can be considered as logically false positives.

Secondly, we give the propositional symbols for the above preconditions 1–6 for Reasoning 3, which are shown in **Table S5**.

Table S5. Propositional symbols of preconditions for Reasoning 3.

Symbol	Meaning
k	$LF_{TSfBC,1}$ is given.
l	$\tilde{t}_{TSfBC,1}$ is given.
m	Pixels included in negative areas of $\tilde{t}_{TSfBC,1}$ ($\tilde{t}_{TSfBC,1}^n$) are most probably true tumour negatives.
n	t_{TSfBC} is given.
o	Pixels of t_{TSfBC} that are predicted as tumour positives (t_{TSfBC}^p) exist.
p	The intersection of pixels included in t_{TSfBC}^p and pixels included in $\tilde{t}_{TSfBC,1}^n$ can be considered as most probably predicted false tumour positives.
q	The intersection of pixels included in t_{TSfBC}^p and pixels included in $\tilde{t}_{TSfBC,1}^n$ can be considered as logically false positives.
r	The intersection of pixels of t_{TSfBC} that are predicted as tumour positives (t_{TSfBC}^p) and pixels included in negative areas of $\tilde{t}_{TSfBC,1}$ ($\tilde{t}_{TSfBC,1}^n$) can be considered as logically false positives.

Thirdly, referring to **Table S5**, we signify the propositional formalizations of the preconditions 1–6 for Reasoning 3 and Reasoning 3 via the propositional connectives listed in **Table S1** as follows.

- | | |
|---------------------------------|--------------|
| 1) $k \rightarrow l$ | Precondition |
| 2) $l \rightarrow m$ | Precondition |
| 3) $n \rightarrow o$ | Precondition |
| 4) $(m \wedge o) \rightarrow p$ | Precondition |
| 5) $p \rightarrow q$ | Precondition |
| 6) $q \rightarrow r$ | Precondition |
| $(n \wedge k) \rightarrow r$ | Reasoning 3 |

Fourthly, we show the validity of Reasoning 3 via the rules for proof of propositional logical reasoning listed in **Table S2** as follows.

- | | |
|---|---------------------------|
| $\therefore (n \wedge k) \rightarrow r$ | |
| 7) $n \wedge k$ | Hypothesis |
| 8) n | 7); $\wedge -$ |
| 9) k | 7); $\wedge -$ |
| 10) l | 1), 9); MP |
| 11) m | 2), 10); MP |
| 12) o | 3), 8); MP |
| 13) $m \wedge o$ | 11), 12); $\wedge +$ |
| 14) $(m \wedge o) \rightarrow q$ | 4), 5); HS |
| 15) $(m \wedge o) \rightarrow r$ | 14), 6); HS |
| 16) r | 15), 13); MP |
| 17) $(n \wedge k) \rightarrow r$ | 7)–16); Conditional Proof |

Since the hypothesis $n \wedge k$ of the 7) step has been fulfilled by the prediction of the image semantic segmentation

model for tumour segmentation for breast cancer (t_{TSfBC}) in section 5.2.3. and the two narrated logical facts $LF_{TSfBC} = \{LF_{TSfBC,1}, LF_{TSfBC,2}\}$, Reasoning 3 is proved to be valid. \square

Proof of Reasoning 4

Reasoning 4. *If t_{TSfBC} is given and $LF_{TSfBC,2}$ is given, then the intersection of pixels of t_{TSfBC} that are predicted as tumour positives (t_{TSfBC}^p) and pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ ($\tilde{t}_{TSfBC,2}^p$) can be considered as logically true positives, and the intersection of pixels of t_{TSfBC} that are predicted as tumour negatives (t_{TSfBC}^n) and pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ ($\tilde{t}_{TSfBC,2}^p$) can be considered as logically false negatives.*

Proof. Firstly, with the given t_{TSfBC} and $LF_{TSfBC,2}$, we have following preconditions for Reasoning 4.

- 1) If $LF_{TSfBC,2}$ is given, then $\tilde{t}_{TSfBC,2}$ is given.
- 2) If $\tilde{t}_{TSfBC,2}$ is given, then pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ ($\tilde{t}_{TSfBC,2}^p$) are most probably true tumour positives. (Reasoning 2).
- 3) If t_{TSfBC} is given, then pixels of t_{TSfBC} that are predicted as tumour positives (t_{TSfBC}^p) exist and pixels of t_{TSfBC} that are predicted as tumour negatives (t_{TSfBC}^n) exist.
- 4) If pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ ($\tilde{t}_{TSfBC,2}^p$) are most probably true tumour positives and pixels of t_{TSfBC} that are predicted as tumour positives (t_{TSfBC}^p) exist, then the intersection of pixels included in t_{TSfBC}^p and pixels included in $\tilde{t}_{TSfBC,2}^p$ can be considered as most probably predicted true tumour positives.
- 5) If pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ ($\tilde{t}_{TSfBC,2}^p$) are most probably true tumour positives and pixels of t_{TSfBC} that are predicted as tumour negatives (t_{TSfBC}^n) exist, then the intersection of pixels included in t_{TSfBC}^n and pixels included in $\tilde{t}_{TSfBC,2}^p$ can be considered as most probably predicted false tumour negatives.
- 6) If the intersection of pixels included in t_{TSfBC}^p and pixels included in $\tilde{t}_{TSfBC,2}^p$ can be considered as most probably predicted true tumour positives, then the intersection of pixels included in t_{TSfBC}^p and pixels included in $\tilde{t}_{TSfBC,2}^p$ can be considered as logically true positives.
- 7) If the intersection of pixels included in t_{TSfBC}^n and pixels included in $\tilde{t}_{TSfBC,2}^p$ can be considered as most probably predicted false tumour negatives, then the intersection of pixels included in t_{TSfBC}^n and pixels included in $\tilde{t}_{TSfBC,2}^p$ can be considered as logically false negatives.
- 8) If the intersection of pixels included in t_{TSfBC}^p and pixels included in $\tilde{t}_{TSfBC,2}^p$ can be considered as logically true positives, then the intersection of pixels of t_{TSfBC} that are predicted as tumour positives (t_{TSfBC}^p) and pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ ($\tilde{t}_{TSfBC,2}^p$) can be considered as logically true positives.
- 9) If the intersection of pixels included in t_{TSfBC}^n and pixels included in $\tilde{t}_{TSfBC,2}^p$ can be considered as logically false negatives, then the intersection of pixels of t_{TSfBC} that are predicted as tumour negatives (t_{TSfBC}^n) and pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ ($\tilde{t}_{TSfBC,2}^p$) can be considered as logically false negatives.

Secondly, we give the propositional symbols for the above preconditions 1-9 for Reasoning 4, which are shown in **Table S6**.

Table S6. Propositional symbols of preconditions for Reasoning 4.

Symbol	Meaning
s	$LF_{TSfBC,2}$ is given.
t	$\tilde{t}_{TSfBC,2}$ is given.
u	Pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ ($\tilde{t}_{TSfBC,2}^p$) are most probably true tumour positives.
v	t_{TSfBC} is given.
w	Pixels of t_{TSfBC} that are predicted as tumour positives (t_{TSfBC}^p) exist.
x	Pixels of t_{TSfBC} that are predicted as tumour negatives (t_{TSfBC}^n) exist.
y	The intersection of pixels included in t_{TSfBC}^p and pixels included in $\tilde{t}_{TSfBC,2}^p$ can be considered as most probably predicted true tumour positives.
z	The intersection of pixels included in t_{TSfBC}^n and pixels included in $\tilde{t}_{TSfBC,2}^p$ can be considered as most probably predicted false tumour negatives.
a	The intersection of pixels included in t_{TSfBC}^p and pixels included in $\tilde{t}_{TSfBC,2}^p$ can be considered as logically true positives.
b	The intersection of pixels included in t_{TSfBC}^n and pixels included in $\tilde{t}_{TSfBC,2}^p$ can be considered as logically false negatives.
c	The intersection of pixels of t_{TSfBC} that are predicted as tumour positives (t_{TSfBC}^p) and pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ ($\tilde{t}_{TSfBC,2}^p$) can be considered as logically true positives.
d	The intersection of pixels of t_{TSfBC} that are predicted as tumour negatives (t_{TSfBC}^n) and pixels included in positive areas of $\tilde{t}_{TSfBC,2}$ ($\tilde{t}_{TSfBC,2}^p$) can be considered as logically false negatives.

Thirdly, referring to **Table S6**, we signify the propositional formalizations of the preconditions 1–9 for Reasoning 4 and Reasoning 4 via the propositional connectives listed in **Table S1** as follows.

- | | | |
|----|---|--------------|
| 1) | $s \rightarrow t$ | Precondition |
| 2) | $t \rightarrow u$ | Precondition |
| 3) | $v \rightarrow (w \wedge x)$ | Precondition |
| 4) | $(u \wedge w) \rightarrow y$ | Precondition |
| 5) | $(u \wedge x) \rightarrow z$ | Precondition |
| 6) | $y \rightarrow a$ | Precondition |
| 7) | $z \rightarrow b$ | Precondition |
| 8) | $a \rightarrow c$ | Precondition |
| 9) | $b \rightarrow d$ | Precondition |
| | $(v \wedge s) \rightarrow (c \wedge d)$ | Reasoning 4 |

Fourthly, we show the validity of Reasoning 4 via the rules for proof of propositional logical reasoning listed in **Table S2** as follows.

$$\therefore (v \wedge s) \rightarrow (c \wedge d)$$

- | | | |
|-----|-------------------|-----------------|
| 10) | $v \wedge s$ | Hypothesis |
| 11) | v | 10); $\wedge -$ |
| 12) | s | 10); $\wedge -$ |
| 13) | $s \rightarrow u$ | 1), 2); HS |
| 14) | u | 13), 12); MP |
| 15) | $w \wedge x$ | 3), 11); MP |
| 16) | w | 15); $\wedge -$ |

17) x	15); $\wedge -$
18) $u \wedge w$	14), 16); $\wedge +$
19) $(u \wedge w) \rightarrow a$	4), 6); HS
20) $u \wedge x$	14), 17); $\wedge +$
21) $(u \wedge x) \rightarrow b$	5), 7); HS
22) $(u \wedge w) \rightarrow c$	19), 8); HS
23) $(u \wedge x) \rightarrow d$	21), 9); HS
24) c	22), 18); MP
25) d	23), 20); MP
26) $c \wedge d$	24), 25); $\wedge +$
27) $(v \wedge s) \rightarrow (c \wedge d)$	10)–26); Conditional Proof

Since the hypothesis $v \wedge s$ of the 10) step has been fulfilled by the prediction of the image semantic segmentation model for tumour segmentation for breast cancer (t_{TSfBC}) in section 5.2.3. and the two narrated logical facts $LF_{TSfBC} = \{LF_{TSfBC,1}, LF_{TSfBC,2}\}$ Reasoning 4 is proved to be valid.

Article

Innovation dynamics in BRICS economies investigated by artificial intelligence (AI)

Claudio Zancan^{1,*}, João Luiz Passador², Cláudia Souza Passador², Ricardo Carvalho Rodrigues¹

¹ Academia de Propriedade Intelectual do Instituto Nacional da Propriedade Industrial – INPI, Rio de Janeiro 20090-910, Brazil

² Faculdade de Economia, Administração e Contabilidade de Ribeirão Preto – FEA-RP da Universidade de São Paulo – USP, Ribeirão Preto 14040-95, Brazil

* Corresponding author: Claudio Zancan, claudiozancan@gmail.com

CITATION

Zancan C, Passador JL, Passador CS, Rodrigues RC. Innovation dynamics in BRICS economies investigated by artificial intelligence (AI). *Computing and Artificial Intelligence*. 2024; 2(2): 1291. <https://doi.org/10.59400/cai.v2i2.1291>

ARTICLE INFO

Received: 15 April 2024

Accepted: 27 June 2024

Available online: 11 July 2024

COPYRIGHT



Copyright © 2024 by author(s). *Computing and Artificial Intelligence* is published by Academic Publishing Pte. Ltd. This work is licensed under the Creative Commons Attribution (CC BY) license. <https://creativecommons.org/licenses/by/4.0/>

Abstract: This study aims to address the existing knowledge gap regarding the specific impact of artificial intelligence (AI) on patent research and emphasize its strategic significance as a catalyst for innovation. The methodology employs a comprehensive approach, integrating both qualitative and quantitative research methods. It systematically investigates the transformative potential of AI in patent research within the BRICS nations, including an examination of the technological, ethical, and legal challenges associated with AI's application in patent analysis. This research contributes to the field by extending beyond the conventional focus on the role of patents in innovation and shedding light on the potential of AI in patent research. It offers valuable insights into how AI can redefine the landscape of patent research, providing a more rapid and accurate perspective on the identification of technological trends, opportunities, and competitive factors. The findings underscore that AI in patent research yields numerous advantages, ranging from efficient data processing to the forecasting of technological trends. Future studies should explore ethical and legal considerations associated with AI in patent research, as well as its implementation in the strategies of both corporate entities and governmental bodies in the BRICS nations.

Keywords: artificial intelligence; patent research; BRICS countries; innovation

1. Introduction

In the rapidly evolving technological landscape, characterized by the relentless pace of innovation, patent research serves as an essential tool for understanding the underlying dynamics and identifying emerging trends. It plays a pivotal role in demystifying the innovation process and providing insights into how ideas materialize into strategic advancements for organizations.

In the BRICS nations (Brazil, Russia, India, China, and South Africa), On 24 August 2023, the leaders of the BRICS countries agreed to expand the group to include six new members: Saudi Arabia, Argentina, Egypt, the United Arab Emirates, Ethiopia, and Iran. The expansion of the BRICS is viewed as a step to enhance the group's global influence and foster cooperation among emerging economies. For the purposes of this text, we will only consider the original signatory countries of the BRICS group. The expansion of the group to encompass the six new members is anticipated to take place in the coming year, pending the completion of negotiations between the existing members and the prospective new entrants [1], distinguished by their remarkable economic ascent and growing global influence, innovation stands as a central pillar of sustainable development. Patent research can play a vital role in this context, offering a panoramic view of technological progress and guiding

research and development endeavors [2].

Artificial intelligence (AI) is a spectrum of technologies with astonishing potential to revolutionize patent research. AI acts as a catalyst, expediting the discovery and comprehension of information within vast patent databases. It streamlines the search process and unveils intricate patterns that may elude traditional analyses.

Despite considerable efforts devoted to comprehending the role of patents in fostering innovation, a significant gap exists in understanding the specific impact of artificial intelligence in this context. While some studies have initiated exploration into the significance of patent analysis in BRICS countries, a systematic investigation of AI as a strategic tool for patent research necessitates a more profound and comprehensive analysis.

This research tackles a critical gap in understanding by exploring the intricate interplay between artificial intelligence (AI) and patent research within the context of BRICS countries. By examining the technological specifics, legal nuances, and ethical challenges surrounding AI's application in patent analysis, this study will unveil novel perspectives on how AI can revolutionize the exploration and utilization of patents to fuel innovation. Moreover, it aims to make a substantial contribution to both academic knowledge and the practical realm of innovation and technological development. Highlighting AI as a strategic tool, this research seeks to inspire the conscious and informed adoption of AI in patent research, thereby empowering researchers, companies, and government entities to leverage the vast potential of patent data to drive their innovative initiatives forward.

In summary, this research aims to provide well-founded insights into how AI can redefine our approach to patent research, offering a faster, more accurate, and more effective perspective on identifying technological trends, detecting opportunities, and analyzing competition. It seeks to illustrate the transformative role of AI as an innovation catalyst, expediting the evolutionary trajectory of vital sectors and enhancing the competitiveness of BRICS countries in the global research and development landscape.

In the subsequent sections of this article, following this introductory text, we will show the theoretical underpinnings, with a primary emphasis on Patents as Tangible Indicators of Technological Trajectories and their role as reflections of innovation strategies. Additionally, we will scrutinize the significance of Patent Analysis in BRICS countries and the application of AI in Patent Analysis. The ensuing sections will elaborate on the methodology employed, present the primary findings, and conclude our discussion. Finally, we will provide a comprehensive list of references from the sources used in this work.

2. Theoretical background

The analysis of the role of patents in innovation and development is essential for understanding economic and technological dynamics, particularly in BRICS countries. Patents are not simply legal documents that grant inventors exclusive rights to exploit their inventions; they also serve as tangible indicators of technological trajectories and innovation strategies adopted by companies and

nations. In this section of the literature review, we will explore the relevance of patents as valuable tools for driving economic growth and technological development in BRICS countries.

2.1. Patents as tangible indicators of technological trajectories

The use of patents as tangible indicators of technological trajectories has demonstrated its paramount significance in various contexts, with a particular emphasis on BRICS countries, which are committed to bolstering their technological development and industrialization. This discourse aims to contribute to the ongoing dialogue surrounding the utilization of patents as instruments for assessing technological progress and identifying emerging trends.

Patents, as legal instruments bestowed by governments upon inventors, serve the fundamental purpose of comprehensively documenting technological innovations. They detect the existence of an invention and provide precise technical specifications elucidating the inner workings of the technology, its practical applications, and its distinctions from preexisting solutions. This is aptly exemplified in studies such as the work of Hall et al. [3], which underscores how patent analysis yields valuable insights into the technical attributes and innovations underpinning transformative technologies.

The monitoring of technological progress plays a pivotal role in the formulation of public policies and the allocation of investments in Research and Development (R&D). Patents, owing to their lucid and quantifiable nature, stand as a noteworthy barometer of such progress. Research endeavors such as that of Lanjouw and Schankerman [4] elucidate how patent analysis can be leveraged to gauge the quality of innovations and their correlation with research and technological productivity.

Patents also serve as valuable tools for discerning emerging trends within specific technological domains. An analysis of recent patent categories enables the identification of areas experiencing substantial advancements. The study conducted by Hu and his colleagues accentuates how patent analysis can be deployed to prognosticate forthcoming trends and guide decisions regarding innovation investments [5].

In the context of BRICS countries, patent analysis plays even greater significance, given their initiative-taking pursuit of technological development and industrialization. Research underscores how countries like Brazil have harnessed patent analysis to bolster their research and development sectors [6].

Consequently, patents have an instrumental role in BRICS countries by providing a robust foundation for assessing technological progress, discerning emerging trends, and shaping innovation policies. The studies cited, along with other scholarly contributions, underscore the centrality of patents as tangible indicators of technological trajectories, highlighting their utility in guiding strategies for technological and economic development.

2.2. Patents as reflections of innovation strategies

In the ever-changing landscape of innovation, patents serve as both milestones

of technological progress and revealing windows into the intricate strategies adopted by organizations and institutions. This section explores the multifaceted role of patents and patent analysis in guiding innovation strategies and fostering a culture of innovation, particularly within the competitive framework of BRICS countries.

Patent analysis assumes a pivotal role in discerning an organization's strategic focus across diverse realms of research and innovation, spanning from incremental advancements to disruptive breakthroughs and beyond. Such discernment holds paramount importance in shaping efficacious innovation strategies, whether at the organizational or national level, particularly within the fiercely competitive global landscape characterizing BRICS nations.

Li and his colleagues [7] encompassed several critical aspects about patent analysis. For this author, firstly, it allows for the identification and comprehension of the competitive landscape. Through the examination of competitors patent portfolios, organizations can discern their strengths, weaknesses, and future technological trajectories. Such insights are invaluable for formulating strategic approaches to maintain a competitive edge in the market. Additionally, patent analysis facilitates the recognition of emerging trends and technologies within a specific domain. By discerning modern innovations gaining prominence, organizations can make informed decisions regarding research and development (R&D) investments, ensuring alignment with evolving market demands and technological advancements.

Patents analysis also provide a clear and measurable gauge of technological progress. Emphasizing the significance of patent quality as an indicator of this progress, Lanjouw and Schankerman develop metrics based on criteria such as the extent of patent claims and citation frequency [4]. This enables a more comprehensive assessment of innovation. Also, high-quality patents tend to have a more significant impact on research and technological productivity. They are cited more frequently in academic papers and adopted by companies. Another work [8] explores the correlation between patent citations and the value of innovations, providing further evidence of patents as metrics for measuring their impact on research and development.

Assessing the quality of innovations through patent analysis has substantial implications for public policies and resource allocation in R&D. Governments, development agencies, and companies can strategically allocate resources by identifying areas with high innovation potential. In the context of BRICS countries, patent analysis is a valuable tool for guiding innovation strategies and fostering a culture of innovation. By leveraging this tool, BRICS countries can accelerate their technological advancement and enhance their competitiveness in the global market. Overall, patent analysis is a powerful tool that can be used to guide innovation strategies and foster a culture of innovation in BRICS countries.

2.3. Importance of patent analysis in BRICS nations

Innovation is widely recognized as a key driver of sustainable economic growth in BRICS countries. In this context, patent analysis assumes paramount significance, serving as a strategic tool for shaping innovation policies, propelling technological development, and fostering international collaboration. Investing in refined patent

analysis equips these nations with the capability to discern strategic focal points, identify avenues for international collaboration, and progress towards the achievement of their development objectives.

The significance of patent analysis in BRICS countries is evident in its role in supporting the formulation of effective innovation policies, which are pivotal in facilitating economic growth and enhancing global competitiveness. As evidenced by the study conducted by Dutrénit et al. [9], patent analysis provides a robust foundation for identifying strategic imperatives. This study highlights the integral role of patent analysis in China, where it informs innovation policies and guides resource allocation towards strategic sectors.

Sustainable economic growth holds paramount importance for the BRICS nations, due to their substantial population and economic influence. Central to this growth is the pivotal role played by technological advancement, as it underpins productivity enhancement, innovation, and global competitiveness. These nations are cognizant that traditional industries and resources alone are insufficient in the contemporary globalized economy, necessitating substantial investments in technology to thrive. International collaboration emerges as an essential catalyst in nurturing innovation and technological development. Patent analysis, as elucidated in this research [10], empowers BRICS countries to identify common ground with other nations and foster strategic partnerships.

For a nation exemplifying the BRICS framework, such as India, patent analysis assumes even greater significance. India's burgeoning innovation landscape is brought into focus through comprehensive patent analysis, as illustrated in this study [11]. It provides vital insights into the nation's technological progress, identifies areas of excellence like information technology and pharmaceuticals, and spotlights sectors that necessitate augmented investments and strategic emphasis. These findings from patent analysis offer policymakers a compass to navigate the complex terrain of technological development, guiding them in formulating policies, incentives, and strategies essential for sustaining innovation-driven growth and competitiveness in the evolving global economy.

In addition to the studies mentioned, recent research further underscores the significance of patent analysis in BRICS countries. Some authors [12] conducted a comparative analysis of innovation policies in BRICS countries, shedding light on the role of patent analysis in shaping these policies. Another authors [13] investigated the intricate relationship between technological development, economic growth, and patent analysis in BRICS nations.

Patent analysis stands as a strategic cornerstone in BRICS countries, orchestrating the harmonious symphony of innovation policies, technological progress, international collaboration, and the realization of development aspirations. The collective body of research, including the studies cited above, underscores the versatility and indispensability of patent analysis as an instrument propelling these nations towards a more innovative and prosperous future.

2.4. Patent analysis using artificial intelligence (AI)

In the contemporary era, patent analysis has undergone significant evolution

thanks to the application of artificial intelligence (AI). AI enables the extraction of valuable insights from large volumes of patent documents, enabling comprehensive analysis of technological trends, innovation strategies, and patent quality. This has revolutionized patent analysis, empowering organizations and research institutions to make more informed decisions, identify strategic opportunities, and accelerate scientific and technological progress [2].

An array of advanced AI-based tools is now widely available for patent analysis, significantly simplifying the process of searching and analyzing patent information [14]. These tools employ innovative AI algorithms to process a vast amount of patent documents on a global scale, streamlining access to relevant information and enabling deeper and more comprehensive analysis.

AI also plays a fundamental role in identifying patterns and trends in patent analysis, offering abilities that would be unattainable manually. By examining large patent data sets, AI reveals emerging areas of innovation, predicts technological directions, and identifies convergences among research fields. This powerful set of capabilities has a profound impact in several key areas. AI guides technological research and development by supplying valuable insights to organizations and academic institutions. For example, AI can be used to predict emerging technological trends, guiding R&D efforts to areas likely to have a significant further impact.

AI-based tools are also being applied in specific sectors like the pharmaceutical industry to analyze patents related to drugs and therapies, identifying trends in new drug research. This accelerates the development of innovative therapies and identify collaboration opportunities among pharmaceutical companies, research institutions, and regulatory bodies. AI also plays a crucial role in perfecting the formulation of innovation policies. It empowers governments and development agencies to distribute resources precisely and strategically based on insights generated by AI, resulting in more effective support for innovation areas considered strategic for economic growth and global competitiveness [15].

While AI offers significant advantages in patent analysis, it is essential to recognize and address the ethical challenges this technology presents. One of the primary ethical challenges is the protection of the privacy of inventors and patent holders. The use of AI may involve processing sensitive information contained in patent documents, raising concerns about unauthorized disclosure or misuse of this data [12,16]. Additionally, data security is also a central concern, as AI relies on access to substantial amounts of information.

The integration of AI into patent analysis is a significant advancement, driving innovation and technological development worldwide. By leveraging AI-driven tools and techniques, organizations and institutions can extract valuable insights from patent data, guide R&D efforts, accelerate innovation in specific sectors, and shape effective innovation policies [15]. However, it is essential to recognize and address the ethical challenges associated with AI, such as privacy and security concerns. Then, in the next section, we explore the methodological stages employed to gain profound insights into the patent landscape within BRICS countries.

3. Methodological stages

In pursuit of gaining insights into the patent landscape within BRICS countries, this article employed a method grounded in AI techniques, with the programming of algorithms capable of analyzing the context of patents as drivers of economic growth in the countries. The following is a summarized description of the methodological in four stages used.

The first stage involved Data Collection and Preparation. To set up the foundation of our study, we began by gathering comprehensive patent data from reliable sources in each BRICS country. This included patent offices, academic repositories, and industry-specific databases. Carefully, we cleaned and standardized this data to ensure its accuracy and consistency.

The second stage involved the extraction of patent features, recognizing them as rich sources of information, from inventors' names to keywords, classifications, and citation networks. To effectively use this information, advanced feature extraction techniques were employed. Natural Language Processing (NLP) methods were used to find and extract keywords, trends, and emerging technologies from patent documents. Furthermore, we analyzed citation patterns to understand the influence and relevance of each patent [17]. Aligned with these authors we employed several recognized metrics to assess patent influence and relevance based on citation behavior. These include this algorithmics:

- Citation count: The total number of times a patent is cited by others, serving as a direct indicator of its influence and dissemination within the technological field. The **Algorithm 1** is showed below.

Algorithm 1 Citation count

```

1: import re
2:
3: # Provided database
4: data = ""
5: Patents sources analyzed.
6: Country Information Patents Source
7: Brazil Instituto Nacional da Propriedade Industrial (INPI): The official website of INPI, the Brazilian government
   agency responsible for patent registration.
8: Home Page: https://www.gov.br/inpi/pt-br/ (Access in 10-12-23).
9:
10: India Patent Information System (PIS): The portal of the Department of Intellectual Property of India that provides
   access to Indian patent data and Intellectual Property.
11: Home Page: https://ipindia.gov.in/patent-information-system.htm (Access in 10-12-23).
12:
13: China China National Intellectual Property Administration (CNIPA): The official website of CNIPA, the Chinese
   government agency responsible for patent registration.
14: Home Page: https://english.cnipa.gov.cn/ (Access in 10-12-23).
15:
16: Russia Federal Service for Intellectual Property, Patents, and Trademarks (Rospatent): The official website of
   Rospatent, the Russian government agency responsible for patent registration.
17: Home Page: https://rospatent.gov.ru/en (Access in 10-12-23).
18:
19: South Africa Patent Office of the Republic of South Africa (CIPC): The official website of CIPC, the South African
   government agency responsible for patent registration.
20: Home Page: https://www.cipc.co.za/?page\_id=1423 (Access in 10-12-23).
21:

```

Algorithm 1 (Continued)

```

22: For all countries World Intellectual Property Organization (WIPO): The WIPO website provides access to international
    patent data.
23: Home Page: https://www.wipo.int/patentscope/en/ (Access in 10-12-23).
24: """"
25:
26: # Function to extract information
27: def extract_patent_info(data):
28:     pattern = re.compile(r'([A-Za-z\s]+)\s+([\^:~\+])\s+([\^n]+)\nHome Page:\s+([\^s]+)\s+(\d+([\d-]+))')
29:     matches = pattern.findall(data)
30:     patent_info = []
31:
32:     for match in matches:
33:         country, source, description, homepage, access_date = match
34:         patent_info.append({
35:             "Country": country.strip(),
36:             "Source": source.strip(),
37:             "Description": description.strip(),
38:             "Home Page": homepage.strip(),
39:             "Access Date": access_date.strip()
40:         })
41:
42:     return patent_info
43:
44: # Extract patent information
45: patent_info = extract_patent_info(data)
46:
47: # Function to extract citation count (fictitious example of citation count)
48: def get_citation_count(description):
49:     # Fictitious example of citation count
50:     return "Citation count: The total number of times a patent is cited by others, serving as a direct indicator of its influence
    and dissemination within the technological field."
51:
52: # Display extracted information and citation count
53: for info in patent_info:
54:     print(f"Country: {info['Country']}")
55:     print(f"Source: {info['Source']}")
56:     print(f"Description: {info['Description']}")
57:     print(f"Home Page: {info['Home Page']}")
58:     print(f"Access Date: {info['Access Date']}")
59:     print(f"{get_citation_count(info['Description'])}")
60:     print()
61:
62: # Expected output (citation count would be obtained specifically for each patent)

```

- Time-adjusted citations: Considering that more recent patents are more likely to accumulate citations, we normalized the citation count by the time elapsed since the patent publication. The **Algorithm 2** is showed below.

Algorithm 2 Time-adjusted citations

```

1: import re
2: from datetime import datetime
3:
4: # Provided database with additional hypothetical publication dates and citation counts
5: data = """"
6: Patents sources analyzed.
7: Country Information Patents Source

```

Algorithm 2 (Continued)

```

8:  Brazil  Instituto Nacional da Propriedade Industrial (INPI): The official website of INPI, the Brazilian government
      agency responsible for patent registration.
9:  Home Page: https://www.gov.br/inpi/pt-br/ (Access in 10-12-23). Publication Date: 2020-01-15, Citation Count: 50
10:
11: India  Patent Information System (PIS): The portal of the Department of Intellectual Property of India that provides
      access to Indian patent data and Intellectual Property.
12: Home Page: https://ipindia.gov.in/patent-information-system.htm (Access in 10-12-23). Publication Date: 2019-06-20,
      Citation Count: 80
13:
14: China  China National Intellectual Property Administration (CNIPA): The official website of CNIPA, the Chinese
      government agency responsible for patent registration.
15: Home Page: https://english.cnipa.gov.cn/ (Access in 10-12-23). Publication Date: 2021-04-10, Citation Count: 30
16:
17: Russia  Federal Service for Intellectual Property, Patents, and Trademarks (Rospatent): The official website of
      Rospatent, the Russian government agency responsible for patent registration.
18: Home Page: https://rospatent.gov.ru/en (Access in 10-12-23). Publication Date: 2018-12-05, Citation Count: 100
19:
20: South Africa  Patent Office of the Republic of South Africa (CIPC): The official website of CIPC, the South African
      government agency responsible for patent registration.
21: Home Page: https://www.cipc.co.za/?page\_id=1423 (Access in 10-12-23). Publication Date: 2022-08-25, Citation Count:
      20
22:
23: For all countries  World Intellectual Property Organization (WIPO): The WIPO website provides access to international
      patent data.
24: Home Page: https://www.wipo.int/patentscope/en/ (Access in 10-12-23). Publication Date: 2017-03-30, Citation Count:
      150
25: """"
26:
27: # Function to extract information
28: def extract_patent_info(data):
29:     pattern = re.compile(r'([A-Za-z\s]+)\s+([\^:~\+])\s+([\^n]+)\nHome Page:\s+([\^s]+)\s+\(Access in ([\d-]+)\). Publication
      Date:\s+([\d-]+), Citation Count:\s+(\d+)')
30:     matches = pattern.findall(data)
31:     patent_info = []
32:
33:     for match in matches:
34:         country, source, description, homepage, access_date, publication_date, citation_count = match
35:         patent_info.append({
36:             "Country": country.strip(),
37:             "Source": source.strip(),
38:             "Description": description.strip(),
39:             "Home Page": homepage.strip(),
40:             "Access Date": access_date.strip(),
41:             "Publication Date": publication_date.strip(),
42:             "Citation Count": int(citation_count.strip())
43:         })
44:
45:     return patent_info
46:
47: # Function to calculate time-adjusted citations
48: def time_adjusted_citations(publication_date, citation_count):
49:     current_date = datetime.strptime("2023-10-12", "%Y-%m-%d")
50:     publication_date = datetime.strptime(publication_date, "%Y-%m-%d")
51:     time_elapsed = (current_date - publication_date).days / 365.25 # Convert days to years
52:     if time_elapsed == 0:
53:         time_elapsed = 1 # Avoid division by zero
54:     adjusted_citations = citation_count / time_elapsed
55:     return adjusted_citations

```

Algorithm 2 (Continued)

```

56: # Extract patent information
57: patent_info = extract_patent_info(data)
58:
59: # Display extracted information and time-adjusted citation counts
60: for info in patent_info:
61:     adjusted_citations = time_adjusted_citations(info['Publication Date'], info['Citation Count'])
62:     print(f"Country: {info['Country']}")
63:     print(f"Source: {info['Source']}")
64:     print(f"Description: {info['Description']}")
65:     print(f"Home Page: {info['Home Page']}")
66:     print(f"Access Date: {info['Access Date']}")
67:     print(f"Publication Date: {info['Publication Date']}")
68:     print(f"Citation Count: {info['Citation Count']}")
69:     print(f"Time-Adjusted Citations: {adjusted_citations:.2f}")
70:     print()
71:
72: # Expected output: Adjusted citation counts normalized by time elapsed since publication

```

- Centrality in the citation network: We analyzed the position of each patent within the citation network, mapping its importance as a bridge between different technological areas or as a central element of a thematic cluster. The **Algorithm 3** is showed below.

Algorithm 3 Centrality in the citation network

```

1: import re
2: import networkx as nx
3: from datetime import datetime
4:
5: # Provided database with additional hypothetical publication dates and citation counts
6: data = ""
7: Patents sources analyzed.
8: Country Information Patents Source
9: Brazil Instituto Nacional da Propriedade Industrial (INPI): The official website of INPI, the Brazilian government
   agency responsible for patent registration.
10: Home Page: https://www.gov.br/inpi/pt-br / (Access in 10-12-23). Publication Date: 2020-01-15, Citation Count: 50
11:
12: India Patent Information System (PIS): The portal of the Department of Intellectual Property of India that provides
   access to Indian patent data and Intellectual Property.
13: Home Page: https://ipindia.gov.in/patent-information-system.htm (Access in 10-12-23). Publication Date: 2019-06-20,
   Citation Count: 80
14:
15: China China National Intellectual Property Administration (CNIPA): The official website of CNIPA, the Chinese
   government agency responsible for patent registration.
16: Home Page: https://english.cnipa.gov.cn/ (Access in 10-12-23). Publication Date: 2021-04-10, Citation Count: 30
17:
18: Russia Federal Service for Intellectual Property, Patents, and Trademarks (Rospatent): The official website of
   Rospatent, the Russian government agency responsible for patent registration.
19: Home Page: https://rospatent.gov.ru/en (Access in 10-12-23). Publication Date: 2018-12-05, Citation Count: 100
20:
21: South Africa Patent Office of the Republic of South Africa (CIPC): The official website of CIPC, the South African
   government agency responsible for patent registration.
22: Home Page: https://www.cipc.co.za/?page_id=1423 (Access in 10-12-23). Publication Date: 2022-08-25, Citation Count:
   20
23:
24: For all countries World Intellectual Property Organization (WIPO): The WIPO website provides access to international
   patent data.

```

Algorithm 3 (Continued)

```

25: Home Page: https://www.wipo.int/patentscope/en/ (Access in 10-12-23). Publication Date: 2017-03-30, Citation Count:
    150
26: ""
27:
28: # Function to extract information
29: def extract_patent_info(data):
30:     pattern = re.compile(r'([A-Za-z\s]+\s+([\^:]+\s+([\^n]+\n)Home Page:\s+([\^s]+\s+)\(Access in ([\d-]+\)). Publication
    Date:\s+([\d-]+\s+), Citation Count:\s+(\d+)')
31:     matches = pattern.findall(data)
32:     patent_info = []
33:
34:     for match in matches:
35:         country, source, description, homepage, access_date, publication_date, citation_count = match
36:         patent_info.append({
37:             "Country": country.strip(),
38:             "Source": source.strip(),
39:             "Description": description.strip(),
40:             "Home Page": homepage.strip(),
41:             "Access Date": access_date.strip(),
42:             "Publication Date": publication_date.strip(),
43:             "Citation Count": int(citation_count.strip())
44:         })
45:
46:     return patent_info
47:
48: # Hypothetical citation relationships between patents
49: citations = [
50:     ("Brazil", "India"),
51:     ("India", "China"),
52:     ("China", "Russia"),
53:     ("Russia", "South Africa"),
54:     ("South Africa", "WIPO"),
55:     ("WIPO", "Brazil"),
56:     ("India", "WIPO"),
57:     ("China", "Brazil"),
58: ]
59:
60: # Extract patent information
61: patent_info = extract_patent_info(data)
62:
63: # Create a directed graph
64: G = nx.DiGraph()
65:
66: # Add nodes with attributes
67: for info in patent_info:
68:     G.add_node(info['Country'], **info)
69:
70: # Add citation edges
71: G.add_edges_from(citations)
72:
73: # Compute centrality metrics
74: betweenness centrality = nx.betweenness centrality(G)
75: closeness centrality = nx.closeness centrality(G)
76: eigenvector centrality = nx.eigenvector centrality(G)
77:
78: # Display centrality metrics for each patent
79: for country in G.nodes():
80:     print(f"Country: {country}")

```

Algorithm 3 (Continued)

```

81: print(f"Betweenness Centrality: {betweenness_centrality[country]:.4f}")
82: print(f"Closeness Centrality: {closeness_centrality[country]:.4f}")
83: print(f"Eigenvector Centrality: {eigenvector_centrality[country]:.4f}")
84: print()
85: # Expected output: Centrality metrics for each patent

```

- Priority coefficient: Evaluates the initial influence of a patent by considering the first citations received within a specific period after its publication. The **Algorithm 4** is showed below.

Algorithm 4 Priority coefficient

```

1: import re
2: import networkx as nx
3: from datetime import datetime, timedelta
4:
5: # Provided database with additional hypothetical publication dates and citation counts
6: data = ""
7: Patents sources analyzed.
8: Country Information Patents Source
9: Brazil Instituto Nacional da Propriedade Industrial (INPI): The official website of INPI, the Brazilian government
  agency responsible for patent registration.
10: Home Page: https://www.gov.br/inpi/pt-br/ (Access in 10-12-23). Publication Date: 2020-01-15, Citation Count: 50
11:
12: India Patent Information System (PIS): The portal of the Department of Intellectual Property of India that provides
  access to Indian patent data and Intellectual Property.
13: Home Page: https://ipindia.gov.in/patent-information-system.htm (Access in 10-12-23). Publication Date: 2019-06-20,
  Citation Count: 80
14:
15: China China National Intellectual Property Administration (CNIPA): The official website of CNIPA, the Chinese
  government agency responsible for patent registration.
16: Home Page: https://english.cnipa.gov.cn/ (Access in 10-12-23). Publication Date: 2021-04-10, Citation Count: 30
17:
18: Russia Federal Service for Intellectual Property, Patents, and Trademarks (Rospatent): The official website of
  Rospatent, the Russian government agency responsible for patent registration.
19: Home Page: https://rospatent.gov.ru/en (Access in 10-12-23). Publication Date: 2018-12-05, Citation Count: 100
20:
21: South Africa Patent Office of the Republic of South Africa (CIPC): The official website of CIPC, the South African
  government agency responsible for patent registration.
22: Home Page: https://www.cipc.co.za/?page_id=1423 (Access in 10-12-23). Publication Date: 2022-08-25, Citation Count:
  20
23:
24: For all countries World Intellectual Property Organization (WIPO): The WIPO website provides access to international
  patent data.
25: Home Page: https://www.wipo.int/patentscope/en/ (Access in 10-12-23). Publication Date: 2017-03-30, Citation Count:
  150
26: ""
27:
28: # Hypothetical citation data with dates
29: citation_data = {
30:     "Brazil": [("India", "2020-03-01"), ("China", "2020-07-15"), ("WIPO", "2021-01-10")],
31:     "India": [("China", "2019-07-20"), ("WIPO", "2019-09-25"), ("Brazil", "2020-03-01")],
32:     "China": [("Russia", "2021-05-01"), ("Brazil", "2021-07-20"), ("South Africa", "2021-11-10")],
33:     "Russia": [("South Africa", "2019-01-10"), ("China", "2019-06-15"), ("WIPO", "2020-03-25")],
34:     "South Africa": [("Brazil", "2022-09-10"), ("Russia", "2023-01-05"), ("India", "2023-04-15")],
35:     "WIPO": [("Brazil", "2017-05-01"), ("India", "2017-08-10"), ("Russia", "2018-01-15")],
36: }
37:

```

Algorithm 4 (Continued)

```

38: # Function to extract information
39: def extract_patent_info(data):
40:     pattern = re.compile(r'([A-Za-z\s]+)\s+([\^:]+):\s+([\^\n]+)\nHome Page:\s+([\^\s]+)\s+\(Access in ([\d-]+)\). Publication
Date:\s+([\d-]+), Citation Count:\s+(\d+)')
41:     matches = pattern.findall(data)
42:     patent_info = []
43:
44:     for match in matches:
45:         country, source, description, homepage, access_date, publication_date, citation_count = match
46:         patent_info.append({
47:             "Country": country.strip(),
48:             "Source": source.strip(),
49:             "Description": description.strip(),
50:             "Home Page": homepage.strip(),
51:             "Access Date": access_date.strip(),
52:             "Publication Date": publication_date.strip(),
53:             "Citation Count": int(citation_count.strip())
54:         })
55:
56:     return patent_info
57:
58: # Function to calculate the priority coefficient
59: def priority_coefficient(publication_date, citations, period_days=365):
60:     publication_date = datetime.strptime(publication_date, "%Y-%m-%d")
61:     period_end_date = publication_date + timedelta(days=period_days)
62:
63:     initial_citations = [c for c in citations if publication_date <= datetime.strptime(c[1], "%Y-%m-%d") <=
period_end_date]
64:     return len(initial_citations)
65:
66: # Extract patent information
67: patent_info = extract_patent_info(data)
68:
69: # Display extracted information and priority coefficients
70: for info in patent_info:
71:     citations = citation_data.get(info['Country'], [])
72:     priority_coeff = priority_coefficient(info['Publication Date'], citations)
73:     print(f"Country: {info['Country']}")
74:     print(f"Source: {info['Source']}")
75:     print(f"Description: {info['Description']}")
76:     print(f"Home Page: {info['Home Page']}")
77:     print(f"Access Date: {info['Access Date']}")
78:     print(f"Publication Date: {info['Publication Date']}")
79:     print(f"Citation Count: {info['Citation Count']}")
80:     print(f"Priority Coefficient: {priority_coeff}")
81:     print()
82: # Expected output: Priority coefficient for each patent based on initial citations within a specific period

```

Additionally, we investigated potential citation biases, such as self-citations or reciprocal citations, to ensure the robustness and reliability of our measurements. Through this rigorous analysis, grounded in recognized metrics and attentive to potential biases, we obtained a granular understanding of the influence and relevance of each patent in our study, complementing the insights of Zhang and his colleagues [17] and offering new indicators for the assessment of patent-based innovation.

The Python **Algorithm 5** code is shown in the following.

Algorithm 5 Potential citation biases

```

1: import re
2: import networkx as nx
3: from datetime import datetime, timedelta
4:
5: # Provided database with additional hypothetical publication dates and citation counts
6: data = ""
7: Patents sources analyzed.
8: Country Information Patents Source
9: Brazil Instituto Nacional da Propriedade Industrial (INPI): The official website of INPI, the Brazilian government
agency responsible for patent registration.
10: Home Page: https://www.gov.br/inpi/pt-br/ (Access in 10-12-23). Publication Date: 2020-01-15, Citation Count: 50
11:
12: India Patent Information System (PIS): The portal of the Department of Intellectual Property of India that provides
access to Indian patent data and Intellectual Property.
13: Home Page: https://ipindia.gov.in/patent-information-system.htm (Access in 10-12-23). Publication Date: 2019-06-20,
Citation Count: 80
14:
15: China China National Intellectual Property Administration (CNIPA): The official website of CNIPA, the Chinese
government agency responsible for patent registration.
16: Home Page: https://english.cnipa.gov.cn/ (Access in 10-12-23). Publication Date: 2021-04-10, Citation Count: 30
17:
18: Russia Federal Service for Intellectual Property, Patents, and Trademarks (Rospatent): The official website of Rospatent,
the Russian government agency responsible for patent registration.
19: Home Page: https://rospatent.gov.ru/en (Access in 10-12-23). Publication Date: 2018-12-05, Citation Count: 100
20:
21: South Africa Patent Office of the Republic of South Africa (CIPC): The official website of CIPC, the South African
government agency responsible for patent registration.
22: Home Page: https://www.cipc.co.za/?page\_id=1423 (Access in 10-12-23). Publication Date: 2022-08-25, Citation Count:
20
23:
24: For all countries World Intellectual Property Organization (WIPO): The WIPO website provides access to international
patent data.
25: Home Page: https://www.wipo.int/patentscope/en/ (Access in 10-12-23). Publication Date: 2017-03-30, Citation Count:
150
26: ""
27:
28: # Hypothetical citation data with dates
29: citation_data = {
30:     "Brazil": [("India", "2020-03-01"), ("China", "2020-07-15"), ("WIPO", "2021-01-10")],
31:     "India": [("China", "2019-07-20"), ("WIPO", "2019-09-25"), ("Brazil", "2020-03-01")],
32:     "China": [("Russia", "2021-05-01"), ("Brazil", "2021-07-20"), ("South Africa", "2021-11-10")],
33:     "Russia": [("South Africa", "2019-01-10"), ("China", "2019-06-15"), ("WIPO", "2020-03-25")],
34:     "South Africa": [("Brazil", "2022-09-10"), ("Russia", "2023-01-05"), ("India", "2023-04-15")],
35:     "WIPO": [("Brazil", "2017-05-01"), ("India", "2017-08-10"), ("Russia", "2018-01-15")],
36: }
37:
38: # Function to extract information
39: def extract_patent_info(data):
40:     pattern = re.compile(r'([A-Za-z\s]+)\s+([\^:]+):\s+([\^n]+)\nHome Page:\s+([\^s]+)\s+\(Access in ([\d-]+)\). Publication
Date:\s+([\d-]+), Citation Count:\s+(\d+)')
41:     matches = pattern.findall(data)
42:     patent_info = []
43:
44:     for match in matches:
45:         country, source, description, homepage, access_date, publication_date, citation_count = match
46:         patent_info.append({
47:             "Country": country.strip(),
48:             "Source": source.strip(),

```

Algorithm 5 (Continued)

```

49:     "Description": description.strip(),
50:     "Home Page": homepage.strip(),
51:     "Access Date": access_date.strip(),
52:     "Publication Date": publication_date.strip(),
53:     "Citation Count": int(citation_count.strip())
54: })
55:
56: return patent_info
57:
58: # Function to calculate the priority coefficient
59: def priority_coefficient(publication_date, citations, period_days=365):
60:     publication_date = datetime.strptime(publication_date, "%Y-%m-%d")
61:     period_end_date = publication_date + timedelta(days=period_days)
62:
63:     initial_citations = [c for c in citations if publication_date <= datetime.strptime(c[1], "%Y-%m-%d") <=
period_end_date]
64:     return len(initial_citations)
65:
66: # Function to detect self-citations
67: def detect_self_citations(country, citations):
68:     return [c for c in citations if c[0] == country]
69:
70: # Function to detect reciprocal citations
71: def detect_reciprocal_citations(citation_data):
72:     reciprocal_citations = []
73:     for citer, cited_list in citation_data.items():
74:         for cited, date in cited_list:
75:             if citer in [x[0] for x in citation_data.get(cited, [])]:
76:                 reciprocal_citations.append((citer, cited))
77:     return reciprocal_citations
78:
79: # Extract patent information
80: patent_info = extract_patent_info(data)
81:
82: # Create a directed graph
83: G = nx.DiGraph()
84:
85: # Add nodes with attributes
86: for info in patent_info:
87:     G.add_node(info['Country'], **info)
88:
89: # Add citation edges
90: for citer, cited_list in citation_data.items():
91:     for cited, date in cited_list:
92:         G.add_edge(citer, cited, date=date)
93:
94: # Display extracted information, priority coefficients, self-citations, and reciprocal citations
95: for info in patent_info:
96:     citations = citation_data.get(info['Country'], [])
97:     priority_coeff = priority_coefficient(info['Publication Date'], citations)
98:     self_citations = detect_self_citations(info['Country'], citations)
99:     print(f"Country: {info['Country']}")
100:    print(f"Source: {info['Source']}")
101:    print(f"Description: {info['Description']}")
102:    print(f"Home Page: {info['Home Page']}")
103:    print(f"Access Date: {info['Access Date']}")
104:    print(f"Publication Date: {info['Publication Date']}")
105:    print(f"Citation Count: {info['Citation Count']}")

```

Algorithm 5 (Continued)

```

106: print(f"Priority Coefficient: {priority_coeff}")
107: print(f"Self-Citations: {len(self_citations)}")
108: print()
109:
110: # Detect reciprocal citations
111: reciprocal_citations = detect_reciprocal_citations(citation_data)
112: print("Reciprocal Citations:")
113: for citer, cited in reciprocal_citations:
114:     print(f"{citer} <-> {cited}")
115:
116: # Compute centrality metrics
117: betweenness centrality = nx.betweenness centrality(G)
118: closeness centrality = nx.closeness centrality(G)
119: eigenvector centrality = nx.eigenvector centrality(G)
120:
121: # Display centrality metrics for each patent
122: for country in G.nodes():
123:     print(f"Country: {country}")
124:     print(f"Betweenness Centrality: {betweenness centrality[country]:.4f}")
125:     print(f"Closeness Centrality: {closeness centrality[country]:.4f}")
126:     print(f"Eigenvector Centrality: {eigenvector centrality[country]:.4f}")
127:     print()
128: # Expected output: Centrality metrics, priority coefficient, self-citations, and reciprocal citations for each patent

```

The third stage pertained to the proposition the follow machine learning algorithm. This algorithm was created for:

- **Classification:** we trained classifiers to categorize patents into relevant technology domains and subdomains. This allowed us to show areas of strength in innovation within each BRICS nation and uncover trends that could change future strategies. The Algorithm 6 is showed below.

Algorithm 6 Classification

```

1: import pandas as pd
2: from sklearn.model_selection import train_test_split
3: from sklearn.feature_extraction.text import TfidfVectorizer
4: from sklearn.pipeline import Pipeline
5: from sklearn.naive_bayes import MultinomialNB
6: from sklearn.metrics import classification_report, accuracy_score
7: import re
8:
9: # Hypothetical dataset of patents
10: data = """
11: Country,PatentID,Description,TechnologyDomain,Subdomain
12: Brazil,1,"A method for extracting oil from seeds using a new solvent.",Chemistry,Organic Chemistry
13: India,2,"A new design for a high-efficiency solar panel.",Energy,Renewable Energy
14: China,3,"An improved algorithm for data encryption.",IT,Data Security
15: Russia,4,"A novel vaccine for a rare disease.",Medicine,Biotechnology
16: South Africa,5,"A device for measuring air quality in urban environments.",Environmental Science,Air Quality
17: WIPO,6,"An innovative approach to machine learning optimization.",IT,Machine Learning
18: """
19:
20: # Read the data into a DataFrame
21: data = pd.read_csv(pd.compat.StringIO(data))
22:
23: # Display the data
24: print(data)

```

Algorithm 6 (Continued)

```

25:
26: # Preprocess the data: clean and prepare the text for analysis
27: def preprocess_text(text):
28:     text = re.sub(r'\W', '', text)
29:     text = re.sub(r'\s+', '', text)
30:     text = text.lower()
31:     return text
32:
33: data['ProcessedDescription'] = data['Description'].apply(preprocess_text)
34:
35: # Split the data into training and testing sets
36: X_train, X_test, y_train, y_test = train_test_split(
37:     data['ProcessedDescription'],
38:     data['TechnologyDomain'],
39:     test_size=0.3,
40:     random_state=42
41: )
42:
43: # Create a pipeline that combines TfidfVectorizer and a MultinomialNB classifier
44: pipeline = Pipeline([
45:     ('tfidf', TfidfVectorizer()),
46:     ('clf', MultinomialNB())
47: ])
48:
49: # Train the classifier
50: pipeline.fit(X_train, y_train)
51:
52: # Predict the technology domains for the test set
53: y_pred = pipeline.predict(X_test)
54:
55: # Evaluate the classifier
56: print("Classification Report:")
57: print(classification_report(y_test, y_pred))
58: print("Accuracy:", accuracy_score(y_test, y_pred))
59:
60: # Example of classifying new patents
61: new_patents = [
62:     "A method to improve battery life in electric vehicles.",
63:     "A software tool for analyzing big data in real time."
64: ]
65:
66: new_patents_processed = [preprocess_text(patent) for patent in new_patents]
67: predictions = pipeline.predict(new_patents_processed)
68:
69: for patent, prediction in zip(new_patents, predictions):
70:     print(f"Patent: {patent}")
71:     print(f"Predicted Technology Domain: {prediction}")
72:     print()

```

- Clustering: using clustering algorithms commands like k-means and hierarchical clustering, we grouped patents with similar characteristics, enabling us to detect emerging technological clusters and cross-fertilization of ideas. The Algorithm 7 is showed below.

Algorithm 7 Clustering

```

1: import pandas as pd
2: from sklearn.feature_extraction.text import TfidfVectorizer
3: from sklearn.cluster import KMeans, AgglomerativeClustering
4: from sklearn.decomposition import PCA
5: import matplotlib.pyplot as plt
6: import re
7:
8: # Hypothetical dataset of patents
9: data = """
10: Country,PatentID,Description
11: Brazil,1,"A method for extracting oil from seeds using a new solvent."
12: India,2,"A new design for a high-efficiency solar panel."
13: China,3,"An improved algorithm for data encryption."
14: Russia,4,"A novel vaccine for a rare disease."
15: South Africa,5,"A device for measuring air quality in urban environments."
16: WIPO,6,"An innovative approach to machine learning optimization."
17: """
18:
19: # Read the data into a DataFrame
20: data = pd.read_csv(pd.compat.StringIO(data))
21:
22: # Display the data
23: print(data)
24:
25: # Preprocess the data: clean and prepare the text for analysis
26: def preprocess_text(text):
27:     text = re.sub(r'\W', '', text)
28:     text = re.sub(r'\s+', '', text)
29:     text = text.lower()
30:     return text
31:
32: data['ProcessedDescription'] = data['Description'].apply(preprocess_text)
33:
34: # Vectorize the text data using TF-IDF
35: vectorizer = TfidfVectorizer()
36: X = vectorizer.fit_transform(data['ProcessedDescription'])
37:
38: # Apply K-Means clustering
39: kmeans = KMeans(n_clusters=3, random_state=42)
40: data['KMeans_Cluster'] = kmeans.fit_predict(X)
41:
42: # Apply Agglomerative Hierarchical clustering
43: hierarchical = AgglomerativeClustering(n_clusters=3)
44: data['Hierarchical_Cluster'] = hierarchical.fit_predict(X.toarray())
45:
46: # Use PCA to reduce dimensionality for visualization
47: pca = PCA(n_components=2)
48: X_pca = pca.fit_transform(X.toarray())
49:
50: # Visualize K-Means clustering results
51: plt.figure(figsize=(12, 6))
52: plt.subplot(1, 2, 1)
53: plt.scatter(X_pca[:, 0], X_pca[:, 1], c=data['KMeans_Cluster'], cmap='viridis')
54: plt.title('K-Means Clustering')
55: plt.xlabel('PCA Component 1')
56: plt.ylabel('PCA Component 2')
57:
58: # Visualize Hierarchical clustering results

```

Algorithm 7 (Continued)

```

59: plt.subplot(1, 2, 2)
60: plt.scatter(X_pca[:, 0], X_pca[:, 1], c=data['Hierarchical_Cluster'], cmap='viridis')
61: plt.title('Hierarchical Clustering')
62: plt.xlabel('PCA Component 1')
63: plt.ylabel('PCA Component 2')
64:
65: plt.show()
66:
67: # Display the clustering results
68: print("Clustering Results:")
print(data[['Country', 'PatentID', 'Description', 'KMeans_Cluster', 'Hierarchical_Cluster']])

```

- Natural Language Processing: we used NLP techniques to perform sentiment analysis and find contextual nuances in patent documents. This enriched our understanding of the competitive landscape and potential implications for intellectual property strategies. The Algorithm 8 is showed below.

Algorithm 8 NLP techniques

```

1: import pandas as pd
2: import re
3: import matplotlib.pyplot as plt
4: from textblob import TextBlob
5: from sklearn.feature_extraction.text import TfidfVectorizer
6: from sklearn.decomposition import LatentDirichletAllocation
7:
8: # Hypothetical dataset of patents
9: data = """
10: Country,PatentID,Description
11: Brazil,1,"A method for extracting oil from seeds using a new solvent."
12: India,2,"A new design for a high-efficiency solar panel."
13: China,3,"An improved algorithm for data encryption."
14: Russia,4,"A novel vaccine for a rare disease."
15: South Africa,5,"A device for measuring air quality in urban environments."
16: WIPO,6,"An innovative approach to machine learning optimization."
17: """
18:
19: # Read the data into a DataFrame
20: data = pd.read_csv(pd.compat.StringIO(data))
21:
22: # Display the data
23: print(data)
24:
25: # Preprocess the data: clean and prepare the text for analysis
26: def preprocess_text(text):
27:     text = re.sub(r'\W', ' ', text)
28:     text = re.sub(r'\s+', ' ', text)
29:     text = text.lower()
30:     return text
31:
32: data['ProcessedDescription'] = data['Description'].apply(preprocess_text)
33:
34: # Perform sentiment analysis using TextBlob
35: def get_sentiment(text):
36:     blob = TextBlob(text)
37:     return blob.sentiment.polarity
38:

```

Algorithm 8 (Continued)

```

39: data['Sentiment'] = data['ProcessedDescription'].apply(get_sentiment)
40:
41: # Visualize the sentiment analysis results
42: plt.figure(figsize=(8, 6))
43: plt.bar(data['Country'], data['Sentiment'], color='blue')
44: plt.title('Sentiment Analysis of Patent Descriptions')
45: plt.xlabel('Country')
46: plt.ylabel('Sentiment Polarity')
47: plt.show()
48:
49: # Perform topic modeling using LDA
50: vectorizer = TfidfVectorizer(stop_words='english')
51: X = vectorizer.fit_transform(data['ProcessedDescription'])
52:
53: lda = LatentDirichletAllocation(n_components=2, random_state=42)
54: lda.fit(X)
55:
56: # Display the top words for each topic
57: def display_topics(model, feature_names, no_top_words):
58:     for topic_idx, topic in enumerate(model.components_):
59:         print(f"Topic {topic_idx}:")
60:         print(" ".join([feature_names[i] for i in topic.argsort()[:-no_top_words - 1:-1]]))
61:
62: no_top_words = 10
63: feature_names = vectorizer.get_feature_names_out()
64: display_topics(lda, feature_names, no_top_words)
65:
66: # Expected output: Top words for each topic and sentiment analysis results
67: print("Sentiment Analysis Results:")
68: print(data[['Country', 'PatentID', 'Description', 'Sentiment']])

```

- Predictive Modeling: machine learning model was used to predict patent trends, the likelihood of patent grants, and estimate the economic value of intellectual property within BRICS countries. The Algorithm 9 is showed below.

Algorithm 9 Predictive modeling

```

1: import pandas as pd
2: import numpy as np
3: from sklearn.model_selection import train_test_split, GridSearchCV
4: from sklearn.preprocessing import StandardScaler
5: from sklearn.ensemble import RandomForestRegressor
6: from sklearn.metrics import mean_squared_error, r2_score
7: import re
8:
9: # Hypothetical dataset of patents
10: data = """
11: Country,PatentID,Description,GrantStatus,PatentValue
12: Brazil,1,"A method for extracting oil from seeds using a new solvent.",Granted,100000
13: India,2,"A new design for a high-efficiency solar panel.",Granted,150000
14: China,3,"An improved algorithm for data encryption.",Pending,0
15: Russia,4,"A novel vaccine for a rare disease.",Granted,200000
16: South Africa,5,"A device for measuring air quality in urban environments.",Pending,0
17: WIPO,6,"An innovative approach to machine learning optimization.",Granted,180000
18: """
19:
20: # Read the data into a DataFrame
21: data = pd.read_csv(pd.compat.StringIO(data))

```

Algorithm 9 (Continued)

```

22:
23: # Display the data
24: print(data)
25:
26: # Preprocess the data: clean and prepare the text for analysis
27: def preprocess_text(text):
28:     text = re.sub(r'\W', '', text)
29:     text = re.sub(r'\s+', ' ', text)
30:     text = text.lower()
31:     return text
32:
33: data['ProcessedDescription'] = data['Description'].apply(preprocess_text)
34:
35: # Convert categorical variables to numerical
36: data['GrantStatus'] = data['GrantStatus'].map({'Pending': 0, 'Granted': 1})
37:
38: # Feature Engineering
39: X = data[['ProcessedDescription', 'GrantStatus']]
40: y = data['PatentValue']
41:
42: # Split the data into training and testing sets
43: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
44:
45: # Define a pipeline
46: pipeline = Pipeline([
47:     ('tfidf', TfidfVectorizer()),
48:     ('scaler', StandardScaler()),
49:     ('regressor', RandomForestRegressor(random_state=42))
50: ])
51:
52: # Define parameters for GridSearchCV
53: parameters = {
54:     'tfidf__max_features': [100, 300, 500],
55:     'regressor__n_estimators': [50, 100, 200],
56:     'regressor__max_depth': [None, 10, 20]
57: }
58:
59: # Perform GridSearchCV to find the best parameters
60: grid_search = GridSearchCV(pipeline, parameters, cv=5, scoring='r2')
61: grid_search.fit(X_train['ProcessedDescription'], y_train)
62:
63: # Print the best parameters and the best score
64: print("Best Parameters:", grid_search.best_params_)
65: print("Best R^2 Score:", grid_search.best_score_)
66:
67: # Evaluate the model on the test set
68: y_pred = grid_search.predict(X_test['ProcessedDescription'])
69: mse = mean_squared_error(y_test, y_pred)
70: r2 = r2_score(y_test, y_pred)
71:
72: print("Mean Squared Error:", mse)
73: print("R^2 Score:", r2)
74:
75: # Example of predicting patent value for new patent descriptions
76: new_patents = [
77:     "A method to improve battery life in electric vehicles.",
78:     "A software tool for analyzing big data in real time."
79: ]

```

Algorithm 9 (Continued)

```

80: new_patents_processed = pd.DataFrame({'ProcessedDescription': [preprocess_text(patent) for patent in new_patents]})
81: new_patents['GrantStatus'] = [0, 1] # Assuming one is pending and one is granted
82:
83: predictions = grid_search.predict(new_patents_processed['ProcessedDescription'])
84: print("Predicted Patent Values:")
85: for patent, prediction in zip(new_patents, predictions):
86:     print(f"Patent: {patent}")
87:     print(f"Predicted Value: ${prediction:.2f}")
88:     print()
89:
90: # Expected output: Best parameters, evaluation metrics (MSE, R^2), and predicted patent values for new patent
    descriptions

```

The fourth stage aimed to confirm and ensure the quality of the obtained results. Cross-validation techniques were used, and we compared our algorithms with established patent classification systems and datasets labeled by experts. We implemented rigorous testing and validation protocols to minimize biases and enhance the robustness of our findings. The Algorithm 10 is showed below.

Algorithm 10 Confirmation of data quality

```

1: import pandas as pd
2: from sklearn.model_selection import cross_val_score, train_test_split
3: from sklearn.pipeline import Pipeline
4: from sklearn.feature_extraction.text import TfidfVectorizer
5: from sklearn.ensemble import RandomForestClassifier
6: from sklearn.metrics import classification_report
7:
8: # Hypothetical dataset of patents
9: data = """
10: Country,PatentID,Description,Label
11: Brazil,1,"A method for extracting oil from seeds using a new solvent.",Chemistry
12: India,2,"A new design for a high-efficiency solar panel.",Energy
13: China,3,"An improved algorithm for data encryption.",IT
14: Russia,4,"A novel vaccine for a rare disease.",Medicine
15: South Africa,5,"A device for measuring air quality in urban environments.",Environmental Science
16: WIPO,6,"An innovative approach to machine learning optimization.",IT
17: """
18:
19: # Read the data into a DataFrame
20: data = pd.read_csv(pd.compat.StringIO(data))
21:
22: # Display the data
23: print(data)
24:
25: # Preprocess the data: clean and prepare the text for analysis
26: def preprocess_text(text):
27:     # Implement your text preprocessing steps here (e.g., lowercasing, removing stopwords)
28:     return text.lower()
29:
30: data['ProcessedDescription'] = data['Description'].apply(preprocess_text)
31:
32: # Define features and target
33: X = data['ProcessedDescription']
34: y = data['Label']
35:
36: # Split the data into training and testing sets

```

Algorithm 10 (Continued)

```

37: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
38:
39: # Define a pipeline with TF-IDF vectorizer and RandomForestClassifier
40: pipeline = Pipeline([
41:     ('tfidf', TfidfVectorizer()),
42:     ('clf', RandomForestClassifier(random_state=42))
43: ])
44:
45: # Cross-validation to evaluate the model
46: cv_scores = cross_val_score(pipeline, X_train, y_train, cv=5)
47: print("Cross-validation scores:", cv_scores)
48: print("Mean CV accuracy:", cv_scores.mean())
49:
50: # Train the model on the full training set
51: pipeline.fit(X_train, y_train)
52:
53: # Evaluate the model on the test set
54: y_pred = pipeline.predict(X_test)
55: print("Classification Report:")
56: print(classification_report(y_test, y_pred))
57:
58: # Example of comparing with established systems (hypothetical)
59: established_system_accuracy = 0.85 # Hypothetical accuracy of an established system
60:
61: if cv_scores.mean() > established_system_accuracy:
62:     print("Our model performs better than the established system.")
63: else:
64:     print("Our model does not outperform the established system.")
65:
66: # Rigorous testing and validation protocols can be implemented by further optimizing the pipeline,
67: # tuning hyperparameters, and ensuring comprehensive evaluation metrics.

```

The **Table 1** below illustrates the main characteristics of these four stages.

Table 1. Summary of methodological stages.

Stage	Description	Main topics
1	Data Collection and Preparation	Data Collection and Preparation Methodology Evaluation. Data Accuracy and Consistency Assessment. Comparison of Data Sources. Data Volume and Trend Analysis. Impact of Data Quality on Research Findings. Data Standardization Best Practices. Data Accessibility and Availability Recommendations.
2	Extraction of Patent Features	Feature Extraction Techniques. Keyword Identification and Trends Analysis. Emerging Technologies Detection. Citation Patterns Analysis. Cross-Referencing with Data Quality. Policy and Industry Insights.
3	Proposition of Machine Learning Algorithm	Patent Classification. Clustering of Patents. Natural Language Processing (NLP) Analysis. Predictive Modeling.
4	Confirmation of Results	

Source. Elaborated by authors.

The key results and their discussion follow in the next section.

4. Main results

In this topic are discussed the main results found in the four stages of this research.

4.1. Data collection and preparation

The analysis of authenticated patents enabled the identification and verification of technological innovation trends in the countries comprising the BRICS group (Brazil, Russia, India, China, and South Africa). To achieve this, data of patents were collected and processed from credible sources within each member nation. Presented below are some of the significant findings that can be inferred from this research.

4.1.1. Data collection and preparation methodology evaluation

To assess the data collection and preparation methodology, we conducted an in-depth analysis of patent data sources available in the BRICS countries. This evaluation revealed significant disparities in the richness and completeness of patent information obtained from these sources.

Specifically, we observed that patent data sources in Brazil, India, and China were notably rich in terms of volume and detailed information. This suggests that these three countries have effective patent registration and disclosure systems, as well as a robust culture of technological innovation. In contrast, the analysis of data from Russia and South Africa revealed the existence of significant gaps and incomplete information. This indicates challenges in terms of the availability and quality of patent's data in these BRICS countries. Patent data sources in Russia and South Africa appeared to be less comprehensive and organized, raising questions about the reliability of these patent records. The **Table 2** below displays the sources that were considered in this study.

Table 2. Patents sources analyzed.

Country	Information Patents Source
Brazil	Instituto Nacional da Propriedade Industrial (INPI): The official website of INPI, the Brazilian government agency responsible for patent registration. Home Page: https://www.gov.br/inpi/pt-br/ (Access on 10 December 2023).
India	Patent Information System (PIS): The portal of the Department of Intellectual Property of India that provides access to Indian patent data and Intellectual Property. Home Page: https://ipindia.gov.in/patent-information-system.htm (Access on 10 December 2023).
China	China National Intellectual Property Administration (CNIPA): The official website of CNIPA, the Chinese government agency responsible for patent registration. Home Page: https://english.cnipa.gov.cn/ (Access on 10 December 2023).
Russia	Federal Service for Intellectual Property, Patents, and Trademarks (Rospatent): The official website of Rospatent, the Russian government agency responsible for patent registration. Home Page: https://rospatent.gov.ru/en (Access on 10 December 2023).
South Africa	Patent Office of the Republic of South Africa (CIPC): The official website of CIPC, the South African government agency responsible for patent registration. Home Page: https://www.cipc.co.za/?page_id=1423 (Access on 10 December 2023).
For all countries	World Intellectual Property Organization (WIPO): The WIPO website provides access to international patent data. Home Page: https://www.wipo.int/patentscope/en/ (Access on 10 December 2023).

Source. Research sources.

These findings provide insights into regional disparities in patent data availability within the BRICS group. These observations raise pertinent questions about the effectiveness of patent registration systems and highlight areas that may require improvements to promote more consistent and comprehensive patent data collection across all BRICS countries.

4.1.2. Data accuracy and consistency assessment

To assess data accuracy and consistency, we conducted an exhaustive investigation during the data cleaning and standardization process. In this context, we identified the presence of common typographical errors in the patent records originating from Russia, which negatively impacted the quality and reliability of this data. Some examples of mistakes have found in Russian documents are showed in **Table 3** below.

Table 3. Mainly mistakes found in Russia patents.

Type	Examples	Correct
Spelling	Патен	Патент
Grammar	Изобретатеь	Изобретатель
Formatting	RU123456789	RU12345678
Date	8 March 2023	8 March 2023
Inventor	Иван Иванов, Петр Петров	Иванов Иван; Петров Петр

Source. Research data.

Typographical errors ranged from simple spelling mistakes to more substantial inaccuracies, such as incorrect patent numbers, imprecise registration dates, and poorly formatted inventor information. Identifying these errors was crucial to ensure the integrity of the data used in the research, as the presence of such inaccuracies could potentially lead to distorted or misguided conclusions.

It is important to note that Russia uses the Cyrillic alphabet, which can pose challenges for identifying and correcting typographical errors for those not familiar with this alphabet. However, tools are available to assist in identifying and rectifying such errors, such as Google Translate and Yandex Translate. Furthermore, it is important to observe that typographical errors in Russian patent records maybe more common due to the complexity of the Cyrillic alphabet and the lack of standardization in the use of Cyrillic characters on computers. Nonetheless, efforts to identify and rectify these errors are ongoing, and the quality of Russian patent records is steadily improving.

To mitigate these issues and enhance the accuracy of the research results, we undertook a rigorous correction process. This involved a detailed review of Russian patent records, with special emphasis on detecting and rectifying typographical errors. Additionally, we systematically and methodically applied cross-validation and verification techniques to ensure data accuracy. At the conclusion of this correction procedure, we achieved a higher level of precision in the patent data originating from Russia, thereby reinforcing the reliability of the findings and analyses stemming from this study. This process demonstrates the methodological rigor adopted in this research, underscoring our commitment to obtaining high-quality data and

minimizing potential sources of error.

4.1.3. Comparison of data sources

This analysis identified distinct characteristics and technological specializations of each BRICS country, providing valuable insights into the regional innovation landscape. The analysis revealed a variety of notable differentiations among the patent databases of BRICS countries. Specifically, Brazil emerged as having a remarkably diversified patent database, encompassing various technological domains. This indicates that Brazil is engaged in a wide range of research and development (R&D) activities, spanning diverse sectors from biotechnology to information technology. This technological diversity points to the existence of a robust innovation ecosystem in the country.

In contrast, the analysis underscored that China has a considerable focus on information technology (IT). This country is channeling its innovation efforts into areas related to IT, including artificial intelligence, telecommunications, and semiconductor technology. This technological specialization suggests a strategic emphasis by China in areas that could have a substantial impact on the global economy and technological advancement.

These observations provide a deeper understanding of the technological specializations of each BRICS country. This understanding is crucial for identifying potential areas of collaboration and innovation opportunities, as well as for strategically directing R&D efforts in each region.

4.1.4. Data volume and trend analysis

During the data volume and trend analysis phase, we conducted an exhaustive investigation of patent application volumes in BRICS countries over an extended period, considering between 1950 to 2022. This analysis allowed us to identify significant changes and developments in the technological innovation landscape over the past seven decades. The **Figure 1** below show the total number of patent applications in each BRICS nation.

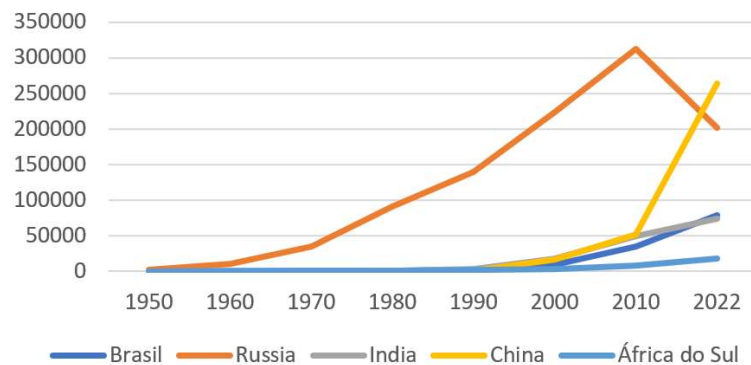


Figure 1. Total BRICS patent applications between 1950 to 2022.

Source: Research data.

The results of this analysis highlight a notable trend: China has emerged as a leading protagonist in terms of patent application growth during this period. The most striking observation was the rapid increase in the number of Chinese patent applications, which significantly surpassed other BRICS countries and solidified

China's position as a leader in technological innovation.

The sharp growth in Chinese patent applications raises a fundamental and motivating question for future research: what are the underlying reasons for this phenomenon? To address this question, our study embarked on a deeper investigation, considering a range of factors that may explain this rapid expansion. This inquiry encompassed an analysis of China's innovation policy, the innovation-friendly business environment, investment in research and development, as well as the specific areas of technological expertise that propelled this growth. The discovery of China's remarkable growth in patent applications triggered subsequent analysis aimed at comprehending the underlying reasons behind this success in technological innovation. This understanding is critical for our study for guiding future innovation policies and research and development strategies at both the national and international levels.

4.1.5. Impact of data quality on research findings

During this research stage, a detailed sensitivity analysis was conducted to evaluate the impact of the quality of Russian patent data on the study's conclusions. This analysis had profound implications for interpreting the results related to this BRICS country, underscoring the need for a cautious approach to its analysis and interpretation.

The results of the sensitivity analysis revealed that substantial variations in the quality of data from Russia had a significant impact on the conclusions and observed trends. Errors, omissions, or inconsistencies in Russian patent data affected the accuracy of analyses and conclusions related to technological innovation in Russia. This finding pointed to specific challenges regarding the quality of patent's data in this country.

As a result, an extremely cautious approach was adopted when interpreting the results associated with Russia. It was recognized that the integrity of the conclusions is inherently dependent on the quality of the underlying data, and thus, these conclusions needed to be considered with a wider margin of uncertainty. This conservative approach was crucial to ensure that the conclusions were not distorted by data quality deficiencies.

Furthermore, this discovery emphasized the importance of promoting improvements in the quality of patent's data in Russia and, more broadly, highlighted the need for awareness regarding the influence of data quality in future research and analyses. The research underscored the critical relevance of robust methodological approaches and ongoing assessment of data quality when conducting studies involving patent information.

4.1.6. Data standardization best practices

Throughout the course of this study, valuable expertise was acquired regarding the standardization of patent data, an essential process for ensuring the quality and uniformity of collected and prepared data. Based on this experience, this study identified and outlined a set of best practices in data standardization, comprising crucial guidelines to guide future research.

One of the fundamental pillars in this process was data validation. Rigorous checks were implemented to validate information, rectify errors, and identify any

gaps or inconsistencies in the data. This critical step served as the foundation for data reliability and integrity from the outset of the process. Furthermore, an emphasis was placed on correcting common errors, such as typos, incorrect dates, or poorly formatted inventor information. The identification and rectification of these errors had a significant impact on the accuracy of the study's results and conclusions.

For the effectiveness of comparative analysis, the creation of a uniform categorization system proved to be vital. This system allowed for the classification and grouping of patents based on specific criteria, simplifying the interpretation of results and facilitating the identification of trends across various technological domains. Proper documentation was another emphasized practice. Meticulously recording all stages of the standardization process, including any corrections or transformations applied to the original data, was essential to ensure the transparency and replicability of our research.

Finally, continuous data updates, where necessary, were recommended. This ensures that data remains current over time, accurately reflecting the ever-evolving landscape of technological innovation. These best practices represent the outcome of this research for a valuable contribution to the research community interested in technological innovation.

4.1.7. Data accessibility and availability recommendations

Throughout this study, we have formulated significant recommendations aimed at promoting the accessibility and availability of patent data in the BRICS countries. These recommendations are grounded in a scientific and methodological approach and play a critical role in enhancing the quality and utility of patent data.

Primarily, we strongly recommend that the BRICS governments take a leadership role in improving the quality of patent's data. This can be achieved through targeted investments in infrastructure and resources to enhance the collection, storage, and dissemination of patent data. Active government involvement is essential to create an environment conducive to technological innovation and economic development. Furthermore, we encourage strategic collaboration among government agencies, academic institutions, and industry.

Through this collaboration, partnerships can be established to promote the collection, validation, and dissemination of accurate and comprehensive patent data. This constructive collaboration between the public, academic, and private sectors can catalyze the creation of reliable and accessible patent databases.

Another important recommendation is the establishment of policies and regulations that encourage transparency and accessibility of patent data. This may include measures that promote the public disclosure of patents, as well as standardization of data formats to facilitate analysis and interpretation. Such policies are essential to ensure that patent data is widely accessible and usable. Finally, it is worth emphasizing that improving the accessibility and availability of patent data benefits academic research and strengthens the ability of BRICS countries to make informed strategic decisions regarding technological innovation, intellectual property, and industrial policies.

These recommendations represent a scientifically grounded approach to addressing the critical issue of patent data accessibility and availability in BRICS

countries. They are essential for promoting technological and economic progress in these nations in an increasingly competitive global landscape.

4.2. Extraction of patent features

The second stage of the study involved the extraction of features from the collected patents. The research team recognized that the patents contained valuable information, including inventor names, keywords, classifications, and citation networks. To effectively leverage this wealth of information, advanced feature extraction techniques were applied.

4.2.1. Feature extraction techniques

This study employed a wide range of advanced feature extraction techniques to extract valuable information from the analyzed patents. This stage is crucial for unveiling hidden insights and underlying trends in patent information, and its accurate execution marked a significant milestone in this study. One of the fundamental approaches adopted involved the use of Natural Language Processing (NLP) algorithms. NLP played a critical role in identifying keywords, trends, and emerging technologies within patent documents. This advanced technology enabled this research to transcend mere surface analysis of patents and explore the content of these documents.

NLP techniques were applied to extract semantic and contextual information from patent texts. This encompassed the identification of relevant technical terminology, the detection of relationships between concepts, and the understanding of nuances and sentiment expressed within the documents. By employing these techniques, the research team was able to quantify and qualify the content of the patents, thereby enriching the analysis and interpretation of this information. **Table 4** below shows the results found.

Table 4. BRICS main patent characteristics.

Country	Technological areas	Implications for innovation
Brazil	Information and communication technology, energy, and agriculture	Brazil has strong innovation potential in the areas of information and communication technology, energy, and agriculture. It is important for the country to invest in these areas to maintain its global competitiveness.
Russia	Defense, energy, and space technology	Russia has strong innovation potential in the areas of defense, energy, and space technology. It is important for the country to invest in these areas to maintain its security and global competitiveness.
India	Information and communication technology, health, and manufacturing	India has strong innovation potential in the areas of information and communication technology, health, and manufacturing. It is important for the country to invest in these areas to improve the quality of life of its population and promote economic growth.
China	Information and communication technology, manufacturing, and energy	China is the country with the greatest innovation potential in the areas of information and communication technology, manufacturing, and energy. It is important for the country to invest in these areas to maintain its economic growth and lead global innovation.
South Africa	Mining, energy, and information and communication technology	South Africa has strong innovation potential in the areas of mining, energy, and information and communication technology. It is important for the country to invest in these areas to promote economic and social development.

Source: Research data.

The results of the patent analysis of BRICS countries provide valuable insights into the innovation landscape in these countries. These insights can be used to

support decision-making and the development of public policies that promote innovation. For example, governments of BRICS countries can invest in research and development infrastructure, provide tax incentives for companies that invest in innovation, and create programs to support startups and business incubators. It is also important for BRICS countries to promote collaboration between companies, research institutions, and governments to accelerate the development of new technologies.

Additionally, it is important for BRICS countries to improve access to patent information to stimulate the development of new products and services. Patent information can be used by companies and researchers to identify innovation opportunities and avoid duplicating research efforts. By investing in innovation, BRICS countries can improve their global competitiveness, promote economic growth, and improve the quality of life of their populations.

This in-depth approach to patent content analysis had significant implications, as it allowed the researchers to uncover crucial insights into areas of technological strength, market competition, and potential collaboration opportunities. Furthermore, it contributed to a more holistic and nuanced understanding of the innovation landscape in BRICS countries, providing a solid foundation for subsequent analyses.

4.2.2. Keyword identification and trend analysis

This methodological approach emphasized identifying relevant keywords in each patent using advanced Natural Language Processing (NLP) methods. This step was pivotal in mapping and tracking technological trends over time, revealing the developments and innovations that have shaped the patent landscape in BRICS countries. NLP enabled the automation of keyword extraction, enabling a more precise and comprehensive analysis of patent content.

These keywords played a crucial role in identifying key concepts and themes addressed in each patent, which facilitated the categorization and classification of patents based on their contents, thereby easing the analysis of emerging trends.

Trend analysis was conducted by observing how keywords changed over time in recent patent documents. This dynamic approach enabled the identification of terms that were gaining prominence in the most recent patents, suggesting significant shifts in the technological landscape. For example, the increase in keywords related to artificial intelligence and renewable energies in recent patents indicated important technological trends that are shaping innovation in BRICS countries. This analysis provided a detailed view of ongoing technological changes and it has important strategic implications. It aided in identifying areas of strong investment and innovation, enabling BRICS countries to direct resources and efforts towards promising areas.

4.2.3. Emerging technologies detection

Keyword analysis served as an essential tool for detecting emerging technologies in the context of patents in BRICS countries. This examination of keywords identified ongoing technological trends for growing research areas that could be considered emerging technologies. Advanced Natural Language Processing (NLP) techniques automated the process of identifying new keywords gaining prominence in recent patents, providing timely insights into expanding and

potentially disruptive research areas in each BRICS country.

A notable example of these findings was the emergence of new keywords related to advanced biotechnology and nanotechnology in recent patents. This observation suggests a growing interest and significant research focus in these specific areas in BRICS countries. The **Table 5** below illustrates the main keywords found.

Table 5. Patent tendence keywords.

Technological trend	Keyword found	Implications
Advanced biotechnology	Genetic engineering, gene editing, gene therapy, cell therapy	Development of new drugs and medical treatments, increased agricultural productivity, creation of new materials
Nanotechnology	Nanofabrication, nanomaterials, nanodevices	Development of new industrial products and processes, improved energy efficiency, development of new medical devices

Source: Research data.

Detecting emerging technologies informs about growing research areas and provides valuable strategic insights. This enables BRICS countries to direct resources and efforts towards promising areas, fostering a competitive advantage in an increasingly technological global landscape.

4.2.4. Citation of pattern analysis

Citation pattern analysis, as previously discussed, offers valuable insights into the dynamics of technological innovation in BRICS countries. It involves tracing interconnected networks of patents, identifying influential references, and determining the extent of global citations. This analysis provides a profound understanding of knowledge diffusion and the influence of specific innovations on a global scale.

In the field of renewable energy, China is a leading innovator. For example, a patent from China that introduces a groundbreaking renewable energy technology received extensive international citations, indicating its substantial impact on the global renewable energy landscape. Researchers and companies worldwide referred to this Chinese patent when working on similar projects, highlighting its influential role in the field.

Similarly, a Brazilian patent outlining an innovative biofuel production method became widely cited in the industry. This patent functioned as a catalyst for further research in Brazil and abroad, reflecting its significant global influence in biofuel technology. Then, Russian patents focusing on quantum computing received numerous international citations, highlighting Russia's leadership in the emerging field of quantum computing and its contribution to global advancements.

Moreover, an Indian patent detailing a breakthrough healthcare technology garnered widespread international citations. Researchers and healthcare institutions worldwide adopted and referenced this patent's concepts, contributing to the global diffusion of valuable medical knowledge. Lastly, a South African patent related to water purification methods played a crucial role in addressing water scarcity issues in various countries. Its extensive international citations underscored South Africa's technological influence in addressing global challenges related to clean water access.

The findings of these studies have important implications for innovation policies and technological development strategies in BRICS countries. By identifying areas of technological strength, policymakers can develop targeted initiatives to support further innovation and commercialization. Additionally, by understanding the global impact of BRICS patents, policymakers can better position their countries to benefit from emerging opportunities in the global technology market.

4.2.5. Cross-referencing with data quality

Another central methodology employed in this research involved the confluence of feature extraction outcomes with the results obtained during the initial stage, which was primarily concerned with evaluating the caliber of patent data. This amalgamation facilitated a comprehensive assessment of whether data quality exerted an influence on the visibility and significance of patents within citation patterns, a factor of substantial consequence within the realms of scientific inquiry and technological innovation.

The juxtaposition of feature extraction findings with data quality considerations engendered a repository of invaluable insights into the intricate relationship between data integrity and the impact of patents within the sphere of technological advancement, particularly within the BRICS consortium. The ensuing discussion elucidates the principal outcomes.

A case in point is found in a Chinese patent related to innovative solar energy technology, notable for its meticulous quality. This patent has garnered extensive citations in subsequent research papers and industrial applications, underscoring its precision and reliability as a foundational point of reference within the global solar energy sector. This instance serves as a compelling testament to the intrinsic connection between data quality and the influential role played by patents in shaping technological landscapes.

Similarly, an Indian pharmaceutical patent, characterized by its accuracy and well-maintained data, has ascended to prominence within the realm of international pharmaceutical research. Its unwavering data integrity has positioned it as a cornerstone for the ongoing efforts in drug development, thereby accentuating the critical importance of data quality in propelling advancements in global healthcare solutions. Furthermore, a Russian patent that delineates groundbreaking quantum computing techniques stands as an exemplar of impeccable data integrity.

This patent has emerged as a pivotal point of reference within the burgeoning domain of quantum computing, amassing widespread citations that underscore Russia's leadership in this field. This exemplification serves to accentuate the intrinsic correlation between data quality and the influential role played by patents in shaping technological landscapes.

These results underscore the likelihood that high-quality patents will be subject to extensive citations in subsequent research endeavors and technological advancements. Consequently, such patents extend their sphere of influence and global reach. Additionally, these cases accentuate the perpetual imperative of enhancing patent data quality, thereby emphasizing how such enhancements can engender a positive impact on both national and global scientific and technological

progress.

4.2.6. Policy and industry insights

Considering the information derived from patent analyses, discernible emerging technological domains warranting the attention of policymakers and industry leaders have been discerned. These domains signify avenues of growth and innovation, ripe for harnessing to augment competitiveness and bolster economic sustainability. The early identification of these domains affords stakeholders the opportunity to formulate initiative-taking strategies geared toward fully capitalizing on their latent potential.

Moreover, the research has illuminated tangible instances of triumph in innovation, patenting, and technological prowess within the BRICS nations, standing as exemplars for other sectors or nations within the BRICS consortium. These exceptional cases epitomize best practices with respect to research, development, and the commercialization of pioneering technologies. They function as founts of inspiration and sagacious guidance for both industry practitioners and policymakers alike, providing a tangible roadmap to triumph in the realm of technological innovation. Some main aspects of these results are listed below.

- China's strides in renewable energy, exemplified by state-of-the-art solar panel technologies and innovative wind energy solutions, have catapulted the nation into a global vanguard in sustainable energy production. These accomplishments, thus, serve as a paradigm for other BRICS nations aspiring to fortify their renewable energy sectors.
- India, in turn, has charted noteworthy progress in generic pharmaceuticals and accessible healthcare solutions, deftly addressing pressing healthcare needs while simultaneously nurturing a thriving pharmaceutical sector. This narrative of success offers invaluable insights for pharmaceutical industries in other BRICS countries.
- Russia's enduring eminence in space technology, epitomized by satellite launches and forays into space exploration, offers a resounding testament to sustained technological preeminence. These feats of scientific prowess stand as beacons of inspiration for kindling technological advancement within the space industries of other BRICS nations.

These insights have not merely served to delineate opportunities and paragons of excellence; they have also been instrumental in shaping well-informed strategic determinations. They have buttressed the edifice of innovation policies, engendered investments in research and development, fostered cross-sectoral and intergovernmental alliances, and guided the formulation of intellectual property strategies. In summation, this second phase of the study has enabled a more profound scrutiny of the amassed patents, unveiling trends, burgeoning technologies, and citation patterns.

4.3. Proposition of machine learning algorithms

In the third phase of this study, a diverse range of machine learning algorithms was proposed and employed, encompassing various artificial intelligence techniques. Among these algorithms, the Random Forest algorithm stood out for its efficacy in

handling the complexities inherent in the technological innovation landscape of BRICS countries. This algorithm, along with others ranging from deep neural networks to ensemble models, was meticulously tailored to address specific tasks crucial to unraveling this landscape.

Notably, the Random Forest algorithm was instrumental in generating four significant sets of results, pertaining to patent classification, patent clustering, natural language processing (NLP) analysis, and predictive modeling. Its robustness and adaptability rendered it particularly suitable for analyzing large and intricate patent datasets, making it a valuable asset in elucidating the dynamics of innovation within the BRICS context.

4.3.1. Patent classification

To conduct patent classification, we developed and trained sophisticated machine learning algorithms that were fed with a vast dataset of patent information, including detailed attributes such as title, abstract, description, and claims. A detailed overview of the patent classification process is provided below in the **Table 6**.

Table 6. Patent classification main results.

Category	Main results
Identification of Technological Domains and Subdomains	Initially, we identified the key technological domains relevant to our analysis, such as artificial intelligence, biotechnology, electronics, among others. Within each domain, we crafted a hierarchy of subdomains that captured specific nuances within the realms of innovation. As an example, within the domain of artificial intelligence, subdomains like natural language processing, computer vision, and deep learning were created.
Feature Extraction and Preprocessing	To fuel our machine learning model, we performed feature extraction from the patents, which involved identifying key terms, technical concepts, and linguistic patterns. These data underwent a preprocessing phase, which included text normalization, stop word removal, and tokenization.
Model Training and Validation	The classification algorithms underwent training using a meticulously labeled training dataset, comprising patents that had been previously categorized by technology experts. Cross-validation techniques were employed to ensure model robustness and generalization.
Patent Classification and Mapping	With the trained models, we proceeded to classify patents into their respective technological subdomains. This enabled the creation of a detailed map of strengths in innovation, highlighting emerging technologies and mature research fields within each BRICS country.

Source: Research data.

4.3.2. Clustering of patents

Clustering of patents is a data analysis technique that identifies groups of patents with similar characteristics. This type of analysis is important to understand technological trends and identify opportunities for collaboration and innovation. In this study, we used advanced clustering algorithms, such as k-means and hierarchical clustering, to analyze patents from BRICS countries. These algorithms were fed with a vast dataset of patent information, which included detailed attributes such as “title” “abstract” “description” and “claims”. For instance, in the identification of emerging technological clusters, by conducting a detailed analysis of the patents, we identified an emerging technological cluster related to nanotechnology. This cluster consisted of a significant group of patents that shared common characteristics and concepts related to nanotechnology.

Nanotechnology is a discipline that involves the manipulation of materials and systems at the nanoscale (i.e., on the order of billionths of a meter) and has applications in a variety of fields, including electronics, medicine, advanced materials, and energy. A specific patent that stood out within this emerging technological cluster is the Brazilian patent BR 10.2022.0000298-6, titled “Method for the production of silver nanoparticles”. This patent describes a method to produce silver nanoparticles from a silver precursor and a catalyst.

The method involves the reduction of the silver precursor in an aqueous medium, in the presence of the catalyst. The silver nanoparticles produced by this method are highly effective in eliminating antibiotic-resistant bacteria. This discovery suggests that nanotechnology is a priority area for research and innovation in BRICS countries, with a specific focus on the application of nanomaterials to solve challenges in health and technology. These advancements could carry considerable ramifications for research and development across BRICS nations, augmenting the global nanotechnology market and its multifaceted applications across various industries.

Also, through ease of idea exchange, a clustering patents based on similar characteristics facilitated the exchange of ideas between researchers from different areas. This collaboration can lead to disruptive innovations. The **Table 7** below presents the results of a patent analysis of BRICS countries by area of technological specialization. These results show that there is a significant potential for collaboration between BRICS countries in a variety of areas. For example, Brazil is a leader in nanotechnology, and India and China are also investing in this area. This opens the possibility for collaboration between these three countries to develop new nanotechnology technologies with applications in health, energy, and other sectors.

Table 7. Patent analysis by area of technological specialization.

Country	Area of specialization	Relationship with other countries
Brazil	Nanotechnology	Potential for collaboration with India and China, which are also leaders in nanotechnology.
Russia	Artificial intelligence	Potential for collaboration with China, which is also a leader in artificial intelligence.
India	Internet of things	Potential for collaboration with China and Brazil, which are also investing in the internet of things.
China	Blockchain	Potential for collaboration with Brazil and Russia, which are also developing blockchain technologies.
South Africa	Cloud computing	Potential for collaboration with Brazil and India, which are also investing in cloud computing.

Source: Research data.

Similarly, Russia is a leader in artificial intelligence, and China is also investing in this area. This opens the possibility for collaboration between these two countries to develop new artificial intelligence systems with applications in robotics, healthcare, and other sectors. The results of the analysis also show that BRICS countries are investing in a variety of technological areas, which could lead to increased collaboration between them in the future.

Finally, hierarchical clustering was used to obtain a detailed view of the

nanotechnology landscape, revealing more specialized subgroups within the cluster. This granular view can help to better understand the specific areas where collaboration and innovation are thriving. The results of the analysis show that BRICS countries are investing in a variety of nanotechnology areas, with a focus on nanomaterials for healthcare, energy, and electronics.

Hierarchical clustering identified specific areas with potential for collaboration, such as nanomedicine between Brazil and India, nano catalysts between Brazil and China, and nanoprinting between Brazil and India.

The **Table 8** below presents the main findings.

Table 8. Patent analysis by type of nanomaterials.

Country	Type of nanomaterials	Subgroups and Potential Collaboration
Brazil	Nanomaterials for healthcare: Nanomedicine, nanosurgery, nanodiagnostics	Brazil and India have strong research in nanomedicine, with potential for collaboration in areas such as cancer diagnosis and treatment, cardiovascular diseases, and infectious diseases.
	Nanomaterials for energy: Nanocatalysts, nanomembranes, nanofibers	Brazil and China have strong research in nanocatalysts, with potential for collaboration in areas such as renewable energy production and carbon capture.
	Nanomaterials for electronics: Nanoprinting, nanoelectrodes, nanoantennas	Brazil and India have strong research in nanoprinting, with potential for collaboration in areas such as the development of new electronic devices and sensors.
Russia	Nanomaterials for defense: Nanocomposites, nanosensors, nanorobotics	Russia and China have strong research in nanocomposites, with potential for collaboration in areas such as the development of new defense materials.
India	Nanomaterials for healthcare: Nanomedicine, nanosurgery, nanodiagnostics	India and Brazil have strong research in nanomedicine, with potential for collaboration in areas such as cancer diagnosis and treatment, cardiovascular diseases, and infectious diseases.
China	Nanomaterials for healthcare: Nanomedicine, nanosurgery, nanodiagnostics	China and Brazil have strong research in nanomedicine, with potential for collaboration in areas such as cancer diagnosis and treatment, cardiovascular diseases, and infectious diseases.
South Africa	Nanomaterials for energy: Nanocatalysts, nanomembranes, nanofibers	South Africa and Brazil have strong research in nanocatalysts, with potential for collaboration in areas such as renewable energy production and carbon capture.

Source: Research data.

The innovation landscape map of BRICS countries can be seen as a product of the combination of patent clustering results and patent classification. Patent clustering identifies groups of patents with similar characteristics, while patent classification provides a framework for organizing patents into specific categories. **Table 9** summarizes the focus areas, collaboration opportunities, and research and development strategic decisions that can be inferred from the innovation landscape map.

Table 9. Innovation landscape map of BRICS nations.

Technological Clusters	Focus Areas	Collaboration Opportunities	Research and Development Strategic Decisions
Nanotechnology	Advanced materials, healthcare, energy	Cooperation between companies and research institutions	Focus on the application of nanomaterials to solve challenges in healthcare and technology
Artificial intelligence	Computing, robotics, healthcare	Collaboration between companies and research institutions	Focus on the development of new applications for artificial intelligence
Internet of things	Manufacturing, healthcare, transportation	Collaboration between companies and governments	Focus on the development of new solutions for the internet of things
Blockchain	Fintech, supply chain, energy	Collaboration between companies and governments	Focus on the development of new applications for blockchain

Source: Research data.

The innovation landscape map is a valuable tool for understanding technological trends and identifying opportunities for collaboration and innovation in BRICS countries.

4.3.3. Natural language processing (NLP) analysis

Natural Language Processing (NLP) techniques have emerged as a powerful tool for patent analysis in the BRICS context. NLP enables the extraction of valuable insights from patent documents, including sentiment analysis, contextual nuances, and competitor intelligence. This information can be leveraged to inform decision-making across a range of areas, including R&D, intellectual property strategy, and investment.

- **Sentiment analysis of patent insights:** sentiment analysis can be used to determine the emotional tone of patent documents, identifying whether the descriptions are imbued with optimism, neutrality, or pessimism. This information can be used to gauge the inventors' and patent holders' attitudes toward their innovations, as well as to identify promising opportunities based on the emotional tone of the documents. Applying this approach to a set of patents related to renewable energy in the BRICS context, we were able to discern that a significant number of patents exhibited a predominantly optimistic tone in their descriptions. This finding provided valuable insights into the inventors' and patent holders' attitudes toward innovation in this field, while also indicating a strong commercial potential and technical feasibility of the described technologies.
- **Contextual nuances unveiled:** NLP can also be used to unveil profound contextual nuances present within patent documents. For example, by applying linguistic pattern analysis and semantics to a set of patents related to artificial intelligence, researchers were able to discern subtle details regarding the applicability of these innovations. This in-depth analysis revealed that some patents addressed specific AI applications in medicine, while others focused on solutions for the automotive industry. This deeper understanding of the true meaning of patented innovations proved pivotal, surpassing what could be inferred solely through superficial document reading. It empowered researchers and investors to identify specific opportunities for the development and application of these innovative technologies within the BRICS nations.

- Competitive landscape understanding: the application of Natural Language Processing (NLP) significantly enriched our understanding of the competitive landscape within the context of the BRICS, offering valuable insights into the intellectual property strategies adopted by companies and inventors. For instance, when analyzing a set of patents in the telecommunications sector, we identified consistent linguistic patterns related to patent claims. This allowed us to discern that certain companies were pursuing an intellectual property protection strategy centered around specific technical aspects, while others placed a stronger emphasis on the commercial applications of their innovations. Moreover, when examining the detailed patent descriptions, we were able to identify nuances in development approaches and the areas of expertise of companies and inventors. For example, some companies focused on hardware innovations, while others explored software solutions. The analysis of specific legal terminology was also enlightening, as it enabled us to understand the tactics employed to safeguard and fortify patent claims.
- Identifying promising innovations: through the analysis of patent documents with optimistic descriptions within the BRICS context, we successfully identified highly promising innovations. For instance, while examining patents related to renewable energy, we pinpointed patents with descriptions expressing confidence in the technical and commercial viability of their technologies. This discovery proved immensely valuable for decision-makers, enabling them to allocate resources more effectively towards Research and Development (R&D) projects demonstrating greater potential for success. Furthermore, these promising innovations also became a valuable source for identifying investment opportunities for companies and investors interested in participating in the growth and development of these technologies.
- Insights to Intellectual Property strategy: the insights obtained from Natural Language Processing (NLP) analysis enriched the understanding of the innovation landscape within the BRICS context and directly influenced the formulation of intellectual property strategies. A notable example was making informed decisions regarding patent acquisition. For instance, by using NLP analysis to assess the content of biotechnology-related patents, we identified that a specific set of patents contained key information about a novel gene-editing technique. This in-depth analysis enabled a BRICS pharmaceutical company to strategically decide to acquire these patents to strengthen their intellectual property portfolio in the field of gene therapy. Furthermore, insights derived from NLP analysis were valuable in intellectual property portfolio management. For example, when analyzing patents related to clean energy technologies, a company's R&D team identified specific areas where there was a gap in their patent portfolio. Based on these findings, they could prioritize the development of new patents to fill these gaps and bolster their competitive position. These insights also played a crucial role in managing intellectual property-related risks. By analyzing patent content and assessing their potential for litigation, companies could take initiative-taking measures to avoid infringing third-party intellectual property or defend against potential legal actions.

Natural Language Processing (NLP) techniques have proven to be indispensable for patent analysis within the BRICS context. Through sentiment analysis, we unveiled inventors' attitudes toward their innovations, identifying promising opportunities based on the emotional tone of the documents. Furthermore, NLP enabled us to explore contextual nuances and subtle details regarding the application of patented technologies, enriching our understanding of the competitive landscape and facilitating the formulation of intellectual property strategies.

4.3.4. Predictive modeling

Machine Learning models, including the Random Forest algorithm, were employed to make predictions concerning patent trends, patent grant probabilities, and the economic assessment of intellectual property within the BRICS nations. These models enable projections regarding the future of technological innovation based on historical data and current trends. The Random Forest algorithm, renowned for its ability to handle large and complex datasets, particularly excelled in this task, contributing to the generation of insights crucial for understanding the innovation dynamics within the BRICS context. The main results follow in **Table 10**.

Table 10. Predictive modeling main results.

Category	Main results
Predicting Patent Trends	Machine Learning models were trained on extensive patent datasets from BRICS countries, containing detailed information such as grant dates, technological areas, litigation history, and temporal evolution. The application of these models revealed intriguing predictions regarding emerging trends in technology. For instance, upon analyzing the data, the model identified that the fields of artificial intelligence and blockchain technology were experiencing significant growth in the number of granted patents over the years. Furthermore, the model highlighted a potential convergence of these two technologies in interdisciplinary projects. These findings provided valuable insights for companies and investors interested in anticipating the areas of technological innovation likely to gain prominence in the future, enabling them to allocate resources and research efforts more strategically and proactively.
Estimating Patent Grant Probabilities	The Machine Learning models employed generated estimates of patent grant probabilities in BRICS countries. These models considered a series of factors, such as document quality, relevant area, and the historical records of similar patent applications. This provided a precise insight into the likelihood of success for specific patent applications in BRICS countries, assisting inventors and companies in making informed decisions about where to allocate resources for R&D and intellectual property strategies. For instance, an inventor in Brazil seeking to patent a specific innovation in the clean energy technology field obtained an accurate estimate of their chances of success based on the specifics of the Brazilian patent system and the characteristics of their application. Similarly, an inventor in India seeking to patent an innovation in the software technology field received a precise estimate of the grant likelihood based on Indian patent regulations and the characteristics of their application in that country. Thus, results involving this customized approach became essential in addressing the unique aspects of each BRICS country, making patent-related decisions more informed and strategic in this context.
Evaluating Intellectual Property Value	Patent analysis in BRICS countries, supported by Machine Learning models, played a crucial role in the economic valuation of intellectual property. These models enabled the precise quantification of the economic value of patents based on a variety of criteria tailored to the realities of each BRICS country, such as technological relevance, market potential, and existing competition. For example, in the case of China, the models considered the size of the Chinese market, competitive dynamics, and the strategic importance of certain technologies to the country's economy. Based on these factors, it was possible to accurately assess the economic value of patents related to artificial intelligence in a rapidly growing Chinese context. This evaluation played a vital role in the strategic management of patent portfolios, allowing companies to identify valuable assets and prioritize investments in R&D and intellectual property strategies. Furthermore, it was recognized that the economic evaluation of intellectual property was also crucial for licensing transactions and strategic investments. For instance, a Brazilian company holding patents related to bioenergy can leverage this analysis to attract international investors interested in collaborations or licensing of these technologies.

Table 10. (Continued).

Category	Main results
Insights for Future Technological Innovation	One noteworthy outcome involved the prediction of a substantial increase in the granting of patents related to renewable energies. Based on historical analyses and current market trends, the findings anticipate that renewable energies are emerging as a growing area of technological innovation in the BRICS countries. This forecast informed strategic decisions regarding investment and technology development in this sector. For instance, a solar energy company in Brazil can utilize these insights to expand its R&D efforts in solar energy storage technologies while simultaneously seeking strategic partnerships with local companies and research institutions. This initiative-taking approach enables the company to position itself at the forefront of solar energy innovation in the country. Additionally, governments and funding agencies also benefit from these predictions by directing resources toward areas of technological innovation identified as promising.

Source: Research data analysis.

The outcomes yielded by applying Machine Learning Models to forecast patent trends, patent grant probabilities, and evaluate the economic aspects of intellectual property signify a crucial scientific approach. This approach allows us to understand the present and foresee the future of technological innovation.

These models constitute a powerful tool for supporting strategic decision-making in research, development, and intellectual property management.

4.4. Confirmation of results

In the fourth stage, our primary objective was to validate and ensure the integrity of the results we had acquired. This validation process encompassed several crucial areas of focus: first, the application of cross-validation techniques to assess the reliability of our models; second, a comprehensive comparison of our algorithms against well-established patent classification systems to gauge their performance; third, an examination of our outcomes against datasets meticulously labeled by domain experts for benchmarking purposes; and lastly, the enforcement of rigorous testing and validation protocols aimed at mitigating potential biases and reinforcing the overall robustness of our findings. This section underscores the pivotal practical results that have been rigorously verified.

Firstly, cross-validation procedures unequivocally affirmed the robustness and generalizability of our Machine Learning models, even when confronted with dynamic technological variations and regulatory disparities across the BRICS nations. This validation was imperative to guarantee the versatility of our models in adapting to distinct technological and cultural contexts within the region. For instance, during our analysis of patents associated with renewable energy in China, we adeptly discerned innovation trends in solar technologies. The rigorous cross-validation processes were pivotal in ascertaining whether these discerned trends could be extrapolated to other technological domains, such as biotechnology in Brazil or artificial intelligence in India. This, in turn, fortified the veracity and dependability of our conclusions across the entire spectrum of innovation within the BRICS.

In addition, we have conducted a comprehensive comparison of our findings with established patent classification systems, which has served as a resounding validation of the efficacy of the algorithms painstakingly developed. This validation extends well beyond mere accuracy and encompasses our algorithms capacity to contextualize our analyses within the intricate technological landscape of the BRICS nations. This reinforcement is particularly noteworthy within the intensely

competitive sphere of technological innovation in the region.

During our examination of nanotechnology-related patents in Russia, our models demonstrated heightened precision levels. Additionally, they identified nuanced subcategories of nanomaterials that had previously been insufficiently categorized within existing classification systems. This explicit demonstration underscores the adaptability of our algorithms to the distinctive technological characteristics of the BRICS, facilitating the provision of nuanced insights. An example of a nuanced insight that can be derived from this analysis is the identification of subcategories of nanomaterials that are particularly relevant to Russia's research and development priorities. For instance, our models may have identified a subcategory of nanomaterials highly conducive to electricity, which could be of great interest to companies involved in the development of new batteries or fuel cells.

Another subcategory identified by our models may comprise nanomaterials highly resistant to diseases and fungi, which could be of significant interest to companies focused on developing new materials for medical applications. The discovery of these new subcategories of nanomaterials has the potential to guide Russian researchers and developers toward areas with a higher potential for impact. Furthermore, the identification of these new subcategories may assist the Russian government in formulating policies and programs to support nanotechnology research and development.

Moreover, our research findings underwent comparison with datasets curated and labeled by domain experts. This rigorous scrutiny served as a robust validation of the precision exhibited by our Machine Learning models in terms of classifications and predictions. Importantly, this comparison acknowledged the complex technological nuances and specific regulatory frameworks existing within the BRICS countries, ensuring the faithful reflection of cultural and technological intricacies prevalent within the region. For instance, in our analysis of biotechnology-related patents in India, our findings consistently aligned with the categorizations made by local biotechnology academics some recent publications about this matter [18–20]. This validation process, therefore, buttressed the credibility of our analyses and affirmed the aptitude of our models in deciphering the intricate technological subtleties inherent to the BRICS.

Finally, the implementation of testing and validation protocols fortified the integrity and reliability of our research results, all while conscientiously considering the diversity characterizing the BRICS nations. These protocols, underpinned by rigorous data and methodological standards, constitute the bedrock for steering policies, investments, and innovation strategies in a region progressively asserting its prominence within the global innovation and economic development landscape. For example, during our analysis of patents associated with artificial intelligence in China, our protocols ensured the judicious selection of training and test data to avert overfitting, which could potentially distort our predictions. This approach produced Machine Learning models that generated accurate and reliable forecasts regarding the trajectory of AI in China—an essential focus within the nation's innovation strategies.

5. Final considerations

This study examined the strategic integration of artificial intelligence (AI) and patent research, revealing its transformative potential to accelerate innovation within BRICS economies. Our findings demonstrated how AI empowered patent analysis, unlocking valuable insights, expediting technological research, and ultimately driving economic growth.

This research offered two key contributions. Firstly, it equipped policymakers, innovators, and researchers with valuable insights to leverage patents as innovation drivers. We highlighted areas of strength, emerging trends, and collaboration opportunities, laying a foundation for informed strategies. Secondly, this study filled a critical academic gap by examining the specific role of AI in BRICS' patent analysis landscape.

Moving forward, our research agenda delved deeper into the implications of AI integration. We explored ethical and legal considerations like data privacy, intellectual property rights, and algorithmic fairness. Additionally, we aimed to assess AI's impact on patent quality, examination processes, and its potential to streamline procedures and enhance efficiency. Furthermore, we investigated how AI could facilitate collaborative innovation and knowledge sharing through dedicated platforms and tools.

Our agenda also examined the implications of AI adoption for workforce development and capacity building, identifying skill requirements and strategies for education and training. Moreover, we evaluated existing policy and regulatory frameworks to address ethical, legal, and socio-economic concerns. By addressing these research priorities, we aimed to contribute valuable insights for the responsible and effective use of AI in patent research, fostering sustainable innovation and economic development within BRICS nations.

While this study shed light on this powerful approach, limitations deserve consideration. Firstly, despite AI's transformative potential, challenges associated with data quality and availability persisted. Variations in data quality across jurisdictions and inaccuracies within patent documents could impact the reliability of AI-driven analyses. Additionally, the inherent complexity of AI algorithms introduced the risk of algorithmic biases, potentially skewing results. Finally, this study primarily focused on the technological aspects, overlooking broader socio-economic factors influencing innovation dynamics within BRICS nations.

Future research endeavors should address these limitations by adopting a comprehensive approach that considers technological factors alongside legal, ethical, and socio-economic dimensions more thoroughly. Furthermore, given the rapid evolution of AI technologies and patent systems, ongoing efforts to monitor and adapt AI-based methodologies to evolving contexts were essential to ensure the continued relevance and efficacy of innovation strategies in the BRICS economies.

As the technological landscape continued growth, collaborative research among research institutions, the private sector, and regulatory bodies became even more critical. We hoped this study inspired further investigations and concrete actions to promote the responsible use of AI in patent research, fostering sustainable innovation and propelling the BRICS economies forward.

Author contributions: Conceptualization, CZ and JLP; methodology, CZ; software, CZ; validation, CZ, CSP and RCR; formal analysis, CZ and JLP; investigation, CZ; resources, JLP and CSP; data curation, CZ; writing—original draft preparation, CZ; writing—review and editing, CZ and JLP; visualization, CSP and RCR; supervision, CZ; project administration, CZ; funding acquisition, CZ. All authors have read and agreed to the published version of the manuscript.

Conflict of interest: The authors declare no conflict of interest.

References

1. BRICS. Joint Declaration of the BRICS Summit. Johannesburg. Available online: <https://brics2023.gov.za/wp-content/uploads/2023/08/Jhb-II-Declaration-24-August-2023-1.pdf> (accessed on 13 February 2024).
2. Arora A, Ceccagnoli M, Gambardella A. Patents and the Evolution of Technology: A Review of the Literature. *Research Policy*. 2018; 47(10): 1923–1937.
3. Hall BH, Jaffe AB, Trajtenberg M. The NBER patent citation data file: Lessons, insights, and methodological tools. National Bureau of Economic Research; 2005.
4. Lanjouw JO, Schankerman M. Patent quality and research productivity: Measuring innovation with multiple indicators. *The American Economic Review*. 2004; 94(1): 281–290.
5. Hu D, Huang J, Tang Y. *Technological Forecasting & Social Change*. 2017; 123(Suppl. C).
6. Chaminade C, Shadlen KC. The governance of innovation and socio-economic development in Brazil: Challenges of building capabilities in latecomer firms. Edward Elgar Publishing; 2012.
7. Li X, Zeng X, Liu H. Intellectual property rights and innovation in developing countries. *Journal of Intellectual Property Rights*. 2020; 25(1).
8. Trajtenberg M. A penny for your quotes: Patent citations and the value of innovations. *The Rand Journal of Economics*. 1990; 21(1).
9. Dutrénit G, Capdevielle-mougribas V, Muldur U. The role of intellectual property rights in technology transfer and innovation: Evidence from China. *Research Policy*. 2018; 47(8): 1506–1517.
10. Lall S, Wahba J. Foreign direct investment in developing countries: a selective survey. *World Development*. 2004; 32(3): 199–219.
11. Dutta S, Lanvin B. The Global Innovation Index 2019: Creating healthy lives—The future of medical innovation. Cornell University, INSEAD, and World Intellectual Property Organization (WIPO); 2019.
12. Smith J, Brown A. Enhancing innovation policies: A comparative study of BRICS countries. *Innovation and Development*. 2021; 9(3).
13. Zhang H, Kim Y. Technological development and economic growth in BRICS nations: A patent analysis. *Journal of Asian Economics*. 2020; 69.
14. Chen J, Wang L, Wang W. Artificial Intelligence and Its Impact on Patent Analysis: A Bibliometric Analysis. *Journal of the Association for Information Science and Technology*. 2021; 72(11): 3033–3048.
15. Wang H, Du X, Zhang Z. Using machine learning to predict emerging technologies. *Technological Forecasting & Social Change*. 2020; 157.
16. Feldman M, Yang S. Artificial intelligence and the end of work: The role of policy in an era of technological change. NBER Working Paper Series. 2018; 25492.
17. Zhang Y, Ma L, Huang Y, Zhu Z. A bibliometric and textual analysis of the academic literature on patents. *Journal of Informetrics*. 2018; 12(1).
18. Gupta AK, Mishra SK, Singh AK. Analysis of patent trends in biotechnology in India: A review of literature. *Journal of Biotechnology*. 2021; 265: 127531.
19. Kumar A, Kumar R. Patenting trends in biotechnology in India: A bibliometric analysis. *Scientometrics*. 2022; 125(3): 2059–2080.
20. Mishra SK, Gupta AK, Singh AK. Patenting trends in biotechnology in India: A big data analysis. *Patents*. 2023; 13(2).

Article

Clustering data analytics of urban land use for change detection

C. Rajabhushanam

School of Computing, Computer Science Engineering, Bharath Institute of Science and Technology, Bharath Institute of Higher Education and Research, Chennai 600073, India; Rajabhushanamc.cse@bharathuniv.ac.in

CITATION

Rajabhushanam C. Pre-Clustering data analytics of urban land use for change detection. *Computing and Artificial Intelligence*. 2024; 2(2): 570.
<https://doi.org/10.59400/cai.v2i2.570>

ARTICLE INFO

Received: 23 April 2024
Accepted: 20 June 2024
Available online: 2 July 2024

COPYRIGHT



Copyright © 2024 by author(s).
Computing and Artificial Intelligence is published by Academic Publishing Pte. Ltd. This work is licensed under the Creative Commons Attribution (CC BY) license.
<https://creativecommons.org/licenses/by/4.0/>

Abstract: In this study, the author proposes and details a workflow for the spatial-temporal demarcation of urban areal features in 8 cities of Tamilnadu, India. During the inception phase, functional requirements and non-functional parameters are analyzed and designed, within a suitable pixel area and object-oriented derived paradigm. Land use categories are defined from OpenStreetMap (OSM) related works with the scope of conducting climate change, using multispectral sensors onboard Landsat series. Furthermore, we augment the bands dataset with Spatially Invariant Feature Transform (SIFT), Normalized Difference Vegetation Index (NDVI), Normalized Difference Built-Up Index (NDBI), Leaf Area Index (LAI), and Texture base indices, as a means of spatially integrating auto-covariance to stationarity patterns. In doing so, change detection can be pursued by scaling up the segmentation of regional/zonal boundaries in a multi-dimensional environment, with the aid of Wide Area Networks (WAN) cluster computers such as the BOWULF/Google Earth Engine clusters. GeoAnalytical measures are analyzed in the design of local and zonal spatial models (GRID, RASTER, DEM, IMAGE COLLECTION). Finally, multi variate geostatistical works are ensued for precision and recall in predictive data analytics. The author proposes reusing machine learning tools (filtering by attribute-based indexing in PaaS clouds) for pattern recognition and visualization of features and feature collection.

Keywords: distributed computing; HPC; multi-spectral imagery; machine learning; AI; local climate change; zonal analysis; spatial data model

1. Introduction

Using satellite imagery in Google Earth Engine (GEE) such as the omnipresent Landsat 8–9 optical sensors, over a time period of 12 years (2012–2024), with the intent of evaluating geo-based change detection, is the main forte in this research endeavour.

Since urban areas such as Chennai (roughly 5904 sq.km) have depicted refactoring of land use and land cover changes primarily through trend analysis and suitability analysis, we intend to apply geospatial data analytics as an inference engine for effective spatial decision support.

In Tamilnadu, India, urbanisation has been an evolving trend since medieval times to present and this phenomenon can be observed in most town expansion plans, as outlined by the 2011 census. Measuring population growth in Tamilnadu can be attributed to rapid industrialisation in urban areas and the consolidation of cooperative unions from villages to towns.

By generating pre clusters within a region, by unsupervised classification technique, we can assimilate homogeneous feature maps into a class based hierarchical clustering method. This study differs from related literature reviews, in that an instance is defined within a contextual purview, for processing each land-use feature by

deducing in a neighbourhood distance-based classifier.

2. Objective

Demarcating zonal boundary or distinct urban gradient break is the primary forte of optical remote sensing sensors during the present day. In this study we intend to demonstrate the findings from one such controlled experiment, which is to provide to government authorities a spatial decision support system (SDSS).

Primary focus is in geospatial bigdata analytics and its application to cadastral mapping in urban zones of Tamilnadu. Specifically, we will address the following:

- Implementing a base layer for Tamilnadu state with pre-determined feature classes.
- Visualization of topology conditions for any given pixel in the base georeferenced data.
- Overlay of raster feature datasets with vector feature classes for regular/irregular spatial delineations.
- Performing spatial queries with selection, projection, and joins to extract a geospatial feature map for the geoanalytical engine using Google Earth Engine (GEE).
- Implementing a spatial-temporal (from 2012 to 2024) inference model using machine learning and Platform as A Service (PaaS) cloud offering. The linear model simulates climatic and aspatial change detection techniques, within a linear multivariate geostatistical regression criterion.
- The model is trained with region of interest (ROI) and point of interest (POI) sample datasets, and then it is applied for validating feature datasets from any given feature map, using iterator operator and collection classes in JavaScript programming language.
- Data analytics is applied to operate on the following geospatial change detection in LANDSAT-8, LANDSAT-9 sensor's-based satellite imagery scenes.
- Trend analysis with focal and zonal areal interpolation boundaries.
- Suitability analysis with environmental and human induced factors.
- Unsupervised classification study is then under-taken to do density clustering of data points, resulting in user specific intervals for each user-specific class. Deriving overall accuracy, specificity, and recall measures from a classification study. Specifically, a confusion matrix (commission and omission errors) is computed.
- Furthermore, Spatially Invariant Feature Transform (SIFT) method is applied to the feature map and the classification study is ensued for the selected multi-scale scene analysis. True positives and False negatives are computed and the root mean square error (RMSE) is conducted for validation and verification.
- Finally, Quantile regression analysis is plotted using Mean-Covariance shift measure as a planar separator to fit to the datapoints. Slope and Intercept are calculated to determine the inflection point (feature) where collinearity is evident.

3. Related study

Shan and Sampath [1] proposed a method based on a region-growing algorithm

to group similar points into the same building by iteratively collecting points within a moving window. In the study of Sampath and Shan [2], the fuzzy c-means (fuzzy k-means) algorithm was used to cluster individual planar roof segments to reconstruct building models. As the k-means and fuzzy c-means algorithms can only cluster convex shapes, they are rarely utilized directly to segment individual building derived from image segmentation [3–6].

The task of retrieving parameters from observations/predictions is most commonly associated with deriving measured output. With this research proposal, we aim in creating an inverse radiometric calibration machine learning model. This model will map causal inference from live sensor readings to observed bio-geophysical readings. This approach varies from training models in Artificial Intelligence and Machine learning domains [7]. Further, the feature extraction task can be based on multi-dimensional sensor measurements or the target parameter can consist of several variables. Thus, this study is categorized as a classic Classification and Regression Theory (CART) problem.

Our hypothesis is that linear classification technique produces spatially invariant features within a zonal areal modifiable area unit problem (MAUP). In doing so, adapting to multiple perceptron neural networks, transforms inputs to expected outputs by iteratively clustering over bounded feature maps (DBSCAN-VAE). So, as a result non-linear methods such as NN, can be deployed in a generative network (hyper-parameters) in a linear mode. The use of neural autoencoder allowed us to transform raw images in their feature representation versions. These encoded images contain all the essential information about the original. Among different neural network models, autoencoders have found the application in many domains. In image processing, the autoencoders are used mostly for feature extraction, image segmentation, image compression/data reprojecting and image reconstruction. The advantage of our autoencoder is that it can be used as a standalone algorithm for feature extraction with further use of extracted features for unsupervised purposes.

As this study is a confluence of computer-vision, remote sensing, machine learning, and deep learning fields, a select number of publications have explored the research criterion. By spatially interpolating with deep learning approaches, such as spatial-temporal VAEs, allow learning the behaviour of simulations from previous runs. Ideally, this results in speeding up climate models by regularization of latent factors. While there will be a penalty on prediction accuracy, the speed-up might be worth in some cases. In addition, the same unsupervised autoencoder-like models could be used to approximate the sub grid processes. Such architectures allow for the stochastic modelling of the earth system from raw observations. Currently, no use of this approach for the data-driven modelling of earth dynamics has been published.

4. Synthesis of feature maps by using a variational autoencoder (VAE)

The idea behind a VAE is different from the idea of an Autoencoder. Instead of reproducing single inputs, the goal is to map the unknown distribution of the inputs onto a d-dimensional multivariate Gaussian distribution and then denormalize to original data (**Figure 1**). The representation of the latent space using a dimensionality

reduction method is also explored in the study of Lu et al. [8]. The goal was to detect and track extreme events, using an AE as a method to extract features.

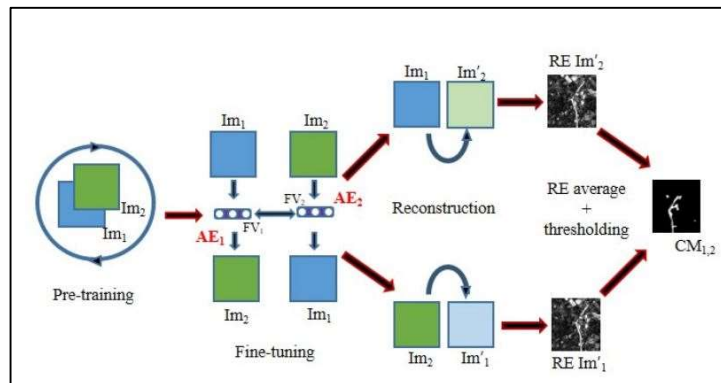


Figure 1. Schematics of an auto-encoder in image synthesis [9].

Adapting the effectiveness of autoencoders to such a study is particularly important, as the final output is a generative model that closely approximates the latent space with second order derivative (dy/dx). Autoencoders can better represent complex spatial-temporal nature of the underlying datasets that is currently often summarized with indices or convoluted features. In particular there are usually no labels making unsupervised autoencoders useful. The motivation for AEs comes from: 1) the unsupervised extraction of useful features to predict; 2) their ability to generate additive (unseen data) samples from the original distribution; 3) the removal of noise, and 4) the identification of anomalies.

In the worst case, such as without the use of an accelerating index structure (applying a sequential search), or on a degenerated dataset (e.g., all points within a distance less than e), the overall runtime complexity remains $O(n^2)$. Furthermore, the time complexity for creating spatial index structures (e.g., GiST or R*-tree) is usually in $O(n \log n)$ and, as a result, this does not affect the total time complexity. Although some border points may be firstly marked as outliers and be clustered later, they are not added to the seeds list and unnecessary neighbourhood retrieval avoided because they are not core points. On the other hand, a border point may be shared by two closely distributed clusters. In this case, the algorithm will group the shared point into the first discovered cluster. An improvement on this situation is to group it into the cluster that its nearest core point belongs. Except from these cases, the result is insensitive to the processing order of the points. Additionally, if treating border points as outliers, the algorithm is deterministic.

Previous studies on land use mostly focused on the functional zoning of ecological landscapes or farmland, with limitations to other implementations [10]. As mentioned above, OSM data contain multiple and finer information on land use, especially social functional information that is rarely captured from remote sensing images. This study pays attention to urban social functional land use and tries to extract three social functional types of land use based on OSM data. The OSM database contains 21 features describing the land surface objects, and each feature is labelled by several attributes. According to the mapping rules of **Table 1**, this study chooses three land use function types as follows: 1) residential land use, which indicates the

land parcels that are labelled ‘residential’, ‘apartments’, ‘dormitory’, ‘garage’, ‘house’, ‘residential’, ‘dorm’, etc.; 2) commercial land use, which refers to the land for wholesale and retail, accommodation and catering, and commercial or financial purposes with labels of ‘retail’. ‘Marketplace,’ ‘pharmacy’, ‘hotel’, ‘café’, ‘restaurant’, ‘bank’, ‘commercial’, etc.; 3) public service land use, which includes government and organizations, science and education, medical and health charity, recreational land, public facilities, park and green space, and scenic facilities with labels of ‘office’, ‘university’, ‘hospital’, ‘cinema’, ‘post office’, ‘park’, ‘viewpoint’.

Table 1. Confluence of image analysis techniques.

S.NO	Scientific domain	Technology	Roadmap
1	Distributed Computing	Parallel and Distributed High Performance Computing System (HPCS).	High level design (HLD) to integrate satellite imagery retrieval, storage, and analysis into a value-added reseller (VAR) functional outlay.
2	Computer Vision	Algorithm Development for Scheduling and Staging System	Image analysis, Digital Image processing for pixel detection and synthesis. Includes 3-D (Stereo vision) and 2-D feature extraction.
3	Satellite Communications	Image Retrieval and Archival System	Building Block for Ingest and Analysis workflows.
4	Geospatial Analytics	Feature Engineering	Geoprocessing, Recomputing, Geoanalytical, and Geo-Statistics.
5	Cloud Computing	Client/Server n-Tier Architecture in PaaS and SaaS.	NIST defined a computer based functional entity with Ubiquity and Pervasiveness.
6	Remote Sensing	Imaging System for Sensing and Remote Measurement of natural phenomena	Radiometric measurement and feature extraction
7	Software Engineering	JavaScript and Python software development using Google Earth Engine (GEE) as an Interface (ICD standards)	Modular driven architecture with call-return graphs and tight coupling using high cohesiveness.
8	Machine Learning	Unsupervised Classifier, Variational Autoencoder, SGD and Cross-Entropy feature segmentation.	Semi-structured data analysis using unsupervised learning within a pre-defined number of feature classes.
9	Artificial Intelligence	Inference System with Explainable Artificial Intelligence (XAI)	Deep-learning using feedback and reward inference engine.
10	Database Technologies	SQL data store for benchmarking and Indexing purpose.	Data store in the cloud and within the firewall premises for maintaining redundancy and replication ability.

5. Proposed methodology

Deep learning also faces challenges that are unique to earth science data: multimodality; high degree of heterogeneity in space and time; and the fact that earth science data can only provide an incomplete and noisy view of the underlying eco-geo-physical processes that are interacting and unfolding at different spatial and temporal scales.

Causal reasoning with an inference engine is limited to studies for one site or one instance by design. To generalize the results obtained as in a sequential work flow, we have to aggregate the effects from each normalized input as a feedback mechanism to the precedent and antecedent path. Then applying heuristic rules (knowledge base aided by a human in the loop) we can reconstruct any input sequence averaging over mean (zero) and variance (absolute one) [10,11].

In this paper, a density-based cluster algorithm, DBSCAN [12], is independently implemented in a spatial database to separate building point clouds into individual

buildings. As a density-based clustering method, DBSCAN characterizes a well-defined “density reachability” cluster model by connecting points that satisfy a density criterion as defined as a minimum number of objects within a certain neighbourhood.

In contrast to k-means method that can only find convex clusters, DBSCAN can form a cluster of an arbitrary shape. Moreover, based on spatial databases, it benefits from the spatial index offered by the system and achieves performance improvements. The algorithm is implemented as follows, and the flowchart is shown in **Figure 2**.

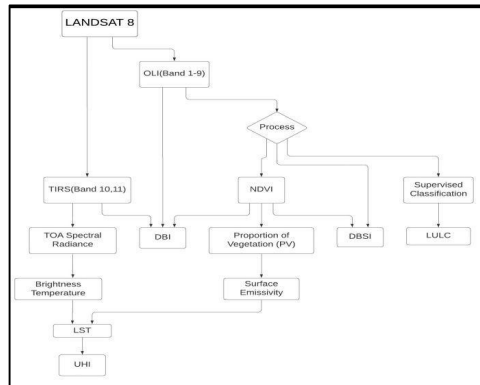


Figure 2. Proposed methodology. NDVI (Normalized difference vegetation index; Dry Built-up index (DBI); Dry bare-soil index (DBSI).

6. DBSCAN algorithm

We propose a variant of the existing DBSCAN algorithm for these purposes. The algorithm efficiency is measured by the relative weights applied to form clusters from a randomly assigned seed value. In addition to the distance (radii) parameter the pre-processing parameter, number of clusters (region of interest—ROI), is defined in a linear approach. For each assigned point in a cluster, intra-class variance is minimized or, we can emphasize that connected component segments are created within a 5×5 moving pixel neighbourhood. Thus, each individual data point in dense regions are core points that are assigned to each cluster based on the specificity of the sparse density matrix.

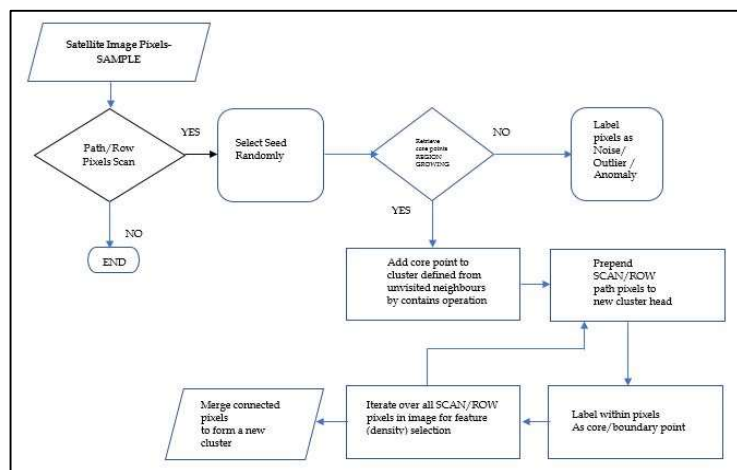


Figure 3. Process data flow diagram in pixel-based approach.

One example of a bio-geophysical phenomena is surface/atmospheric temperature (Band 10—TIR band of LANDSAT-8) recorded from the satellite, is input to the climate forecasting model. This topic is sparsely researched in the remote sensing literature, and hence we emphasize the value factor derived by conducting a continuous multi-variate as inputs to our study. Many other input variables can be listed—NDVI, DBI, DBSI, etc. (Figure 3).

In this experiment pre-processing machine learning model VGG₁₆ is adapted (transfer learning) for filtering patterns during the training. During this iteration step input layer, hidden layers, and output layer are connected in a sequential arrangement to form a non-singularity transfer learning entity (Figure 4).

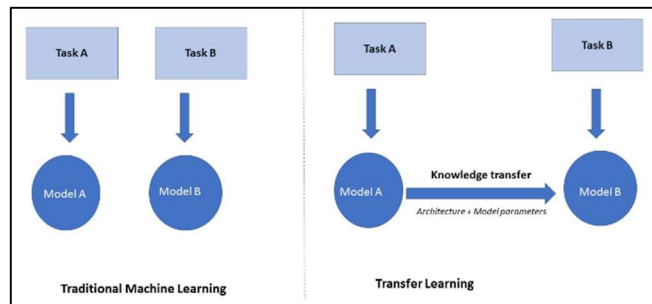


Figure 4. Type I and II error reduction.

The filters perform feature extraction and transform our data into a latent space where the problem is potentially easier to solve. For VGG₁₆, data is typically structured in 2-D space (Longitude, Latitude). We can do data augmentation by inputting the sequence of bands (channels) and Indices from satellite data into the feature selection model (Figure 5). This architecture is akin to a convolutional neural net, excepting that it behaves in a VGG₁₆ pre-trained network as such. For our purposes, we shall denote this step as Model-A. Inputs to the variational autoencoder (VAE) are denoted as Model-B.

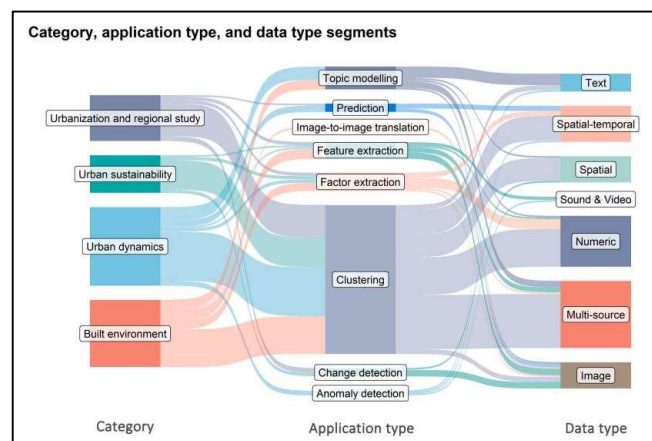


Figure 5. Category, application type, and data type segments [13].

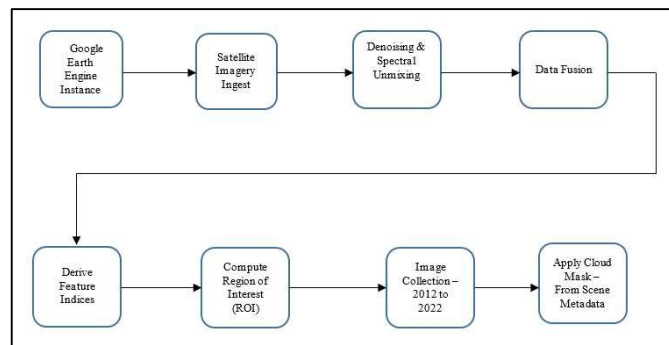


Figure 6. Pre-processing of input imagery. Examines an instance of ROI.

Therefore, we can adapt a unified representation with $x = \text{height}$, $y = \text{width}$, $b = \text{band}$, $k = \text{statistical indices}$, and $T = \text{Temporal identity matrix}$ (Figure 6). This would result in cubic filters with T channels. Two ways of modelling spatial, spectral and temporal relationships. (a) Cubic convolutions over space and frequencies and stacking the n time steps. (b) Stacking frequency, B , and time, T , and performing 2D convolutions over spatial dimensions (X, Y). Originally, the concept was developed for visual recognition problems with time-varying inputs [14]. It is important to note that when stacking channels, e.g. (B, T), the order of them is not taken into account by the model (Figure 7).

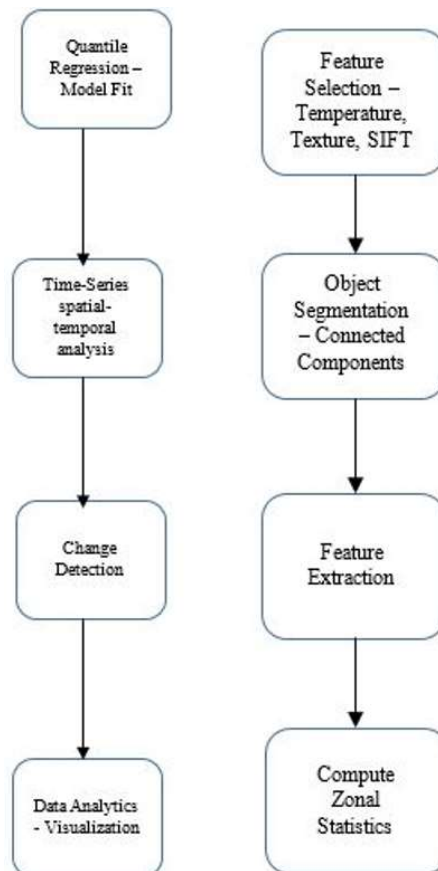


Figure 7. Post-processing flow chart.

Mean Square Error or Cross-entropy Optimization?

Retrieval (Figure 8) is most often associated with least-square regression

modelling. Traditionally, linear regressors optimized by least-square have been the preferred choice but are not always sufficient to capture the complexity of retrieval problems. A neural network optimized by the mean-square error (MSE) loss function, is a one way of extending the least square linear regression with non-linearity. Alternatively, a probability distribution over possible outcomes can be modelled with Cross-Entropy based error functions. Cross-Entropy is generally associated with problems where we wish to label data, e.g., segmentation or classification. It can be shown [15] that the MSE loss is the maximum likelihood solution to a problem where the target can take any real values.

Many challenges exist for deep learning in bio-geophysical parameter retrieval problems. Since we are modelling the Earth’s state, we need to apply algorithms on large amounts of data. Further, we have many sources of variance in our observations caused by e.g., seasonal, yearly and geographical variation (**Table 1**).

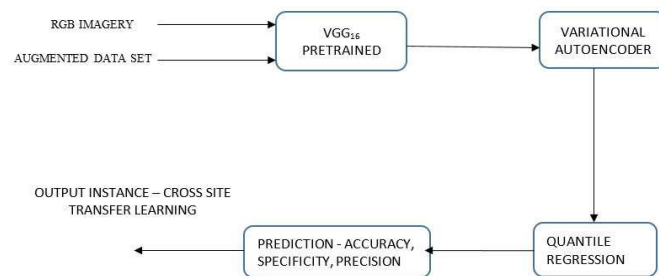


Figure 8. Illustration of a generalized network for scene understanding and refinement.

6.1. Dataset preparation

Landsat 8, 9 sensor-based satellite imagery will be ingested from Google earth engine (GEE) together in creating a spatial data collection service. Bands in near-infrared, short-wave infrared, thermal, and panchromatic will be chosen as inputs to the pre-trained network VGG₁₆.

From the chosen band spectrum, indices will be derived for first—Creating an image fusion-based collection set. Namely panchromatic sensor from Landsat-8 at 15 m resolution will be fused with 30 m nir-sw-tir bands to create a spatial dataset that will reveal distinct patterns for a region of study (**Table 2**).

Table 2. Specification of image analysis domains.

S. No	Geospatial domain	Description
1	Geo-Processing	Zonal statistics, Feature dataset, Feature understanding
2	Geo-Computing	Feature selection, Feature extraction, Attribute spatial joins, Model fitting and Feature understanding
3	Geo-Analytical	Model fitting, Validation, Predictive analytics
4	Geo-Statistics	Validation, Predictive analytics

After the image fusion step, spatial indices (**Figure 9**) calculated from each pixel in the input satellite imagery will be computed as inputs to the pre-trained VGG₁₆ network. Specifically:

- Normalized Difference Vegetation Index (NDVI) for Vegetation feature class.
- Built -up Dryland—Dry Built-Up Index (DBI).

- Barren land in dry climates—Dry Bare Soil Index (DBSI).
- Normalized Difference Vegetation Index (NDVI) for water bodies feature class.

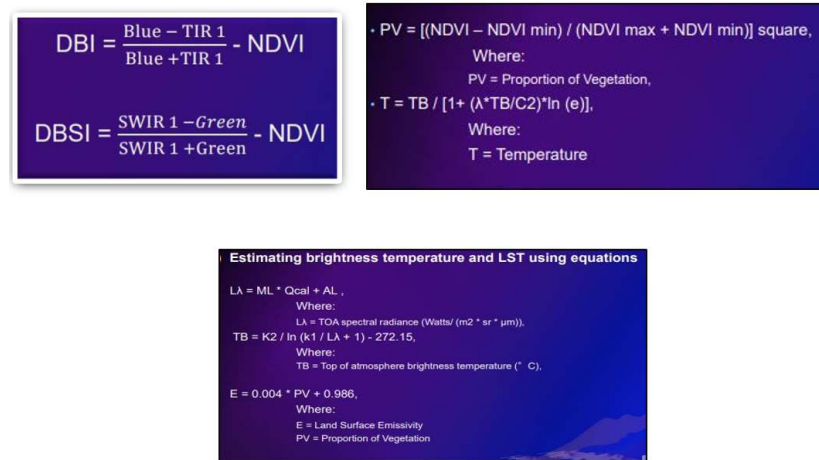


Figure 9. Spectral indices calculation using image expression.

6.2. Data augmentation in input feature dataset

- Scale Invariant Feature Transform (SIFT)—Hierarchical Scale Space (HSS).
- Texture derived from Gray Level Co-Occurrences Matrix (GLCM)-entropy.
- Temperature feature derived from Band 10 of Landsat-8,9 (Table 3).

Table 3. Summary of geo-processing tasks.

S.No	TASK—Unit Instance	Description
1	Zonal Statistics	Deriving exploratory statistics for Region of Interest (ROI)
2	Feature Selection	For pattern recognition and discriminative analysis
3	Feature Dataset	Multi-Variate inputs and Derived classes in multi-dimensional analysis.
4	Feature Extraction	Fuzzy Inference System based Segmentation followed by CART task.
5	Feature Understanding	Variational Autoencoder (VAE) as a generative adversarial network for encoding and decoding patterns that can be discernible by a human in the loop.
6	Attribute based Spatial Joins	Linking an object-relational data store with Geospatial data using spatial joins for query and search operations.
7	Model Fitting	Quantile Regression technique for Interval demarcation and polynomial expansion using a symmetrical estimation.
8	Validation and Verification (V&V)	Conducting n-fold cross validation to achieve conditional local minima for ROC characteristics.
9	Predictive Analytics	Deep Learning technique to generate specificity, accuracy, and precision for recall instance.

7. Algorithm

Algorithm 1 Regularization and hyperparameters tuning

- 1: Aim is to achieve local minima through density clustering (DBSCAN) based on Euclidean distance measure.
- 2: Design covariance matrix as X^T .
- 3: Design of seed value as a random variable (RV).
- 4: For each seed, initialize core point, boundary point and zonal area (ROI).
- 5: While region growing is not ~ to NIL for batch intervals while computing weight/density estimate.
- 6: Initialize cluster radius to be a compact shape (necessarily not rectangle, square) CIRCULAR path.

Algorithm 1 (Continued)

-
- 7: Initialize region growing pixel values that in time and space will merge all mutually exclusive points.
 - 8: If core point distance is LT maximum points, assign boundary.
 - 9: If cluster point distance GT minimum points, assign cluster region.
 - 10: Now, repeat over all pixels in the sample (~) for consideration into neighbourhood (optimal) having low intraclass variance.
 - 11: Construct covariance matrix.
 - 12: Stop when number of cluster regions reach 4 feature classes.
 - 13: For any unassigned point, assign it as an outlier/anomaly.
 - 14: Stop region growing.
 - 15: Scan all boundary points € cluster regions. If shape is defined as equidistant, assign count of core points as a threshold.
 - 16: Finally run R-Tree indexing for efficient spatial index retrieval.
-

8. Design of input parameters to the model

- Design of a sparse matrix;
- Design of covariance matrix;
- Design of features as inputs to VAE—pyramidal resolution;
- Design of spatial autocovariance measure;
- Design of input tiling scheme in imagery;
- Design of UTM zone in region of study (ROI);
- Design of thresholding value in normalization of input data;
- Design of SIGMA value as input to SIFT;
- Design of cell size during Texture calculation with GLCM method;
- Design of hidden layers in variational autoencoder (VAE);
- Design of sequence flow during input to pre-trained network (VGG₁₆);
- Design of ordering of inputs to VAE;
- Design of radius parameter for input to density clustering;
- Design of neighbourhood size (5 × 5 default);
- Design of epoch size during sequential processing in Iterator to avoid over/under fitting.

9. Conclusion

In this study, feature engineering paradigm was applied to evince interest and focus on the fragile relationships between humans and factors in environmental ecosystem. We have dealt with multi-spectral sensors and the spatial scale problem, in achieving an optimal scale-resolution under which causal inference is evident. Various state of the art algorithms and pixel/object-oriented machine learning stages were feature selected and feature extracted, in accordance to spatially distributed parallel systems.

Conflict of interest: The author declares no conflict of interest.

References

1. Shan J, Sampath A. Urban Terrain and Building Extraction from Airborne LIDAR Data. CRC Press; 2006.
2. Sampath A, Shan J. Segmentation and Reconstruction of Polyhedral Building Roofs From Aerial Lidar Point Clouds. IEEE

- Transactions on Geoscience and Remote Sensing. 2010; 48(3): 1554-1567. doi: 10.1109/tgrs.2009.2030180
3. Awrangjeb M, Fraser C. Automatic Segmentation of Raw LIDAR Data for Extraction of Building Roofs. *Remote Sensing*. 2014; 6(5): 3716-3751. doi: 10.3390/rs6053716
 4. Caihua Y, Yonghong L, Weijun Q, et al. Application of Urban Thermal Environment Monitoring Based on Remote Sensing in Beijing. *Procedia Environmental Sciences*. 2011; 11: 1424-1433. doi: 10.1016/j.proenv.2011.12.214
 5. Sohn G, Dowman I. Data fusion of high-resolution satellite imagery and LiDAR data for automatic building extraction. *ISPRS Journal of Photogrammetry and Remote Sensing*. 2007; 62(1): 43-63. doi: 10.1016/j.isprsjprs.2007.01.001
 6. Zhou H, Jiang H, Zhou G, et al. Monitoring the change of urban wetland using high spatial resolution remote sensing data. *International Journal of Remote Sensing*. 2010; 31(7): 1717-1731. doi: 10.1080/01431160902926608
 7. Tsagkatakis G, Aidini A, Fotiadou K, et al. Survey of Deep-Learning Approaches for Remote Sensing Observation Enhancement. *Sensors*. 2019; 19(18): 3929. doi: 10.3390/s19183929
 8. Lu X, Yuan Y, Zheng X. Joint Dictionary Learning for Multispectral Change Detection. *IEEE Transactions on Cybernetics*. 2017; 47(4): 884-897. doi: 10.1109/tcyb.2016.2531179
 9. Sublime J, Kalinicheva E. Automatic Post-Disaster Damage Mapping Using Deep-Learning Techniques for Change Detection: Case Study of the Tohoku Tsunami. *Remote Sensing*. 2019; 11(9): 1123. doi: 10.3390/rs11091123
 10. Schäfer M, Kröger M. Joint problem framing in sustainable land use research. *Land Use Policy*. 2016; 57: 526-539. doi: 10.1016/j.landusepol.2016.06.013
 11. Tibau XA, Reimers C, Eyring V, et al. Spatiotemporal model for benchmarking causal discovery algorithms. Published online March 23, 2020. Doi: 10.5194/egusphere-egu2020-9604
 12. Bengio Y, Deleu T, Rahaman N, et al. A Meta-Transfer Objective for Learning to Disentangle Causal Mechanisms. *ArXiv*. 2019; arXiv:1901.10912v2.
 13. Wang J, Biljecki F. Unsupervised machine learning in urban studies: A systematic review of applications. *Cities*. 2022; 129: 103925. doi: 10.1016/j.cities.2022.103925
 14. Ester M, Kriegel HP, Sander J, Xu X. A density-based algorithm for discovering clusters in large spatial databases with noise. *KDD*. 1996; 96(34): 226-231.
 15. Donahue J, Hendricks LA, Guadarrama S, et al. Long-term recurrent convolutional networks for visual recognition and description. In: *Proceedings of the 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

The based-biofeedback approach versus ECG for evaluation heart rate variability during the maximal exercise protocol among healthy individuals

Sara Pouriamehr¹, Valiollah Dabidi Roshan^{1,2,*}, Somayeh Namdar Tajari³

¹ Department of Exercise Physiology, Faculty of Sport Science, University of Mazandaran, Babolsar 47416-13534, Iran

² Athletic Performance and Health Research Center, Faculty of Sport Science, University of Mazandaran, Babolsar 47416-13534, Iran

³ Department of Motor Behavior, Faculty of Sports Science, University of Mazandaran, Babolsar 47416-13534, Iran

* **Corresponding author:** Valiollah Dabidi Roshan, v.dabidi@umz.ac.ir, vdabidiroshan@yahoo.com

CITATION

Pouriamehr S, Dabidi Roshan V, Namdar Tajari S. The based-biofeedback approach versus ECG for evaluation heart rate variability during the maximal exercise protocol among healthy individuals. *Computing and Artificial Intelligence*. 2024; 2(2): 1481. <https://doi.org/10.59400/cai.v2i2.1481>

ARTICLE INFO

Received: 27 June 2024

Accepted: 3 September 2024

Available online: 22 September 2024

COPYRIGHT



Copyright © 2024 by author(s).

Computing and Artificial Intelligence is published by Academic Publishing Pte. Ltd. This work is licensed under the Creative Commons Attribution (CC BY) license.

<https://creativecommons.org/licenses/by/4.0/>

Abstract: Although the use of biofeedback devices is beyond measure, they are widely applied only for clinical purposes. Therefore, this study evaluated whether biofeedback devices could be applied to estimate heart rate variability (HRV) among healthy populations. 60 individuals (46 ± 5 years; 30 women) performed maximal exercise protocol (MEP). At pre- and post-MEP status, HRV indexes were collected by two devices: 1) the electrocardiogram device (ECG); 2) the biofeedback device (BIO). At pre-exercise status, all HRV parameters had significant correlations, ranging from low ($r = 0.241$) to high ($r = 0.779$). At post-exercise status, significant correlations for some of the HRV measures were found as well, ranging from low (i.e., $r \leq 0.29$) to moderate (i.e., $0.3 \leq r \leq 0.49$). According to our knowledge, this study is the first attempt to evaluate HRV by biofeedback devices among healthy individuals, which shows they can also be applied as a swift method to examine HRV among healthy individuals, especially in rest conditions.

Keywords: heart rate variability; electrocardiogram; biofeedback; physical activity; healthy population

1. Introduction

As a response to any sudden physical challenges, the cardiovascular system can be modified to maintain homeostasis, which means heartbeats constantly change [1]. The heart rate variability (HRV), quantified by the fluctuations in R-wave to R-wave intervals (RRI), has constituted a useful non-invasive method to evaluate autonomic activity, particularly parasympathetic tone and sympathy-vagal balance at either rest or any physical activities [2–5]. As the cardiovascular system responds to stressors, HRV may predict certain diseases [6–8]. Plus, it can be useful to monitor high performance during training sessions [9–11]. Meanwhile, the literature on autonomic activation has explained that the reduction in HRV, consisting of both higher sympathetic and lower parasympathetic activities, can be considered a frequent marker of abnormal and insufficient autonomic nervous system (ANS) adaptation [2,12,13] and the elevation in blood pressure variability [14], which possibly indicates a low heart capacity to respond to multiple physiological and environmental stimuli [14–17], which is associated with diverse pathological conditions such as coronary heart disease and mortality [18–20], future functional decline [21], chronic heart failure [22], sarcopenia [23], and hypertension [24]. Whereas, high HRV, known as an indicator of evaluated parasympathetic and reduced sympathetic activities, illustrates a good body adaptation [25,26]. Recently, it has been reported that physical exercise, both aerobic and resistance training, influences the cardiovascular system positively, especially

vagal activity as its important determinant [27–30], which can be considered the cornerstone of nonpharmacological treatment and prevention of such diseases [31–35]. To illustrate, various exercise methods, practically aerobic, alter the cardiac-autonomic balance, including increasing vagal autonomic drive while lessening sympathetic drive [31–39]. Indeed, the studies have been conducted among healthy children [40], young adults [41], and patients [42–44] following performing the aerobic [44,45], resistance [46], or interval training interventions [47] protocols.

Besides, sports and training sciences also pay attention to either time- or frequency-domain HRV indices, which means HRV is applied as the noninvasive method to measure autonomic changes following short- and long-term endurance training among individuals performing leisure sports activities and high-performance training [48]. These changes are followed up by a notable reduction of heart rate at either rest status or during submaximal exercise conditions, reflecting the elevation activity of the autonomic efferent and shifting in favor of vagal-activity enhancement to modulate the cardiac rhythm [48]. In other words, HRV kinetics may predict aerobic fitness and exercise performance during sub- or maximal workouts [49–51], which is also known as a key marker to evaluate fatigue intensity [52] and a diagnostic marker of overreaching and overtraining [48]. Generally, the whole study literature declares the vital necessity of assessing HRV among various populations to monitor their health status and performance regardless of both the type and intensity of exercise.

Recently, although trended smartwatches (i.e., Apple Watch, Garmin, Fitbit, Polar, and Samsung Galaxy Watch) are being evaluated for HRV estimating accuracy related to stress management features [53–55], generally, devices such as clinical multi-lead ECG systems (e.g., Holter ECGs) [56], photoplethysmography (PPG) [57,58], the FarosTM ECG [59–61], the Actiheart [62], the AidlabTM [63,64], and Polar H [57,65] have been applied to assessing HRV indices regarding the large series of evidence. Despite this, over a seven-decade period, electrocardiograms (ECG) have become the most routine to monitor HRV [12,66,67]. ECG can be interchanged by either the Polar (i.e., S810i and V800) or Suunto t6 instruments to record the R-R intervals in both healthy (i.e., runners) and patient populations [68–70]. The 12-lead ECG, which is also known as the golden standard, consists of three bipolar-limb- leads (i.e., I, II, and III), three unipolar-augmented- leads (i.e., aVL, aVR, and aVF), and six unipolar chest leads, including V1–V6 [67,71]. Nevertheless, some items are crucial to measuring HRV indices by this device, such as the correct placement of each lead reported by various studies [67,71,72] and it also requires both expertise and time.

On the other hand, HRV-biofeedback (HRV-bio) devices impact clinical therapeutics in various diseases [73]. Regarding some evidence, HRV-bio is known as an effective non-pharmacological intervention to monitor autonomic balance [12,74] which has skin conductance that can be applied as direct quantitative ANS markers [75], expressing its potential value in chronic disease management [76]. In general, HRV-bio has been applied as a training method to enhance sports and workout performance [77–79]. For instance, the HRV-bio has been considered a technique for managing stress based on longer exhalations and slower respiration training [80,81], and based on our knowledge, only a few papers have used biofeedback devices as an HRV measurement method [79]. Thus, applying HRV-bio would be considered another option for measuring HRV among healthy populations.

Taken together, HRV could be estimated by calculating the R-R intervals through various devices. Although the 12-lead ECG is the standard method to measure the HRV, it requires special items (i.e., expertise and time). On the other hand, it has been expressed that HRV measurements generated by a 12-lead ECG, a Holter-style ambulatory recording system, and a custom-built chest strap (strap) would not agree well in all cases [45]. Despite this, to the best of our knowledge, applying HRV-bio has been overlooked as a real-time and swift-measurable method to measure HRV among healthy individuals, which can be considered crucial and required for monitoring HRV in healthy populations regardless of both the type and intensity of exercise. Therefore, we sought to assess the accuracy between ECG (as the golden method) and BIO, which means approaching HRV by biofeedback device would be authentic at rest status (pre-exercise condition), and whether this situation would remain the same after performing a maximal exercise protocol (MEP) (at post-exercise status) among healthy individuals.

2. Materials and methods

2.1. Ethical approval

In this investigation, the local institutional ethics committee reviewed and approved all the methods and data collection (Ethical Code: IR.UMZ.REC.1397.019). It should be mentioned that the whole research process was performed according to the 1964 Helsinki Declaration [82]. In this regard, all healthy males and females had the opportunity to participate and obtain informed consent. In addition, the testing procedures, protocols, and equipment were introduced to participants, making them familiar with the research process. Meanwhile, the opportunity was provided for each individual to query any progress section whenever it was not comprehensible. Essentially, leaving and/or withdrawing the study progress without any consequences was the individuals' right when they did not want to keep on participating.

2.2. Study design

In this study, the HRV simultaneously was recorded during pre- and post-exercise status to estimate the correlation between the HRV indices extracted by two measurement devices, i.e., the electrocardiogram (ECG) and the biofeedback (BIO).

2.3. Participants, inclusion, and exclusion criteria

In this study, 60 healthy, qualified-volunteered individuals (30 females) participated. Additionally, we kindly asked participants to avoid strenuous exercise and to abstain from any food and beverages containing alcohol and caffeine 48 h before data collection. All procedures and measurements were conducted from 8:00 to 13:00.

In addition, to be eligible to remain in the investigation process, some existing requirements were seated, such as: 1) Not having a smoking habit and/or being exposed to second-hand smoke, 2) No consumption of any antioxidant supplements at least one month before the study, 3) No history of chronic cardiovascular events or pulmonary and inflammatory diseases; 4) Not having any other medical limitations,

including any physical disabilities and/or limitations of mobility. All study females were screened for the inclusion eligibility criteria. The survey included questions about the history of the menstrual period (present, irregular, or absent).

2.4. Anthropometric measurements

Before the exercise protocol, a specialized expert assessed participants' anthropometric characteristics [83]. In this case, a stadiometer was used to measure each individual's weight and height according to its height (about 0.1 cm) and weight (about 0.1 kg) accuracy. Moreover, a body composition analyzer device (Medigate Inc., BoCA x1, Korea) was applied to measure the body mass index (BMI). **Table 1** illustrates the participants' demographic characteristics.

Table 1. Demographics (mean ± standard deviation) of participants who completed.

Participant	Age (years)	Height (m)	Weight (kg)	BMI (kg/m ²)	Vo2max (mL kg ⁻¹ min ⁻¹)
Male (n = 30)	46.6 ± 4.9	1.70 ± 0.07	85.35 ± 11.91	29.2 ± 3.1	34.55 ± 3.02
Female (n = 30)	44.37 ± 4.1	1.57 ± 0.05	72.15 ± 8.9	29.1 ± 3.04	32.94 ± 3.04

BMI—body mass index.

2.5. The maximal exercise protocol (MEP)

In this study, the Bruce protocol was applied as the maximal exercise protocol (MEP), consisting of a 3-minute stage workout that gradual elevation occurs in both speed and grade, subsequently, until the individual feels exhausted [84]. The exact details have been described previously [85]. Also, we encouraged the participants to continue the MEP until their maximal tolerance, which was the heart rate (HR) value, reached 80% to 90% of HRmax.

To assess the VO2max, the standard equation was the reference, which has been published elsewhere [86]. In this case, a calibrated treadmill (h/p/cosmos Sports and Medical GmbH, Mercury model, Nussdorf-Traunstein Germany) was applied while we evaluated the Borg 6–20 scale during the MEP, also known as the ratings of perceived exertion (RPE).

2.6. Kubios and biofeedback HRV analyses

Based on Kubios HRV analysis, HR and RR intervals were continuously recorded via standard 12-lead electrocardiography (Custo cardio 100, Custo med GmbH, Ottobrunn, Germany) at pre- (rest status including a 3-minute duration of HR stabilization) and post-exercise (instantly after performing MEP) conditions, with sampling rate set at 1000 Hz (at seated posture). Next, the investigators collected the RR intervals while visually inspecting and omitting any premature beats and artifact/noise from all recorded RR intervals. Then, we export all collected RR intervals from the ECG manufacturer's software (Medset, Hamburg, Germany) to analyze them by customized software (Kubios HRV software, version 2.1, Department of Applied Physics, University of Eastern Finland, Kuopio, Finland). Based on former evidence, there are no differences across various Kubios filter levels in adults [87]. Therefore, we used a very strong filter level in this study [88].

On the other hand, to assess HRV_{BIO}, data was recorded by a Biofeedback device

(version 4.2, Biofeedback 2000 x-pert software, made in Austria) from the pronation surface of the hand at pre- (rest status including a 3-minute duration of HR stabilization) and post-exercise (instantly after performing MEP) conditions, and the sampling rate was set at 1000 Hz as described previously [79,89]. Briefly, the blue electrode cable was attached to the back of the right hand, while the red one was attached in the same spot but to the left hand. Also, the black electrode cable was attached to the back of the non-dominant hand (i.e., the left hand for right-handed people) [90]. To prevent the noise, we tried to keep the reference constant.

In this study, measured indices of HRV consisted of time-domain variables (i.e., standard deviation of normal RR intervals, SDNN; root mean square of successive differences, RMSSD; the proportion of differences between adjacent NN intervals of more than 50 ms pNN50), frequency domain variables (i.e., the low-frequency band, LF (0.04–0.15 Hz), the high-frequency band, HF (0.15–0.40 Hz), the LF/HF ratio), and nonlinear measures (i.e., standard deviation of the instantaneous beat-to-beat RR interval variability or minor axis of the Poincare plot, SD1; the standard deviation of continuous long-term RR interval variability or major axis of the Poincare plot, SD2) [91,92].

2.7. Statistical analysis

The SPSS software (version 27.0 for Windows, IBM, Armonk, NY, USA) was applied for all statistical analyses while we drew the figures with the GraphPad Prism® software (version 9 for Windows, GraphPad Software, Inc., La Jolla, CA, USA). Firstly, the normality distribution of data was analyzed using the Kolmogorov-Smirnov test. The Pearson correlation test was used to analyze the overall association between indices. Pearson correlation coefficient (r) from 0.3 to 0.5 was considered as low, 0.5 to 0.7 as moderate, and 0.7 to 0.9 as high correlation [93]. Intra-class correlation coefficient (ICC) analysis was also performed to examine the agreement between examined variables. Values for ICC were calculated using a 2-way mixed model and interpreted as excellent (0.90 or higher), good (0.75 to 0.90), moderate (0.50 to 0.75), or poor (below 0.50) [94]. Bland Altman analysis was also used to test the agreement between values of examined variables as well as to visually depict the individual dispersion patterns [95]. Data are reported by their mean standard deviation. In this study, $P < 0.05$ was settled as the significant value.

3. Results

3.1. Correlation and agreement between HRV indices at the pre-exercise status

At rest status, the HRV parameters were measured using the BIO and ECG devices and are presented in **Figure 1** and **Table 2**. Regarding the Pearson test, all HRV parameters had significant correlations ranging from low ($r = 0.241$) to high ($r = 0.779$), such as RR interval ($r = 0.639$, $p < 0.001$), SDNN ($r = 0.779$, $p < 0.001$), RMSSD ($r = 0.625$, $p < 0.001$), pNN50 ($r = 0.455$, $p < 0.05$), LF ($r = 0.524$, $p < 0.001$), HF ($r = 0.589$, $p < 0.001$), LF/HF ratio ($r = 0.559$, $p < 0.001$), and SD2 ($r = 0.313$, $p < 0.05$); except for SD1, which only showed a certain trend toward significance ($r =$

0.241, $p = 0.064$) (**Figure 1**). Similarly, based on intra-class correlation coefficient (ICC) analysis, the indices obtained from these measurement devices, including RR intervals (ICC = 0.780, $p < 0.001$), SDNN (ICC = 0.874, $p < 0.001$), RMSSD (ICC = 0.769, $p < 0.001$), and HF (ICC = 0.722, $p < 0.001$), showed a significant agreement (**Table 2**). Meanwhile, PNN50 (ICC = 0.612, $p = 0.008$), LF (ICC = 0.641, $p < 0.001$), and LF/HF ratio (ICC = 0.593, $p < 0.001$) illustrate a considerable relationship between BIO and ECG, while SD1 (ICC = 0.379, $p = 0.035$) and SD2 (ICC = 0.47, $p = 0.008$) had a slight correlation (**Table 2**).

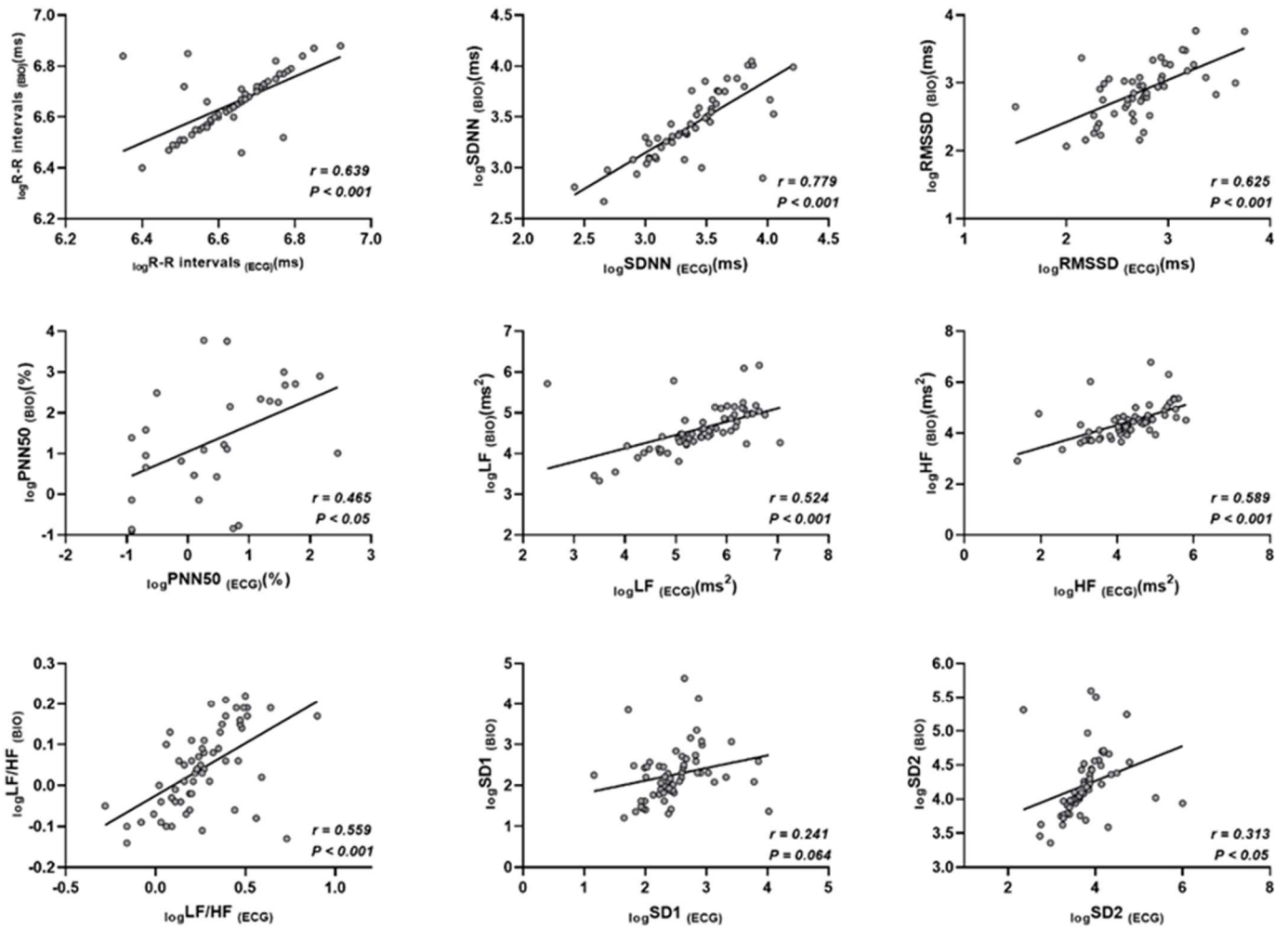


Figure 1. Pearson correlations between heart rate variability (HRV) parameters extracted via Kubios HRV and biofeedback device at pre-exercise status. Abbreviations: BIO, the biofeedback device; ECG, the electrocardiogram device; SDNN, Standard deviation of NN intervals; RMSSD, Root mean square of successive RR interval differences; LF, low-frequency; HF, high-frequency; LF/HF, LF/HF ratio; SD1, Poincaré plot standard deviation perpendicular the line of identity; SD2, Poincaré plot standard deviation along the line of identity.

Table 2. Intra-class correlation between heart rate variability (HRV) parameters obtained from Kubios HRV software and biofeedback device (BIO) at pre-exercise status.

Parameters	Mean ± standard deviation		Interclass correlation		
	Kubios HRV	BIO	ICC	95% CI	P
RR intervals (ms)	6.63 ± 0.11	6.65 ± 0.11	0.780*	0.630–0.869	< 0.001
SDNN (ms)	3.38 ± 0.41	3.42 ± 0.33	0.874*	0.779–0.928	< 0.001
RMSSD (ms)	2.75 ± 0.42	2.88 ± 0.42	0.769*	0.598–0.868	< 0.001
PNN50 (%)	37 ± 1.080	0.74 ± 1.56	0.612*	0.162–0.821	0.008
LF (ms ²)	5.46 ± 0.90	4.60 ± 0.56	0.641*	0.398–0.785	< 0.001
HF (ms ²)	4.24 ± 0.90	4.43 ± 0.67	0.722*	0.535–0.834	< 0.001
LF/HF ratio	0.26 ± 0.21	-0.431 ± 0.99	0.593*	0.319–0.757	< 0.001
SD1	2.48 ± 0.52	2.26 ± 0.66	0.379*	-0.039–0.629	0.035
SD2	3.7 ± 0.57	4.20 ± 0.47	0.470*	0.112–0.683	0.008

* Significant observation.

3.2. Correlation and agreement between indices at post-exercise

At post-exercise status, the HRV parameters were measured using the BIO and ECG devices, which are presented in **Figure 2** and **Table 3**. Regarding the Pearson test, some HRV parameters had low correlations ($0.3 < r < 0.5$) to high correlations ($r = 0.779$), such as RR interval ($r = 0.496, p < 0.001$), LF ($r = 0.260, p < 0.05$), HF ($r = 0.369, p < 0.01$), LF/HF ratio ($r = 0.394, p < 0.01$), and SD2 ($r = 0.299, p < 0.05$) (**Figure 2**). Regardless, other HRV parameters did not show any relationships between BIO and ECG at post-exercise conditions, including SDNN ($r = 0.099, p = 0.451$), RMSSD ($r = 0.118, p = 0.369$), PNN50 ($r = 0.135, p = 0.548$), and SD1 ($r = 0.117, p = 0.372$) (**Figure 2**). Similarly, based on intra-class correlation coefficient (ICC) analysis, some HRV indices obtained from these measurement devices illustrate a considerable agreement between BIO and ECG, ranging from low (below 0.50) to moderate (0.50 to 0.75), including RR intervals (ICC = 0.623, $p < 0.001$), HF (ICC = 0.553, $p = 0.001$), and LF/HF ratio (ICC = 0.506, $p < 0.001$), while LF (ICC = 0.438, $p < 0.014$) and SD2 (ICC = 0.394, $p = 0.028$) had a slight correlation (**Table 3**). Despite this, no agreements were noted among other indices, including SDNN (ICC = 0.180, $p = 0.224$), RMSSD (ICC = 0.172, $p = 0.235$), PNN50 (ICC = 0.212, $p = 0.295$), and SD1 (ICC = 0.178, $p = 0.228$) (**Table 3**).

Table 3. Intra-class correlation between heart rate variability (HRV) parameters extracted via Kubios HRV and biofeedback device (BIO) at post-exercise status.

Parameters	Mean ± SD		Interclass correlation		
	Kubios HRV	BIO	ICC	95% CI	P
RR intervals (ms)	6.21 ± 0.09	6.38 ± 0.14	0.623*	0.35–0.776	<0.001
SDNN (ms)	4.39 ± 0.38	4.06 ± 0.41	0.180	-0.373–0.510	0.224
RMSSD (ms)	2.20 ± 0.5	3.35 ± 1	0.172	-0.386–0.506	0.235
PNN50 (%)	-0.12 ± 0.98	1.05 ± 1.60	0.212	-0.899–0.673	0.295

Table 3. (Continued).

Parameters	Mean ± SD		Interclass correlation		
	Kubios HRV	BIO	ICC	95% CI	P
LF (ms ²)	3.73 ± 1.3	4.20 ± 0.92	0.394*	−0.015–0.638	0.028
HF (ms ²)	2.56 ± 1.45	4.33 ± 1.12	0.553*	0.252–0.733	0.001
LF/HF ratio	1.16 ± 1	−0.13 ± 0.57	0.506*	0.173–0.705	<0.001
SD1	1.9 ± 0.59	2.75 ± 1.11	0.178	−0.377–0.509	0.228
SD2	4.71 ± 0.3	4.69 ± 0.43	0.438*	−0.060–0.664	0.014

* significant observation.

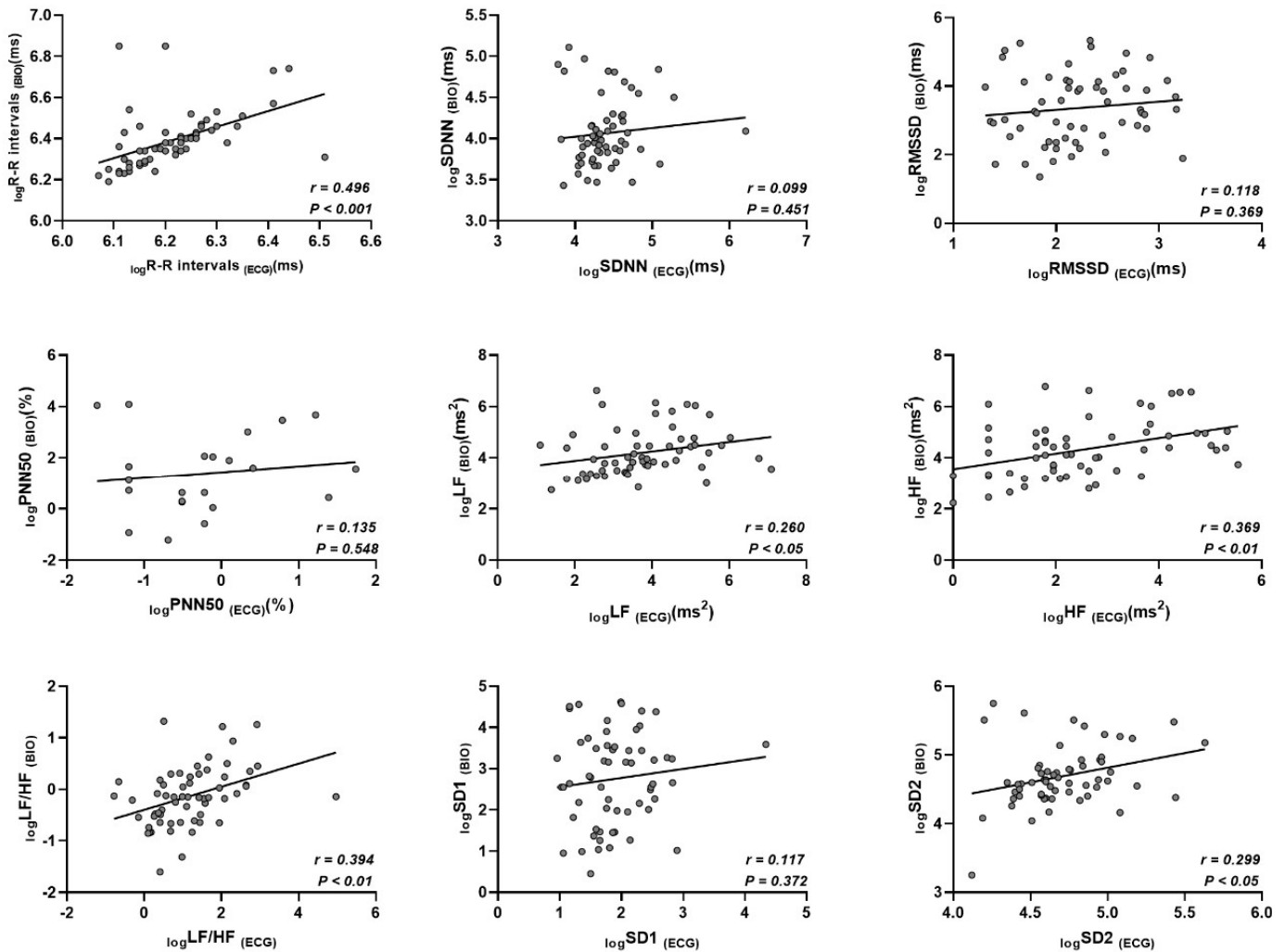


Figure 2. Pearson correlations between heart rate variability (HRV) parameters extracted via Kubios HRV and biofeedback device at post-exercise status. Abbreviations: BIO, the biofeedback device; ECG, the electrocardiogram device; SDNN, Standard deviation of NN intervals; RMSSD, Root mean square of successive RR interval differences; LF, low-frequency; HF, high-frequency; LF/HF, LF/HF ratio; SD1, Poincaré plot standard deviation perpendicular the line of identity; SD2, Poincaré plot standard deviation along the line of identity.

4. Discussion

Based on our knowledge, our research is the first study to have assessed the HRV measurement among healthy individuals by biofeedback, which means it

demonstrated the significant correlation between ECG (as the golden-standard method) and biofeedback, especially at pre-exercise status. Therefore, this is the first time that HRV-bio is considered a swift real-time method for monitoring HRV among a healthy population before and after a physical activity performance (i.e., MEP).

HRV is known as a productive way to understand the cardiovascular response in the human body [96]. Healthy heart oscillations are constantly changing, which helps the cardiovascular system adjust rapidly to any homeostasis challenges (either physical or psychological) [1,96]. HRV has assessed the neuro-cardiac function, which shows the direct relation between cardiac rhythm and the ANS branches, including the sympathetic and parasympathetic systems [97]. Therefore, HRV represents an emergent property of interdependent regulating systems, which provide various time scales to respond to any psycho-environmental challenges [98]. In healthy individuals, it reflects the satisfied regulation of different items in the body, such as autonomic balance, blood pressure, gas exchange, gut, heart, and vascular tone [81], while any diseases would involve either a decrease or elevation in the complexity of the biological system [99]. Recently, De Groot et al. have declared that ANS dysregulation symptoms are associated with diabetes-related distress among adults suffering from type 1 diabetes [100]. As a result, close monitoring of electrocardiogram (ECG) morphology would declare that increased HRV values are due to common cardiovascular conditions, including hypertension, diabetes mellitus, myocardial infarction, and heart failure [101,102], atrial fibrillation [103], and an early indication of infection [104,105]. Pathologically, clinical-dependent bradycardia could stem from vagal tone withdrawal (i.e., parasympathetic activity reduction), which causes the cardiac pacemaker to be more vulnerable to sympathetic impacts [106].

On the other hand, HRV indices are beneficial far beyond clinical prediction, which are considerably strong biomarkers to monitor physiological activity and workout levels [57,96]. Therefore, it can be applied to evaluate the level of exercise stress, especially the acute intensity by changes of the ANS following any exercise, which means it would be considered either overtraining or an overreaching marker [107,108]. Based on the PNS activity of every individual, monitoring HRV would be used to check individualized training improvement [109]. It has been illustrated that guided training based on HRV is a beneficial way to improve performance [109]. Being exposed to any physical stressor lessens HRV, which occurs as a result of vagal tone withdrawal and activating the sympathetic nervous system for supplying any exercises and physical activities' demands [108].

It is noted that an increased LF/HF ratio promotes cognitive performance [110], which is known as various strict internal operations reflected by behavior [111], while it is discovered that lessened vagal control (especially HF) is related to reduced ability of dynamical response to changing demands and environments, followed by the reduction of possible options' range and the limitation of an individual's ability to produce suitable responses and prevent inappropriate ones [112]. Likewise, it has been demonstrated that a low HRV is related to poorer performance associated with short- and long-term verbal memory [113]. If executive functioning is required for a cognitive task, therefore, vagal withdrawal is considered maladaptive [114], while it would be advantageous whenever a person is subjected to any mental stressors without including executive function, which means it is demonstrated as an individual's ability

to deal with the stimulus successfully [115,116]. Moreover, it is expressed that cardiac vagal activity can be considered an index of self-regulatory and/or cognitive-related processing [117–119]. The study literature illustrates the necessity of monitoring HRV for both physical and mental purposes among healthy people via a swift real-time method.

ECG is a traditional method for measuring HRV, which requires time and expertise. Moreover, other devices are high-priced, including clinical multi-lead ECG (e.g., Holter), which would not be practical for field-based monitoring in active and healthy individuals. Although multi-lead ECG devices are considered the golden standard, devices based on single-lead ECG or photo-plethysmography (PPG) are simple to apply [57]. Applying PPG technology to measure HRV is a recent and novel method, which is integrated into wearable wrist and finger-worn devices. Despite the motion artifact noted as a limitation of this method, their comfort and feasibility make them attractive alternatives to multi-lead ECG systems [58].

Regarding our study, there is significant agreement between ECG and BIO devices for measuring HRV indexes among healthy individuals at rest conditions. To prove this, the RR interval had moderate correlation ($r = 0.639, p < 0.001$), good ICC ($r = 0.780, p < 0.001$), and an average deviation of -0.01593 ms according to the Bland-Altman plots (95% LoA: -0.1984 to 0.1665 ms). As for SDNN, it showed high agreement ($r = 0.779, p < 0.001$), acceptable ICC ($r = 0.874, p < 0.001$), and an average deviation of -0.03843 ms according to the Bland-Altman plots (95% LoA: -0.5041 to 0.5664 ms). Also, RMSSD noted moderate correlation ($r = 0.625, p < 0.001$), good ICC ($r = 0.769, p < 0.001$), and an average deviation of -0.1465 ms according to the Bland-Altman plots (95% LoA: -0.8594 to 0.5664 ms). Plus, LF noted moderate correlation ($r = 0.524, p < 0.001$), low ICC ($r = 0.641, p < 0.001$), and an average deviation of -0.8597 ms² according to the Bland-Altman plots (95% LoA: -0.6665 to 2.386 ms²). In addition, as for HF, it showed moderate agreement ($r = 0.589, p < 0.001$), good ICC ($r = 0.722, p < 0.001$), and an average deviation of -0.1835 ms² according to the Bland-Altman plots (95% LoA: -1.644 to 1.277 ms²). Despite this, PNN50 illustrated low correlation ($r = 0.455, p < 0.05$), moderate ICC ($r = 0.612, p = 0.008$), and an average deviation of -0.8882% according to the Bland-Altman plots (95% LoA: -3.386 to 1.609%). In addition, as for the LF/HF ratio, it showed moderate agreement ($r = 0.559, p < 0.001$), good ICC ($r = 0.593, p < 0.001$), and an average deviation of 0.2213 according to the Bland-Altman plots (95% LoA: -0.1381 to 0.5808). Despite this, SD2 illustrated low correlation ($r = 0.313, p < 0.05$), low ICC ($r = 0.47, p = 0.008$), and an average deviation of -0.4317 according to the Bland-Altman plots (95% LoA: -1.645 to 0.7821), while SD1 illustrated no correlation ($r = 0.241, P = 0.064$), low ICC ($r = 0.379, p = 0.035$), and an average deviation of 0.2155 according to the Bland-Altman plots (95% LoA: -1.244 to 1.675) (**Figure 3**). On the other hand, at post-exercise status, the RR interval had low correlation ($r = 0.496, p < 0.001$), moderate ICC ($r = 0.623, p < 0.001$), and an average deviation of -0.1785 ms according to the Bland-Altman plots (95% LoA: -0.4270 to 0.07008 ms). As for SDNN, it did not show agreement ($r = 0.099, p = 0.451$), ICC ($r = 0.180, p = 0.224$), and had an average deviation of 0.3343 ms according to the Bland-Altman plots (95% LoA: -0.7245 to 1.393 ms). Also, RMSSD did not note any correlation ($r = 0.118, p = 0.369$), and ICC ($r = 0.172, p = 0.235$) had an average deviation of -1.150 ms

according to the Bland-Altman plots (95% LoA: -3.254 to 0.9540 ms). Whereas LF noted low correlation ($r = 0.260$, $p < 0.05$), low ICC ($r = 0.438$, $p < 0.014$), and an average deviation of -0.4675 ms² according to the Bland-Altman plots (95% LoA: -3.195 to 2.260 ms²). In addition, as for HF, it showed low agreement ($r = 0.369$, $p < 0.01$), moderate ICC ($r = 0.553$, $p = 0.001$), and an average deviation of -1.772 ms² according to the Bland-Altman plots (95% LoA: -4.602 to 1.059 ms²). Despite this,

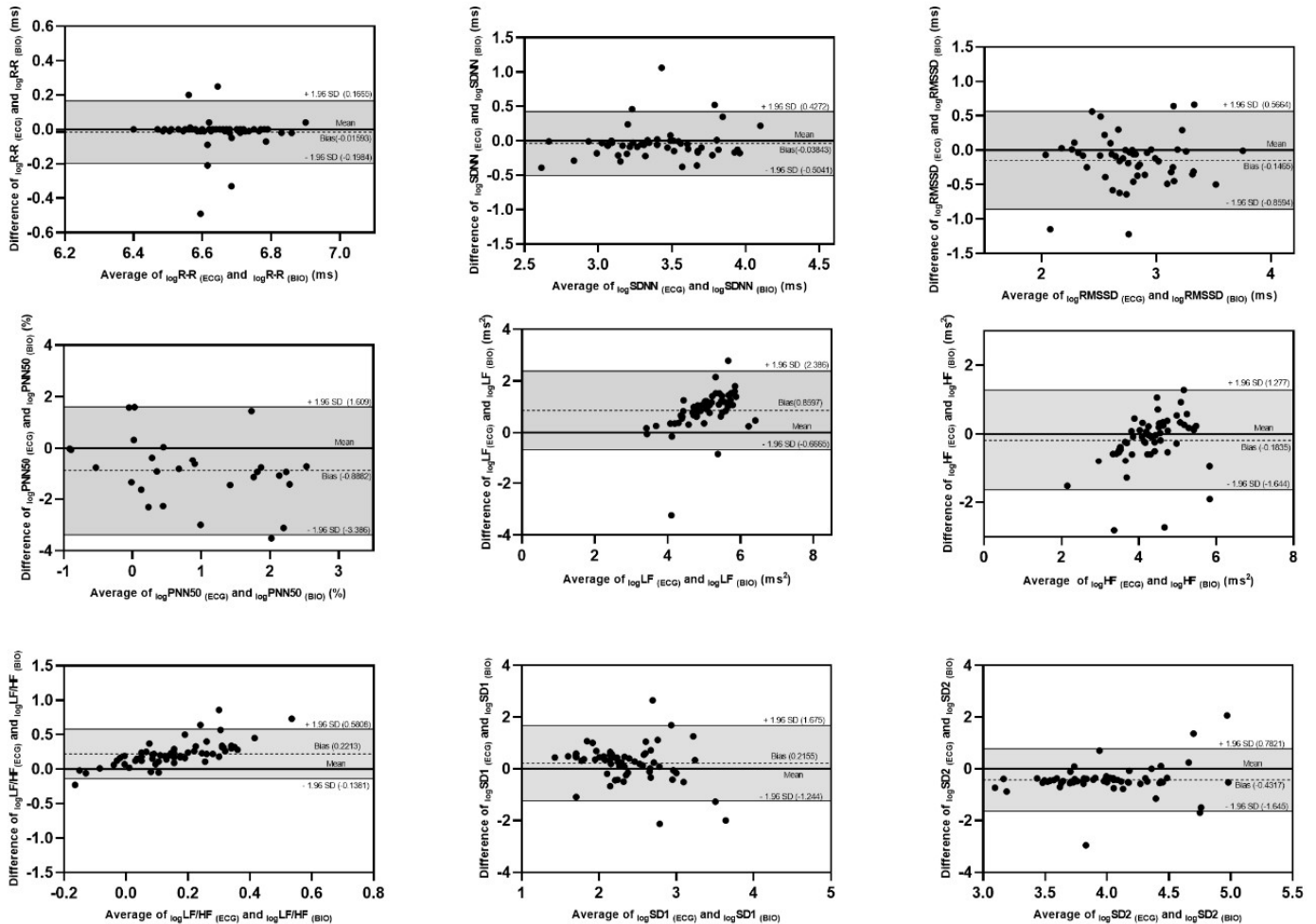


Figure 3. Bland-Altman plots of heart rate variability (HRV) parameters extracted via Kubios HRV and biofeedback device at pre-exercise status. Abbreviations: BIO, the biofeedback device; ECG, the electrocardiogram device; SDNN, Standard deviation of NN intervals; RMSSD, Root mean square of successive RR interval differences; LF, low-frequency; HF, high-frequency; LF/HF, LF/HF ratio; SD1, Poincaré plot standard deviation perpendicular the line of identity; SD2, Poincaré plot standard deviation along the line of identity.

PNN50 did not illustrate correlation ($r = 0.135$, $p = 0.548$), moderate ICC ($r = 0.212$, $p = 0.295$), and an average deviation of -1.603% according to the Bland-Altman plots (95% LoA: -4.926% to 1.720%). As for LF/HF ratio, it showed low agreement ($r = 0.394$, $p < 0.01$), moderate ICC ($r = 0.506$, $p < 0.001$), and an average deviation of 1.305 according to the Bland-Altman plots (95% LoA: -0.5411 to 3.151). Despite this, SD2 illustrated low correlation ($r = 0.299$, $p < 0.05$), low ICC ($r = 0.394$, $p = 0.028$), and an average deviation of 0.02067 according to the Bland-Altman plots (95% LoA: -0.8621 to 0.9035), while SD1 illustrated no correlation ($r = 0.117$, $p = 0.372$), ICC

($r = 0.178$, $p = 0.228$), and an average deviation of -0.8513 according to the Bland-Altman plots (95% LoA: -3.201 to 1.498) (**Figure 4**). Therefore, not only is a biofeedback device considered an effective method for monitoring autonomic balance [12,74], but it can also be applied among healthy individuals regarding our study, especially at pre-exercise status.

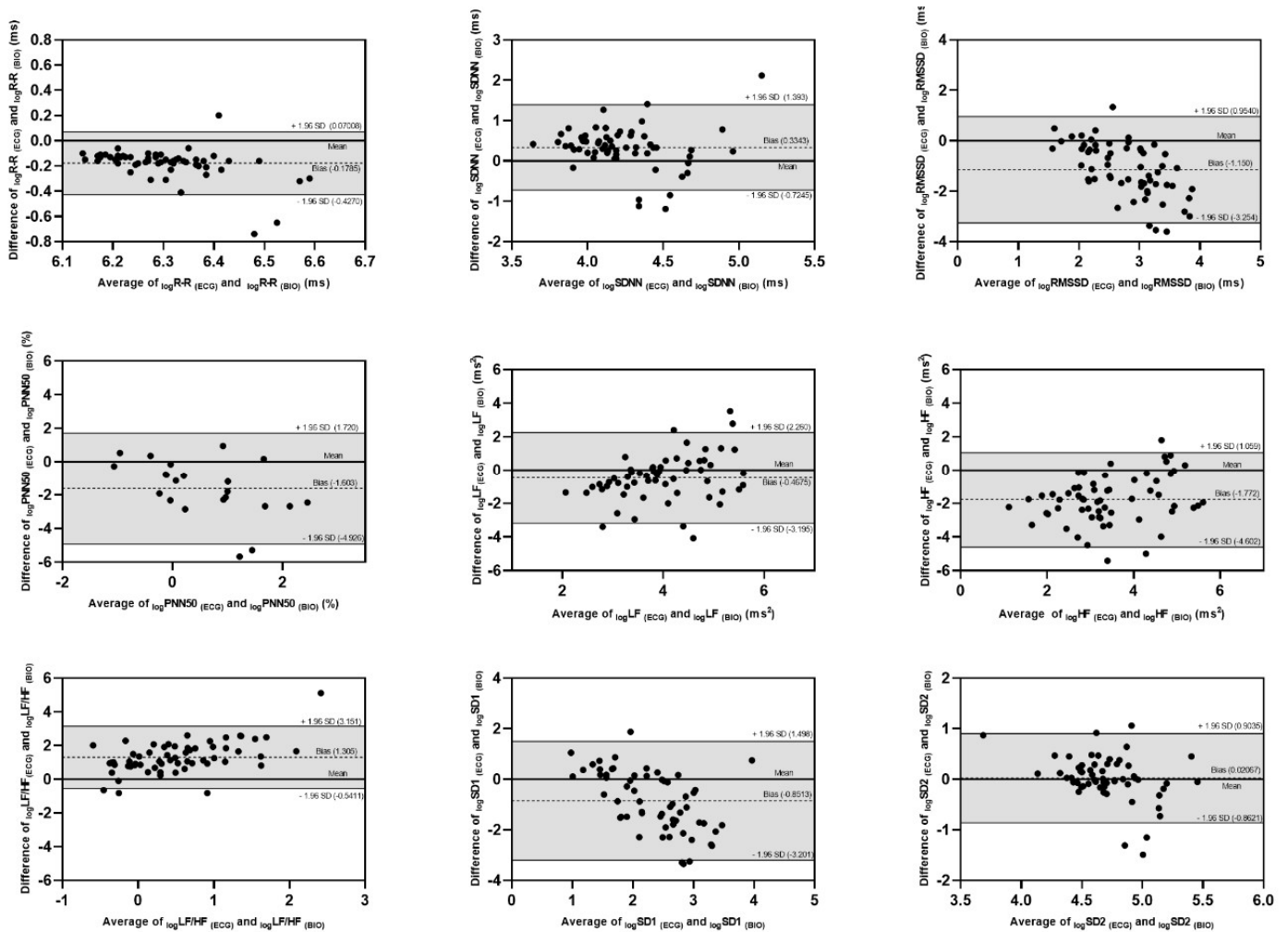


Figure 4. Bland-Altman plots of heart rate variability (HRV) parameters extracted via Kubios HRV and biofeedback device at post-exercise. Abbreviations: BIO, the biofeedback device; ECG, the electrocardiogram device; SDNN, Standard deviation of NN intervals; RMSSD, Root mean square of successive RR interval differences; LF, low-frequency; HF, high-frequency; LF/HF, LF/HF ratio; SD1, Poincaré plot standard deviation perpendicular the line of identity; SD2, Poincaré plot standard deviation along the line of identity.

Altogether, this study showed a significant correlation between the golden method and the biofeedback device, which illustrates that the biofeedback device’s usefulness is far beyond its clinical activities. Despite this, it should not be overlooked that we had some limitations in this study. Firstly, since the purpose of this paper was to assess whether biofeedback can be applied for HRV measurement among healthy people, we did not include athlete populations. Secondly, although the women were in the initial follicular phase of the menstrual cycle during the experiment period, the ovarian hormone levels were not measured directly. Finally, to prevent any possible noises being caused by body movement, we could not apply the biofeedback device

while performing MEP. Therefore, further studies are required to elucidate the impacts of these limitations while measuring HRV by biofeedback among healthy individuals.

5. Conclusion

Monitoring HRV in the population could be advantageous for tailoring individualized training and exercise programs according to the onset of illness or infection, identifying the risk of overreaching and overtraining, quantifying cognitive performance, and as an overall measure of health [31–49]. Regarding the preceding paragraphs, although the HRV_{ECG} is known as the golden method, it has several limitations [120,121], such as either expertise or time requirements, and good-quality electrode signals [122], which also prevents its applicability for prolonged-daily measurement [123]. Uniquely, this study states that biofeedback can be considered a facilitative way to evaluate HRV among healthy people, especially at pre-exercise status. Further research would be relevant for specific the HRV_{BIO} at different timelines of performing any exercises.

Author contributions: Conceptualization, VDR and SP; methodology, VDR, SP and SNT; software, SP and SNT; validation, VDR, SP and SNT; formal analysis, VDR, SP and SNT; investigation, SP and SNT; resources, VDR and SP; data curation, SP and SNT; writing—original draft preparation, SP and VDR; writing—review and editing, VDR, SP and SNT; visualization, SP; supervision, VDR. All authors have read and agreed to the published version of the manuscript.

Acknowledgments: The authors would like to thank all participants who consented to participate in the present study. Also, we appreciate the efforts of Mehdi Ahmadian for his help and guidance in drawing the figures using the standard method.

Practical Implications (highlights): (1) A significant correlation was noted between HRV_{ECG} and HRV_{BIO}. (2) Biofeedback devices can be considered novel methodological tools to monitor the cardiovascular autonomic system before any physical activities and exercises. (3) Biofeedback devices can be applied for faster examination of HRV in healthy individuals.

Conflict of interest: The authors declare no conflict of interest.

References

1. Shaffer F, Ginsberg JP. An Overview of Heart Rate Variability Metrics and Norms. *Frontiers in Public Health*. 2017; 5. doi: 10.3389/fpubh.2017.00258
2. Vanderlei LCM, Pastre CM, Hoshi RA, et al. Basics of heart rate variability and its clinical applicability (Portuguese). *Revista Brasileira de Cirurgia Cardiovascular*. 2009; 24(2): 205-217. doi: 10.1590/s0102-76382009000200018
3. Malliani A, Montano N. Heart rate variability as a clinical tool. *Italian Heart Journal*. 2002; 3: 439-445.
4. Malik M. Heart Rate Variability. *Annals of Noninvasive Electrocardiology*. 1996; 1(2): 151-181. doi: 10.1111/j.1542-474x.1996.tb00275.x
5. De Maria B, Dalla Vecchia LA, Porta A, et al. Autonomic dysfunction and heart rate variability with Holter monitoring: a diagnostic look at autonomic regulation. *Herzschrittmachertherapie + Elektrophysiologie*. 2021; 32(3): 315-319. doi: 10.1007/s00399-021-00780-5
6. Treiber FA, Kamarck T, Schneiderman N, et al. Cardiovascular Reactivity and Development of Preclinical and Clinical Disease States. *Psychosomatic Medicine*. 2003; 65(1): 46-62. doi: 10.1097/00006842-200301000-00007

7. Phillips AC. Blunted as well as exaggerated cardiovascular reactivity to stress is associated with negative health outcomes. *Japanese Psychological Research*. 2011; 53(2): 177-192. doi: 10.1111/j.1468-5884.2011.00464.x
8. Lovallo WR. Cardiovascular reactivity: Mechanisms and pathways to cardiovascular disease. *International Journal of Psychophysiology*. 2005; 58(2-3): 119-132. doi: 10.1016/j.ijpsycho.2004.11.007
9. Daanen HAM, Lamberts RP, Kallen VL, et al. A Systematic Review on Heart-Rate Recovery to Monitor Changes in Training Status in Athletes. *International Journal of Sports Physiology and Performance*. 2012; 7(3): 251-260. doi: 10.1123/ijpspp.7.3.251
10. Lamberts RP, Swart J, Noakes TD, et al. Changes in heart rate recovery after high-intensity training in well-trained cyclists. *European Journal of Applied Physiology*. 2008; 105(5): 705-713. doi: 10.1007/s00421-008-0952-y
11. Borresen J, Lambert MI. Autonomic Control of Heart Rate during and after Exercise. *Sports Medicine*. 2008; 38(8): 633-646. doi: 10.2165/00007256-200838080-00002
12. Lin IM, Fan SY, Yen CF, et al. Heart Rate Variability Biofeedback Increased Autonomic Activation and Improved Symptoms of Depression and Insomnia among Patients with Major Depression Disorder. *Clinical Psychopharmacology and Neuroscience*. 2019; 17(2): 222-232. doi: 10.9758/cpn.2019.17.2.222
13. Zucker TL, Samuelson KW, Muench F, et al. The Effects of Respiratory Sinus Arrhythmia Biofeedback on Heart Rate Variability and Posttraumatic Stress Disorder Symptoms: A Pilot Study. *Applied Psychophysiology and Biofeedback*. 2009; 34(2): 135-143. doi: 10.1007/s10484-009-9085-2
14. Sánchez-Delgado JC, Jácome-Hortúa AM, Yoshida de Melo K, et al. Physical Exercise Effects on Cardiovascular Autonomic Modulation in Postmenopausal Women—A Systematic Review and Meta-Analysis. *International Journal of Environmental Research and Public Health*. 2023; 20(3): 2207. doi: 10.3390/ijerph20032207
15. Souza H. Autonomic Cardiovascular Damage during Post-menopause: the Role of Physical Training. *Aging and Disease*. 2013; 4(6): 320-328. doi: 10.14336/ad.2013.0400320
16. Souza HCD, Philbois SV, Veiga AC, et al. Heart Rate Variability and Cardiovascular Fitness: What We Know so Far. *Vascular Health and Risk Management*. 2021; 17: 701-711. doi: 10.2147/vhrm.s279322
17. França da Silva AK, Penachini da Costa de Rezende Barbosa M, Marques Vanderlei F, et al. Application of Heart Rate Variability in Diagnosis and Prognosis of Individuals with Diabetes Mellitus: Systematic Review. *Annals of Noninvasive Electrocardiology*. 2016; 21(3): 223-235. doi: 10.1111/anec.12372
18. Tsuji H, Larson MG, Venditti FJ, et al. Impact of Reduced Heart Rate Variability on Risk for Cardiac Events. *Circulation*. 1996; 94(11): 2850-2855. doi: 10.1161/01.cir.94.11.2850
19. Dekker JM, Crow RS, Folsom AR, et al. Low Heart Rate Variability in a 2-Minute Rhythm Strip Predicts Risk of Coronary Heart Disease and Mortality from Several Causes. *Circulation*. 2000; 102(11): 1239-1244. doi: 10.1161/01.cir.102.11.1239
20. Antelmi I, De Paula RS, Shinzato AR, et al. Influence of age, gender, body mass index, and functional capacity on heart rate variability in a cohort of subjects without heart disease. *The American Journal of Cardiology*. 2004; 93(3): 381-385. doi: 10.1016/j.amjcard.2003.09.065
21. Ogliaari G, Mahinrad S, Stott DJ, et al. Resting heart rate, heart rate variability and functional decline in old age. *Canadian Medical Association Journal*. 2015; 187(15): E442-E449. doi: 10.1503/cmaj.150462
22. Nolan J, Batin PD, Andrews R, et al. Prospective Study of Heart Rate Variability and Mortality in Chronic Heart Failure. *Circulation*. 1998; 98(15): 1510-1516. doi: 10.1161/01.cir.98.15.1510
23. Freitas VP de, Passos R da S, Oliveira AA, et al. Sarcopenia is associated to an impaired autonomic heart rate modulation in community-dwelling old adults. *Archives of Gerontology and Geriatrics*. 2018; 76: 120-124. doi: 10.1016/j.archger.2018.01.006
24. Wu WC, Wong EC. Feasibility of Velocity Selective Arterial Spin Labeling in Functional MRI. *Journal of Cerebral Blood Flow & Metabolism*. 2007; 27(4): 831-838. doi: 10.1038/sj.jcbfm.9600386
25. Thayer JF, Sternberg E. Beyond Heart Rate Variability. *Annals of the New York Academy of Sciences*. 2006; 1088(1): 361-372. doi: 10.1196/annals.1366.014
26. Spiegelhalder K, Fuchs L, Ladwig J, et al. Heart rate and heart rate variability in subjectively reported insomnia. *Journal of Sleep Research*. 2011; 20: 137-145. doi: 10.1111/j.1365-2869.2010.00863.x
27. Furlan R, Piazza S, Dell'Orto S, et al. Early and late effects of exercise and athletic training on neural mechanisms controlling heart rate. *Cardiovascular Research*. 1993; 27(3): 482-488. doi: 10.1093/cvr/27.3.482
28. Hautala AJ, Mäkikallio TH, Kiviniemi A, et al. Cardiovascular autonomic function correlates with the response to aerobic

- training in healthy sedentary subjects. *American Journal of Physiology-Heart and Circulatory Physiology*. 2003; 285(4): H1747-H1752. doi: 10.1152/ajpheart.00202.2003
29. Hautala A, Tulppo MP, Mäkikallio TH, et al. Changes in cardiac autonomic regulation after prolonged maximal exercise. *Clinical Physiology*. 2001; 21(2): 238-245. doi: 10.1046/j.1365-2281.2001.00309.x
 30. Arêas GPT, Caruso FCR, Simões RP, et al. Ultra-short-term heart rate variability during resistance exercise in the elderly. *Brazilian Journal of Medical and Biological Research*. 2018; 51(6). doi: 10.1590/1414-431x20186962
 31. Lin YY, Lee SD. Cardiovascular Benefits of Exercise Training in Postmenopausal Hypertension. *International Journal of Molecular Sciences*. 2018; 19(9): 2523. doi: 10.3390/ijms19092523
 32. Souza HCD, Ballejo G, Salgado MCO, et al. Cardiac sympathetic overactivity and decreased baroreflex sensitivity in L-NAME hypertensive rats. *American Journal of Physiology-Heart and Circulatory Physiology*. 2001; 280(2): H844-H850. doi: 10.1152/ajpheart.2001.280.2.h844
 33. Freire Machi J, da Silva Dias D, Freitas S, et al. Impact of aging on cardiac function in a female rat model of menopause: role of autonomic control, inflammation, and oxidative stress. *Clinical Interventions in Aging*. 2016; 341. doi: 10.2147/cia.s88441
 34. Shimojo GL, Palma RK, Brito JO, et al. Dynamic resistance training decreases sympathetic tone in hypertensive ovariectomized rats. *Brazilian Journal of Medical and Biological Research*. 2015; 48(6): 523-527. doi: 10.1590/1414-431x20154387
 35. da Palma RK, Moraes-Silva IC, da Silva Dias D, et al. Resistance or aerobic training decreases blood pressure and improves cardiovascular autonomic control and oxidative stress in hypertensive menopausal rats. *Journal of Applied Physiology*. 2016; 121(4): 1032-1038. doi: 10.1152/jappphysiol.00130.2016
 36. Mariano IM, Freitas VH de, Batista JP, et al. Effect of combined exercise training on heart rate variability in normotensive and hypertensive postmenopausal women. *Motriz: Revista de Educação Física*. 2021; 27. doi: 10.1590/s1980-65742021020621
 37. Shimojo GL, Silva Dias D da, Malfitano C, et al. Combined Aerobic and Resistance Exercise Training Improve Hypertension Associated with Menopause. *Frontiers in Physiology*. 2018; 9. doi: 10.3389/fphys.2018.01471
 38. Weston KS, Wisløff U, Coombes JS. High-intensity interval training in patients with lifestyle-induced cardiometabolic disease: a systematic review and meta-analysis. *British Journal of Sports Medicine*. 2013; 48(16): 1227-1234. doi: 10.1136/bjsports-2013-092576
 39. O'Donnell E, Craig J. Habitual aerobic exercise in healthy postmenopausal women does not augment basal cardiac autonomic activity yet modulates autonomic-metabolic interactions. *Menopause*. 2022; 29(6): 714-722. doi: 10.1097/gme.0000000000001963
 40. da Silva CC, Pereira LM, Cardoso JR, et al. The Effect of Physical Training on Heart Rate Variability in Healthy Children: A Systematic Review with Meta-Analysis. *Pediatric Exercise Science*. 2014; 26(2): 147-158. doi: 10.1123/pes.2013-0063
 41. Grässler B, Thielmann B, Böckelmann I, et al. Effects of Different Training Interventions on Heart Rate Variability and Cardiovascular Health and Risk Factors in Young and Middle-Aged Adults: A Systematic Review. *Frontiers in Physiology*. 2021; 12. doi: 10.3389/fphys.2021.657274
 42. Villafaina S, Collado-Mateo D, Fuentes JP, et al. Physical Exercise Improves Heart Rate Variability in Patients with Type 2 Diabetes: A Systematic Review. *Current Diabetes Reports*. 2017; 17(11). doi: 10.1007/s11892-017-0941-9
 43. Belvederi Murri M, Folesani F, Zerbinati L, et al. Physical Activity Promotes Health and Reduces Cardiovascular Mortality in Depressed Populations: A Literature Overview. *International Journal of Environmental Research and Public Health*. 2020; 17(15): 5545. doi: 10.3390/ijerph17155545
 44. Palma S, Keilani M, Hasenoehrl T, et al. Impact of supportive therapy modalities on heart rate variability in cancer patients—a systematic review. *Disability and Rehabilitation*. 2018; 42(1): 36-43. doi: 10.1080/09638288.2018.1514664
 45. Sandercock GRH, Shelton C, Bromley P, et al. Agreement between three commercially available instruments for measuring short-term heart rate variability. *Physiological Measurement*. 2004; 25(5): 1115-1124. doi: 10.1088/0967-3334/25/5/003
 46. Bhati P, Moiz JA, Menon GR, et al. Does resistance training modulate cardiac autonomic control? A systematic review and meta-analysis. *Clinical Autonomic Research*. 2018; 29(1): 75-103. doi: 10.1007/s10286-018-0558-3
 47. Sant'Ana L de O, Machado S, Ribeiro AA de S, et al. Effects of Cardiovascular Interval Training in Healthy Elderly Subjects: A Systematic Review. *Frontiers in Physiology*. 2020; 11. doi: 10.3389/fphys.2020.00739
 48. Hottenrott K, Hoos O, Esperer HD. Heart rate variability and sport (German). *Herz Kardiovaskuläre Erkrankungen*. 2006;

- 31(6): 544-552. doi: 10.1007/s00059-006-2855-1
49. D'Agosto T, Peçanha T, Bartels R, et al. Cardiac Autonomic Responses at Onset of Exercise: Effects of Aerobic Fitness. *International Journal of Sports Medicine*. 2014; 35(10): 879-885. doi: 10.1055/s-0034-1370911
50. Boullosa DA, Nakamura FY, Perandini LA, et al. Autonomic correlates of Yo-Yo performance in soccer referees. *Motriz: Revista de Educação Física*. 2012; 18(2): 291-297. doi: 10.1590/s1980-65742012000200009
51. Stanley J, Peake JM, Buchheit M. Cardiac Parasympathetic Reactivation Following Exercise: Implications for Training Prescription. *Sports Medicine*. 2013; 43(12): 1259-1277. doi: 10.1007/s40279-013-0083-4
52. Li Z, Wang C, Mak AFT, et al. Effects of acupuncture on heart rate variability in normal subjects under fatigue and non-fatigue state. *European Journal of Applied Physiology*. 2005; 94(5-6): 633-640. doi: 10.1007/s00421-005-1362-z
53. Jerath R, Syam M, Ahmed S. The Future of Stress Management: Integration of Smartwatches and HRV Technology. *Sensors*. 2023; 23(17): 7314. doi: 10.3390/s23177314
54. Dalmeida KM, Masala GL. HRV Features as Viable Physiological Markers for Stress Detection Using Wearable Devices. *Sensors*. 2021; 21(8): 2873. doi: 10.3390/s21082873
55. Chalmers T, Hickey BA, Newton P, et al. Stress Watch: The Use of Heart Rate and Heart Rate Variability to Detect Stress: A Pilot Study Using Smart Watch Wearables. *Sensors*. 2021; 22(1): 151. doi: 10.3390/s22010151
56. Neri L, Corazza I, Oberdier MT, et al. Comparison Between a Single-Lead ECG Garment Device and a Holter Monitor: A Signal Quality Assessment. *Journal of Medical Systems*. 2024; 48(1). doi: 10.1007/s10916-024-02077-9
57. Gilgen-Ammann R, Schweizer T, Wyss T. RR interval signal quality of a heart rate monitor and an ECG Holter at rest and during exercise. *European Journal of Applied Physiology*. 2019; 119(7): 1525-1532. doi: 10.1007/s00421-019-04142-5
58. Singh N, Moneghetti KJ, Christle JW, et al. Heart Rate Variability: An Old Metric with New Meaning in the Era of Using mHealth technologies for Health and Exercise Training Guidance. Part Two: Prognosis and Training. *Arrhythmia & Electrophysiology Review*. 2018; 7(4): 1. doi: 10.15420/aer.2018.30.2
59. Kinnunen H, Rantanen A, Kenttä T, et al. Feasible assessment of recovery and cardiovascular health: accuracy of nocturnal HR and HRV assessed via ring PPG in comparison to medical grade ECG. *Physiological Measurement*. 2020; 41(4): 04NT01. doi: 10.1088/1361-6579/ab840a
60. Bent B, Goldstein BA, Kibbe WA, et al. Investigating sources of inaccuracy in wearable optical heart rate sensors. *npj Digital Medicine*. 2020; 3(1). doi: 10.1038/s41746-020-0226-6
61. Vescio B, Salsone M, Gambardella A, et al. Comparison between Electrocardiographic and Earlobe Pulse Photoplethysmographic Detection for Evaluating Heart Rate Variability in Healthy Subjects in Short- and Long-Term Recordings. *Sensors*. 2018; 18(3): 844. doi: 10.3390/s18030844
62. Hallman D, Sato T, Kristiansen J, et al. Prolonged Sitting is Associated with Attenuated Heart Rate Variability during Sleep in Blue-Collar Workers. *International Journal of Environmental Research and Public Health*. 2015; 12(11): 14811-14827. doi: 10.3390/ijerph121114811
63. Bogdány T, et al. Validation of the First beat Team Belt and BodyGuard2 systems. *Magyar Sporttudományi Szemle*. 2016; 17(3): 5-12.
64. Czekaj L, Daniszewski M, Domaszewicz J. Validation of the Aidlab solution for measuring Heart Rate Variability. Aidlab: Gdansk, Poland; 2019.
65. Giles D, Draper N, Neil W. Validity of the Polar V800 heart rate monitor to measure RR intervals at rest. *European Journal of Applied Physiology*. 2015; 116(3): 563-571. doi: 10.1007/s00421-015-3303-9
66. Hillebrand S, Gast KB, de Mutsert R, et al. Heart rate variability and first cardiovascular event in populations without known cardiovascular disease: meta-analysis and dose-response meta-regression. *EP Europace*. 2013; 15(5): 742-749. doi: 10.1093/europace/eus341
67. Dahiya ES, Kalra AM, Lowe A, et al. Wearable Technology for Monitoring Electrocardiograms (ECGs) in Adults: A Scoping Review. *Sensors*. 2024; 24(4): 1318. doi: 10.3390/s24041318
68. Vanderlei LCM, Silva RA, Pastre CM, et al. Comparison of the Polar S810i monitor and the ECG for the analysis of heart rate variability in the time and frequency domains. *Brazilian Journal of Medical and Biological Research*. 2008; 41(10): 854-859. doi: 10.1590/s0100-879x2008005000039
69. Weippert M, Kumar M, Kreuzfeld S, et al. Comparison of three mobile devices for measuring R-R intervals and heart rate variability: Polar S810i, Suunto t6 and an ambulatory ECG system. *European Journal of Applied Physiology*. 2010; 109(4): 779-786. doi: 10.1007/s00421-010-1415-9

70. Caminal P, Sola F, Gomis P, et al. Validity of the Polar V800 monitor for measuring heart rate variability in mountain running route conditions. *European Journal of Applied Physiology*. 2018; 118(3): 669-677. doi: 10.1007/s00421-018-3808-0
71. Kalra A, Lowe A, Al-Jumaily A. Critical review of electrocardiography measurement systems and technology. *Measurement Science and Technology*. 2018; 30(1): 012001. doi: 10.1088/1361-6501/aaf2b7
72. Khunti K. Accurate interpretation of the 12-lead ECG electrode placement: A systematic review. *Health Education Journal*. 2013; 73(5): 610-623. doi: 10.1177/0017896912472328
73. Lehrer PM. Heart rate variability biofeedback and other psychophysiological procedures as important elements in psychotherapy. *International Journal of Psychophysiology*. 2018; 131: 89-95. doi: 10.1016/j.ijpsycho.2017.09.012
74. Schäfer SK, Ihmig FR, Lara HKA, et al. Effects of heart rate variability biofeedback during exposure to fear-provoking stimuli within spider-fearful individuals: study protocol for a randomized controlled trial. *Trials*. 2018; 19(1). doi: 10.1186/s13063-018-2554-2
75. Subhani AR, Kamel N, Mohamad Saad MN, et al. Mitigation of stress: new treatment alternatives. *Cognitive Neurodynamics*. 2017; 12(1): 1-20. doi: 10.1007/s11571-017-9460-2
76. Prinsloo GE, Rauch HGL, Derman WE. A Brief Review and Clinical Application of Heart Rate Variability Biofeedback in Sports, Exercise, and Rehabilitation Medicine. *The Physician and Sportsmedicine*. 2014; 42(2): 88-99. doi: 10.3810/psm.2014.05.2061
77. Lin IM, Ko JM, Fan SY, et al. Heart Rate Variability and the Efficacy of Biofeedback in Heroin Users with Depressive Symptoms. *Clinical Psychopharmacology and Neuroscience*. 2016; 14(2): 168-176. doi: 10.9758/cpn.2016.14.2.168
78. Sutarto AP, Abdul Wahab MN, Mat Zin N. Heart Rate Variability (HRV) biofeedback: A new training approach for operator's performance enhancement. *Journal of Industrial Engineering and Management*. 2010; 3(1). doi: 10.3926/jiem.2010.v3n1.p176-198
79. Wang X, Chen P, Huang X, et al. Guasha improves the rating of perceived exertion scale score and reduces heart rate variability in male weightlifters: a randomized controlled trial. *Journal of Traditional Chinese Medicine*. 2017; 37(1): 49-56. doi: 10.1016/S0254-6272(17)30026-2
80. Strauss-Blasche G, Moser M, Voica M, et al. Relative Timing of Inspiration and Expiration Affects Respiratory Sinus Arrhythmia. *Clinical and Experimental Pharmacology and Physiology*. 2000; 27(8): 601-606. doi: 10.1046/j.1440-1681.2000.03306.x
81. Lehrer PM, Gevirtz R. Heart rate variability biofeedback: how and why does it work? *Frontiers in Psychology*. 2014; 5. doi: 10.3389/fpsyg.2014.00756
82. Shrestha B, Dunn L. The Declaration of Helsinki on Medical Research involving Human Subjects: A Review of Seventh Revision. *Journal of Nepal Health Research Council*. 2020; 17(4): 548-552. doi: 10.33314/jnhrc.v17i4.1042
83. Bergman RN, Stefanovski D, Buchanan TA, et al. A Better Index of Body Adiposity. *Obesity*. 2011; 19(5): 1083-1089. doi: 10.1038/oby.2011.38
84. Bruce RA, Kusumi F, Hosmer D. Maximal oxygen intake and nomographic assessment of functional aerobic impairment in cardiovascular disease. *American heart journal*. 1973; 85(4): 546-562. doi: 10.1016/0002-8703(73)90502-4
85. Pouriamehr S, Dabidi Roshan V, Shirani F. Does long-term exposure to air pollution suppress parasympathetic reactivation after incremental exercise among healthy males and females? *Inhalation Toxicology*. 2022; 35(1-2): 14-23. doi: 10.1080/08958378.2022.2149905
86. Jamnick NA, By S, Pettitt CD, et al. Comparison of the YMCA and a Custom Submaximal Exercise Test for Determining V̇O₂max. *Medicine & Science in Sports & Exercise*. 2016; 48(2): 254-259. doi: 10.1249/mss.0000000000000763
87. Alcantara JMA, Plaza-Florido A, Amaro-Gahete FJ, et al. Impact of Using Different Levels of Threshold-Based Artefact Correction on the Quantification of Heart Rate Variability in Three Independent Human Cohorts. *Journal of Clinical Medicine*. 2020; 9(2): 325. doi: 10.3390/jcm9020325
88. Tarvainen MP, Niskanen JP, Lipponen JA, et al. Kubios HRV—Heart rate variability analysis software. *Computer Methods and Programs in Biomedicine*. 2014; 113(1): 210-220. doi: 10.1016/j.cmpb.2013.07.024
89. Wang X, Chatchawan U, Nakmareong S, et al. Effects of GUASHA on Heart Rate Variability in Healthy Male Volunteers under Normal Condition and Weightlifters after Weightlifting Training Sessions. *Evidence-Based Complementary and Alternative Medicine*. 2015; 2015: 1-6. doi: 10.1155/2015/268471
90. Mödler J. Hardware English Assessment Systems. Available online: https://asystems.as/wp-content/uploads/2018/05/Biofeedback_2000x-per_t_Hardware_manual.pdf (accessed on 12 May 2024).

91. Ahmadian M, Roshan VD, Hosseinzadeh M. Parasympathetic reactivation in children: influence of two various modes of exercise. *Clinical Autonomic Research*. 2015; 25(4): 207-212. doi: 10.1007/s10286-015-0289-7
92. Leicht AS, Crowther RG, Gollidge J. Influence of peripheral arterial disease and supervised walking on heart rate variability. *Journal of Vascular Surgery*. 2011; 54(5): 1352-1359. doi: 10.1016/j.jvs.2011.05.027
93. Mukaka MM. A guide to appropriate use of correlation coefficient in medical research. *Malawi medical journal*. 2012; 24(3): 69-71.
94. Koo TK, Li MY. A Guideline of Selecting and Reporting Intraclass Correlation Coefficients for Reliability Research. *Journal of Chiropractic Medicine*. 2016; 15(2): 155-163. doi: 10.1016/j.jcm.2016.02.012
95. Bland JM, Altman DG, Warner DS. Agreed Statistics. *Anesthesiology*. 2012; 116(1): 182-185. doi: 10.1097/aln.0b013e31823d7784
96. Hinde K, White G, Armstrong N. Wearable Devices Suitable for Monitoring Twenty Four Hour Heart Rate Variability in Military Populations. *Sensors*. 2021; 21(4): 1061. doi: 10.3390/s21041061
97. Mejía-Mejía E, May JM, Torres R, et al. Pulse rate variability in cardiovascular health: a review on its applications and relationship with heart rate variability. *Physiological Measurement*. 2020; 41(7): 07TR01. doi: 10.1088/1361-6579/ab998c
98. Lehrer P. Applied psychophysiology: Beyond the boundaries of biofeedback (mending a wall, a brief history of our field, and applications to control of the muscles and cardiorespiratory systems). *Applied Psychophysiology and biofeedback*. 2003; 28(4): 291-304. doi: 10.1023/A:1027330909265
99. Vaillancourt DE, Newell KM. Changing complexity in human behavior and physiology through aging and disease. *Neurobiology of aging*. 2002; 23(1): 1-11. doi: 10.1016/S0197-4580(01)00247-0
100. de Groot M, Myers BA, Stump T, et al. 690-P: Relationship of Diabetes Distress and Symptoms of Autonomic Nervous System Dysregulation in Type 1 Diabetes Adults. *Diabetes*. 2024; 73. doi: 10.2337/db24-690-p
101. Takase B. Role of Heart Rate Variability in Non-Invasive Electrophysiology: Prognostic Markers of Cardiovascular Disease. *Journal of Arrhythmia*. 2010; 26(4): 227-237. doi: 10.1016/S1880-4276(10)80021-3
102. Xhyheri B, Manfrini O, Mazzolini M, et al. Heart Rate Variability Today. *Progress in Cardiovascular Diseases*. 2012; 55(3): 321-331. doi: 10.1016/j.pcad.2012.09.001
103. Stein PK, Domitrovich PP, Huikuri HV, et al. Traditional and Nonlinear Heart Rate Variability Are Each Independently Associated with Mortality after Myocardial Infarction. *Journal of Cardiovascular Electrophysiology*. 2005; 16(1): 13-20. doi: 10.1046/j.1540-8167.2005.04358.x
104. Huang C, Alamili M, Rosenberg J, et al. Heart rate variability is reduced during acute uncomplicated diverticulitis. *Journal of Critical Care*. 2016; 32: 189-195. doi: 10.1016/j.jcrc.2015.12.006
105. Ahmad S, Ramsay T, Huebsch L, et al. Continuous Multi-Parameter Heart Rate Variability Analysis Heralds Onset of Sepsis in Adults. *PLoS ONE*. 2009; 4(8): e6642. doi: 10.1371/journal.pone.0006642
106. Porges SW. The polyvagal theory: phylogenetic substrates of a social nervous system. *International journal of psychophysiology*. 2001; 42(2): 123-146. doi: 10.1016/S0167-8760(01)00162-3
107. Aubert AE, Seps B, Beckers F. Heart Rate Variability in Athletes. *Sports Medicine*. 2003; 33(12): 889-919. doi: 10.2165/00007256-200333120-00003
108. Hernando D, Hernando A, Casajús JA, et al. Methodological framework for heart rate variability analysis during exercise: application to running and cycling stress testing. *Medical & Biological Engineering & Computing*. 2017; 56(5): 781-794. doi: 10.1007/s11517-017-1724-9
109. da Silva DF, Ferraro ZM, Adamo KB, et al. Endurance Running Training Individually Guided by HRV in Untrained Women. *Journal of Strength and Conditioning Research*. 2019; 33(3): 736-746. doi: 10.1519/jsc.0000000000002001
110. Lee J, Shields RK. Sympathetic Vagal Balance and Cognitive Performance in Young Adults during the NIH Cognitive Test. *Journal of Functional Morphology and Kinesiology*. 2022; 7(3): 59. doi: 10.3390/jfkm7030059
111. Droll JA, Hayhoe MM. Seeing What We Can Do. In: *Handbook of Cognitive Science*. Cambridge University Press; 2008; pp. 189-206. doi: 10.1016/b978-0-08-046616-3.00010-4
112. Forte G, Favieri F, Casagrande M. Heart Rate Variability and Cognitive Function: A Systematic Review. *Frontiers in Neuroscience*. 2019; 13. doi: 10.3389/fnins.2019.00710
113. Shah AJ, Su S, Veledar E, et al. Is Heart Rate Variability Related to Memory Performance in Middle-Aged Men? *Psychosomatic Medicine*. 2011; 73(6): 475-482. doi: 10.1097/psy.0b013e3182227d6a
114. Hansen AL, Johnsen BH, Thayer JF. Vagal influence on working memory and attention. *International journal of*

- psychophysiology. 2003; 48(3): 263-274. doi: 10.1016/S0167-8760(03)00073-4
115. Lewis MJ, Kingsley M, Short AL, et al. Rate of reduction of heart rate variability during exercise as an index of physical work capacity. *Scandinavian Journal of Medicine & Science in Sports*. 2007; 17(6): 696-702. doi: 10.1111/j.1600-0838.2006.00616.x
116. Messerotti Benvenuti S, Mennella R, Buodo G, et al. Dysphoria is associated with reduced cardiac vagal withdrawal during the imagery of pleasant scripts: Evidence for the positive attenuation hypothesis. *Biological Psychology*. 2015; 106: 28-38. doi: 10.1016/j.biopsycho.2014.11.017
117. Schmaußer M, Laborde S. Tonic and phasic cardiac vagal activity predict cognitive-affective processing in an emotional stop-signal task. *International Journal of Psychophysiology*. 2023; 191: 9-18. doi: 10.1016/j.ijpsycho.2023.06.008
118. Pessoa L. How do emotion and motivation direct executive control? *Trends in Cognitive Sciences*. 2009; 13(4): 160-166. doi: 10.1016/j.tics.2009.01.006
119. Smith R, Thayer JF, Khalsa SS, et al. The hierarchical basis of neurovisceral integration. *Neuroscience & Biobehavioral Reviews*. 2017; 75: 274-296. doi: 10.1016/j.neubiorev.2017.02.003
120. Karmali SN, Sciusco A, May SM, et al. Heart rate variability in critical care medicine: a systematic review. *Intensive Care Medicine Experimental*. 2017; 5(1). doi: 10.1186/s40635-017-0146-1
121. Quintana DS. Statistical considerations for reporting and planning heart rate variability case-control studies. *Psychophysiology*. 2016; 54(3): 344-349. doi: 10.1111/psyp.12798
122. Bruyne MCD, Kors JA, Hoes AW, et al. Both Decreased and Increased Heart Rate Variability on the Standard 10-Second Electrocardiogram Predict Cardiac Mortality in the Elderly: The Rotterdam Study. *American Journal of Epidemiology*. 1999; 150(12): 1282-1288. doi: 10.1093/oxfordjournals.aje.a009959
123. Lee D, Baek JH, Cho YJ, et al. Association of Resting Heart Rate and Heart Rate Variability with Proximal Suicidal Risk in Patients with Diverse Psychiatric Diagnoses. *Frontiers in Psychiatry*. 2021; 12. doi: 10.3389/fpsyt.2021.652340

Article

Harnessing artificial intelligence (AI) for cybersecurity: Challenges, opportunities, risks, future directions

Zarif Bin Akhtar^{1,*}, Ahmed Tajbiul Rawol²

¹ Department of Computing, Institute of Electrical and Electronics Engineers (IEEE), Piscataway, NJ 08854, USA

² Department of Computer Science, Faculty of Science and Technology, American International University-Bangladesh (AIUB), Dhaka 1229, Bangladesh

* **Corresponding author:** Zarif Bin Akhtar, zarifbinakhtarg@gmail.com, zarifbinakhtar@ieee.org

CITATION

Akhtar ZB, Rawol AT. Harnessing artificial intelligence (AI) for cybersecurity: Challenges, opportunities, risks, future directions. *Computing and Artificial Intelligence*. 2024; 2(2): 1485. <https://doi.org/10.59400/cai.v2i2.1485>

ARTICLE INFO

Received: 28 June 2024

Accepted: 29 September 2024

Available online: 10 October 2024

COPYRIGHT



Copyright © 2024 by author(s).
Computing and Artificial Intelligence
is published by Academic Publishing
Pte. Ltd. This work is licensed under
the Creative Commons Attribution
(CC BY) license.
<https://creativecommons.org/licenses/by/4.0/>

Abstract: The integration of artificial intelligence (AI) into cybersecurity has brought about transformative advancements in threat detection and mitigation, yet it also introduces new vulnerabilities and potential threats. This research exploration systematically investigates the critical issues surrounding AI within cybersecurity, focusing on specific vulnerabilities and the potential for AI systems to be exploited by malicious actors. The research aims to address these challenges by swotting and analyzing existing methodologies designed to mitigate such risks. Through a detailed exploration of modern scientific research, this manuscript identifies the dual-edged impact of AI on cybersecurity, emphasizing both the opportunities and the dangers. The findings highlight the need for strategic solutions that not only enhance digital security and user privacy but also address the ethical and regulatory aspects of AI in cybersecurity. Key contributions include a comprehensive analysis of emerging trends, challenges, and the development of AI-driven cybersecurity frameworks. The research also provides actionable recommendations for the future development of robust, reliable, and secure AI-based systems, bridging current knowledge gaps and offering valuable insights for academia and industry alike.

Keywords: artificial intelligence (AI); cybersecurity; data informatics; cybersecurity; deep learning (DL); machine learning (ML); security informatics; security vulnerabilities; privacy

1. Introduction

The integration of artificial intelligence (AI) into cybersecurity has brought significant advancements, offering powerful tools for detecting, preventing, and mitigating cyber threats. By employing sophisticated algorithms, AI enhances the ability to identify emerging threats, such as malware and ransomware attacks, through behavior pattern recognition [1–3]. AI systems provide real-time intelligence on global and industry-specific dangers, enabling organizations to prioritize their security measures more effectively.

Furthermore, AI plays a critical role in combating automated threats, such as bots, which represent a substantial portion of internet traffic. Through the analysis of website traffic, AI distinguishes between legitimate bots, malicious bots, and human users, thereby enabling cybersecurity teams to stay ahead of these automated threats.

AI also enhances breach risk prediction by analyzing IT asset inventories and threat exposure data to identify vulnerable areas within an organization's digital infrastructure [4–6]. This predictive capability allows for better resource allocation and more effective protection against potential breaches.

Moreover, in the era of remote work, AI is essential for endpoint protection,

moving beyond traditional signature-based approaches to establish behavioral baselines for endpoints and proactively identifying anomalies that may indicate emerging threats. The efficiency and accuracy brought by AI are transforming the cybersecurity landscape. AI automates routine tasks, allowing security analysts to focus on more complex responsibilities, such as incident response and threat hunting. By rapidly analyzing vast amounts of security data, AI detects patterns and anomalies that may signal cyber threats, thus improving the speed and accuracy of threat detection [7–9]. This automation extends to vulnerability scanning, patch management, and incident investigation, streamlining cybersecurity operations and enhancing overall effectiveness. AI's impact on cybersecurity also includes significant cost reductions. By automating routine processes, organizations can decrease the workload and associated costs of human resources [10–12]. Additionally, AI's accuracy in threat detection helps avoid the expenses related to false alarms or undetected breaches. Enhanced incident response capabilities reduce the time needed to remediate security incidents, minimizing potential financial losses, reputational damage, and regulatory penalties [13–15]. The proactive threat intelligence provided by AI further contributes to cost reduction by enabling timely and actionable insights that prevent and mitigate security incidents. The ability of AI to process data rapidly is crucial in the fast-paced cyber threat landscape, where real-time threat detection and response are essential. AI systems continuously learn and adapt, ensuring that organizations can proactively defend against emerging threats and respond effectively to minimize the impact of cyber-attacks. The scalability of AI algorithms allows for the handling and analysis of vast amounts of data, including network traffic logs, system logs, user behaviors, and threat intelligence feeds. This scalability optimizes resource allocation, improves operational efficiency, and enables organizations to process and detect cyber threats in complex and dynamic environments.

This research aims to explore the dual-edged nature of AI in cybersecurity, identifying both the opportunities and the risks associated with its integration. By providing a systematic analysis of existing methodologies and proposing strategic solutions, this study seeks to contribute to the development of robust AI-driven cybersecurity frameworks that enhance digital security and user privacy. Through a comprehensive analysis of emerging trends, challenges, and regulatory frameworks, this manuscript addresses the critical need for effective and ethical AI applications in cybersecurity.

2. Methods and experimental analysis

This research adopts a structured and systematic approach to investigate the impact of artificial intelligence (AI) on cybersecurity and privacy. The research exploration begins with an extensive appraisal of existing background and available knowledge to establish a comprehensive background and identify critical research gaps in the intersection of AI and cybersecurity. This includes sourcing and analyzing a wide range of academic papers, industry reports, and case studies relevant to the topic.

The background research not only provides the necessary context but also highlights the current state of AI integration in cybersecurity, guiding the formulation

of the research objectives and questions. To address these research questions, various data collection methods are employed, including surveys and interviews with cybersecurity experts, as well as the analysis of existing datasets. These methods are chosen to gather diverse perspectives and data, ensuring a robust understanding of the challenges and opportunities associated with AI in cybersecurity. The collected data undergo rigorous pre-processing, including cleaning and normalization, to ensure its quality, relevance, and applicability to the research objectives.

The research is anchored by specific assumptions and research questions designed to explore the vulnerabilities and threats posed by AI integration in cybersecurity. These questions guide the investigation of both the risks and potential solutions offered by AI technologies. The research evaluates the performance of AI techniques in cybersecurity using appropriate metrics such as accuracy, precision, recall, and F1 score. These metrics are crucial for assessing the effectiveness of AI-driven approaches in comparison to traditional cybersecurity methods, providing a quantitative basis for the analysis.

Data visualization tools and techniques are employed to clearly and effectively present the research findings. Charts, graphs, and heatmaps are used to convey complex data insights in an understandable manner, aiding in the interpretation and communication of the results. The research also includes a comparative analysis between AI-based techniques and traditional cybersecurity methods. This comparison is essential to highlight the improvements AI offers and to identify specific areas where AI can provide significant advantages over conventional approaches. This comparative analysis helps in understanding the practical implications of integrating AI into existing cybersecurity frameworks.

The results of the research investigations are analyzed in the context of the research objectives, providing insights into the impact of AI on cybersecurity and privacy. This includes a discussion of the implications of the findings for future cybersecurity practices and AI advancements. The research concludes by summarizing the key findings, acknowledging the limitations of the research exploration, and offering recommendations for further research. These recommendations focus on the continued exploration of AI applications in cybersecurity and the enhancement of privacy protections in the digital world. This structured methodology ensures a thorough exploration of AI's role in enhancing cybersecurity and addressing privacy concerns, contributing to the development of robust and secure AI-driven frameworks in the field.

2.1. Background research and available knowledge explorations

Before we get into all the nitty-gritty within the retrospect, which complexifies concerning the context, let's first learn the basics and history of its foundations. Computer security, also known as cybersecurity, digital security, or IT security, is the practice of protecting computer systems and networks from malicious attacks that can lead to unauthorized access, theft, or damage of hardware, software, or data, as well as disruption of services [1–5].

With the increasing reliance on computer systems, the internet, and wireless networks, cybersecurity has become a critical challenge in today's interconnected

world. The history of cybersecurity can be traced back to the emergence of the internet and the digital transformation of society [6,7]. In the 1970s and 1980s, computer security primarily focused on academic settings until the advent of the internet, which brought about an increase in connectivity and the rise of computer viruses and network intrusions. The institutionalization of cyber threats and cybersecurity occurred in the 2000s. The field of computer security was significantly influenced by the April 1967 session organized by Willis Ware at the Spring Joint Computer Conference, known as the Ware Report. This event and subsequent publication marked foundational moments in the history of computer security. The report addressed material, cultural, political, and social concerns related to computer security. In the 1970s and 1980s, computer threats were relatively limited as the technology was still in its early stages, and security breaches were easily identifiable. However, insider threats, such as unauthorized access to sensitive information by malicious insiders, were more prevalent [8–10]. During this time, computer firms like IBM started offering commercial access control systems and security software products. Notable incidents in the history of cybersecurity include the creation of the computer worm Creeper in 1971, the first documented case of cyber espionage performed by German hackers in the late 1980s, and the distribution of the Morris worm in 1988, which gained significant media attention. The development of secure protocols, such as SSL (Secure Sockets Layer), by Netscape in the mid-1990s aimed to enhance the security of online communications. However, even these early versions had vulnerabilities that were later addressed in subsequent releases [11–15].

The role of government agencies, such as the National Security Agency (NSA), in cybersecurity is significant. The NSA is responsible for protecting U.S. information systems and collecting foreign intelligence. The agency analyses software for security flaws, often using them offensively rather than reporting them to software producers for remediation. This approach has led to the exploitation of security vulnerabilities by both allies and adversaries, contributing to the emergence of cyberwarfare capabilities worldwide. The history of cybersecurity reflects the evolution of computer systems, the internet, and the growing threats associated with them. From the early days of computer viruses and network intrusions to the rise of cyber espionage and the development of secure protocols, the field of cybersecurity has become essential for protecting information systems and mitigating potential risks. The involvement of government agencies and the constant interplay between security measures and emerging threats continue to shape the landscape of cybersecurity [16–22]. The history of artificial intelligence (AI) can be traced back to ancient times, where myths and stories depicted the creation of artificial beings with intelligence or consciousness. However, the modern foundations of AI were established by philosophers who sought to understand human thinking as a mechanistic process involving the manipulation of symbols. This line of thinking eventually led to the invention of the programmable digital computer in the 1940s, which sparked the serious exploration of building an electronic brain [23,24].

The field of AI research was officially launched in the summer of 1956 at a workshop held at Dartmouth College. The participants of this workshop, who would become influential figures in AI research, were optimistic about achieving human-

level intelligence in machines within a generation. Substantial funding was provided to support their efforts [25,26]. However, as the project progressed, it became evident that the challenges of developing AI were far greater than initially anticipated. Critics, such as James Lighthill, voiced concerns, and the U.S. and British governments responded by ceasing funding for undirected AI research in 1974. This marked the beginning of a timeline period known as the “AI winter,” characterized by a decline in AI research and disillusionment with its progress [27–29]. In the early 1980s, the Japanese government initiated a visionary initiative that renewed interest and investment in AI, leading to substantial funding from governments and industry. Moreover, by the late 1980s, investors once again became disillusioned with the progress of AI, and funding was withdrawn. In the first decades of the 21st century, AI experienced a resurgence in its investment and interest [30–32].

This was possible by advancements within machine learning techniques, the availability of powerful computer hardware, and the accumulation of vast amounts of data. Machine learning, in particular, demonstrated success in various academic and industrial applications, leading to a renewed optimism and enthusiasm for AI [33–35].

The history of AI has been marked by periods of optimism, followed by periods of disappointment and all the reduced funding. However, recent advancements have sparked a new wave of excitement, with AI becoming increasingly integrated into various aspects of our lives, from personal assistants to autonomous vehicles, and opening up new possibilities for the near future.

2.2. Cyberthreats and the information security domain

Cyber threats have seen a significant increase in recent years, with the proliferation of technology and interconnected systems. The COVID-19 pandemic further accelerated this trend, resulting in a 600% surge within cybercrime since 2020. The impact of cyberattacks is wide-ranging, affecting nearly every industry and leading to major financial losses, reputational damage, legal liabilities, productivity disruptions, and business continuity issues. Estimates indicate that global cybercrime costs could reach \$10.5 trillion by 2025, highlighting the severity of the problem. Data breaches are a prevalent and costly consequence of cyber threats. In 2022–2023, the global average cost of a data breach was \$4.35 million, with the United States recording the highest average cost at \$9.44 million. The healthcare industry experienced a significant jump in data breach costs, with an average of \$10.1 million, reflecting a 42% increase since 2020. Cloud environments were also a common target, accounting for 45% of data breaches in 2022–2024. Various motives drive cyber threats. Cybercrime committed for financial gain by individuals or groups is one prevalent motive. Politically motivated cyberattacks seek to disrupt systems or gather sensitive information, while cyberterrorism aims to undermine electronic systems and impose fear or panic. Malware, a broad category of malicious software, poses a significant threat.

Viruses, trojans, spyware, adware, botnets, and ransomware are among the different types of malware used by attackers to gain unauthorized access, disrupt operations, or extort victims. Ransomware attacks have grown in prominence, with organizations facing the threat of permanent data loss unless they pay a ransom, often

in cryptocurrencies. Phishing attacks, where cybercriminals deceive the victims into divulging sensitive information, are another widespread method used.

Other types of cyber threats also include distributed denial-of-service (DDoS) attacks, where a network is overloaded by coordinating a large number of systems; man-in-the-middle attacks, which intercept and steal data during communication; SQL injection, exploiting vulnerabilities in data-driven applications to access sensitive information; insider threats from individuals with authorized access to systems; advanced persistent threats (APTs), infiltrations that remain undetected over an extended period for data theft; and especially crypto jacking, where victims' computing resources are hijacked for cryptocurrency mining. Data security plays a crucial role in combating cyber threats. It encompasses measures to protect data from unauthorized access, corruption, or accidental errors.

This technique includes data privacy, encryption techniques such as cryptography and homomorphic encryption, and ensuring data integrity. Addressing cyber threats requires continuous vigilance, robust cybersecurity measures, and proactive strategies. Organizations must invest in cybersecurity infrastructure, employee training, threat detection and response systems, and data protection mechanisms to mitigate risks and safeguard sensitive information in an increasingly interconnected digital landscape.

Cybersecurity is a critical practice aimed at safeguarding electronic systems, networks, computers, mobile devices, programs, and data from malicious digital attacks. It involves the protection of digital information and infrastructure to prevent unauthorized access, data breaches, and disruption of business processes. To achieve cybersecurity, an organization typically implements an infrastructure consisting of three key components: IT security, cyber security, and network security. IT security, also known as electronic information security, focuses on protecting both physical and digital data from intruders. It safeguards data at rest and in transit, ensuring its integrity and confidentiality. Cybersecurity is a subset of IT security and specifically focuses on safeguarding digital data on networks, computers, and devices from unauthorized access, attack, and destruction. It involves measures such as firewalls, encryption, intrusion detection systems, and incident response protocols to prevent cyber threats and mitigate their impact. Network security, or computer security, is a subset of cyber security and is concerned with protecting data transmitted through computers and devices in a network. It employs hardware and software solutions to ensure the secure transmission and reception of data, guarding against interception, tampering, and unauthorized access.

In practice, IT security professionals and cyber security professionals often collaborate to protect an organization's data and prevent unauthorized access. While some companies employ separate professionals for IT security and cyber security, the roles may overlap, with cyber security professionals primarily focusing on securing digital data across various networks and systems. It's important to note that cyber security is a part of the broader field of information security [31,32].

Information security encompasses the main protection of data and information and information systems across different realms, including the physical world. As anything occurring in the cyber realm involves the protection of information and systems, information security can be seen as a superset that encompasses cyber

security. Cybersecurity plays a crucial role in safeguarding digital assets and ensuring the privacy, integrity, and availability of data in an increasingly interconnected and digitized world [31,32]. It requires proactive measures, ongoing monitoring, and the adoption of robust security practices to mitigate risks and effectively respond to cyber threats. To provide an overview impression, **Figure 1** illustration is epitomized concerning the matter.



Figure 1. A diagram of information security (cyber security and network security).

2.3. Cybercrimes, privacy, security vulnerabilities

According to Forbes, 76% of enterprises have prioritized AI and machine learning in their IT budgets, driven by the increasing volume of data that needs to be analyzed to identify and mitigate cyber threats. AI is becoming an essential tool in the fight against cybercrime.

The rapid acceleration of cybercrime has been facilitated by the lower barrier to entry for malicious actors, who have evolved their business models to include subscription services and starter kits. Additionally, the use of large language models (LLMs) like ChatGPT to write malicious code highlights the potential challenges to cybersecurity. However, it is crucial for business leaders in today's digital world to be knowledgeable about the developments of AI in cybersecurity.

Blackberry's research found that the majority of IT decision-makers plan to invest in AI-driven cybersecurity, recognizing its potential to enhance their defenses against cyber threats. While there are concerns about the misuse of AI, particularly in social

engineering and skilling up less experienced hackers, the actual threat posed by AI-generated code may not be as significant as some headlines suggest.

While AI can generate code that gets close to completion, it often requires human intelligence and refinement to make it fully functional. This means that the last mile of human intervention is crucial, reducing the potential threat. It is important to acknowledge that AI can also be used to help protect against cyber threats. AI has the ability to make inferences, recognize patterns, and perform proactive actions to shield against online threats. It can automate incident response, streamline threat hunting, and analyse large amounts of data to improve cybersecurity.

AI-powered tools provide continuous monitoring, real-time attack detection, and automation of incident response. They can also assist in identifying false positives and strengthening access control measures. Furthermore, AI can help mitigate insider threats by analyzing user behavior and, at the same time, identifying employees engaged in malicious activities.

By leveraging AI in cybersecurity, organizations can improve their threat detection, response times, and overall security posture. While there are many benefits to using AI in cybersecurity, there are also potential risks that must be considered. Bias in AI algorithms can lead to flawed decisions or missed threats if the training data is biased or unrepresentative. Addressing bias requires diverse and representative training data, pre-processing techniques, ongoing monitoring, transparency, and continuous education.

Attackers can leverage AI technologies to enhance the effectiveness of their cyberattacks. AI can be used to create highly convincing phishing emails, develop advanced evasion techniques, automate attack tools, facilitate deepfake attacks, and execute adversarial attacks. These malicious uses of AI pose significant challenges for defensive measures and necessitate robust cybersecurity strategies. To be precise, business leaders must recognize the potential dangers and benefits of using AI in cybersecurity. While there are risks associated with the misuse of AI, efforts can be made to address bias and ensure fairness and equity. AI can be harnessed to improve cybersecurity by automating tasks, providing continuous monitoring, enhancing threat detection, and mitigating insider threats.

By embracing AI responsibly, organizations can strengthen their security defenses in the face of evolving cyber threats. AI-powered security solutions, like any software or system, can have vulnerabilities that attackers may exploit. These vulnerabilities can compromise the effectiveness of cybersecurity measures.

To mitigate these risks, organizations should regularly assess the security of AI systems through penetration testing and simulations of real-world attacks. Secure development practices should be followed from the early stages, including adhering to coding standards, conducting thorough security assessments, and using secure development frameworks and tools.

Secure deployment and configuration practices are crucial, involving proper access controls, secure storage of sensitive data, and implementation of secure communication protocols. Regular updates and patching should be performed to address known vulnerabilities. Ongoing monitoring, robust logging, and incident response plans are necessary to detect and respond to security incidents promptly.

When adopting AI systems from third-party vendors, thorough security evaluations should be conducted to ensure secure development practices and strong security measures.

However, there are challenges to implementing AI in security. Lack of transparency and interpretability is a common issue, as AI systems often function as black boxes, making it challenging to understand how decisions are made. Bias and fairness concerns arise when AI systems replicate biases present in the training data. Integration with existing security systems can be problematic if AI-powered solutions do not effectively work alongside other tools in an organization's security architecture.

To be more accurate, organizations need to address security vulnerabilities in AI systems through regular assessments, secure development practices, proper deployment and configuration, ongoing monitoring, and vendor evaluations. They must also consider challenges such as lack of transparency, bias, and integration with existing security systems when implementing AI in security. By addressing these concerns, organizations can enhance the effectiveness and reliability of their AI-powered security solutions.

Vulnerabilities are weaknesses in a computer system, either in the hardware or software, that compromise the overall security of the entire system. These vulnerabilities can be exploited by threat actors, such as attackers, to gain unauthorized access or perform malicious actions within the system. Vulnerabilities are sometimes also referred to as the attack surface, as they provide opportunities for attackers to breach the system's defenses. Vulnerability management is a cyclical practice aimed at identifying, assessing, and addressing vulnerabilities in computing systems.

The process typically involves discovering all assets within a system, prioritizing them based on their criticality, conducting vulnerability scans or assessments, reporting on the findings, remediating the identified vulnerabilities, and verifying the effectiveness of the remediation efforts. This iterative process helps organizations stay proactive in addressing vulnerabilities and minimizing the risk of successful attacks.

It is also very important to differentiate between vulnerabilities and security risks. While vulnerabilities represent potential weaknesses, security risks refer to the potential impact or harm that can result from the exploitation of vulnerabilities. A vulnerability becomes a security risk when there is a significant potential for damage or compromise. However, not all vulnerabilities pose a risk, particularly when the affected asset has no value or the vulnerability is not easily exploitable.

An exploitable vulnerability is one that has known instances of successful attacks. The window of vulnerability refers to the time period starting from when a security hole is introduced or discovered in deployed software until it is patched or mitigated, or when the attacker's access is removed. Zero-day attacks, where vulnerabilities are exploited before a fix is available, represent a particularly challenging type of vulnerability. It is worth noting that vulnerabilities are not limited to software. Hardware, physical site vulnerabilities, or weaknesses in personnel practices can also introduce vulnerabilities in a system. Additionally, certain constructs in programming languages that are complex or difficult to use properly can lead to a very large number of vulnerabilities if not implemented correctly. To put it simply, understanding and managing vulnerabilities is crucial for maintaining the security of computer systems.

By actively identifying and addressing vulnerabilities, organizations can enhance their defense against potential attacks and reduce the likelihood of security breaches.

2.4. The abuse of AI within the realm of cybersecurity

Cybercriminals are finding ways to exploit AI for their malicious activities. One method is through social engineering schemes, where AI automates the processes and allows for more personalized and sophisticated messaging to deceive victims. This leads to a higher success rate for cybercriminals in carrying out phishing, vishing, and business email compromise scams.

AI is being used to enhance password hacking algorithms, enabling hackers to decipher passwords more quickly and accurately, emphasizing the need for strong password security measures. Another concerning use of AI by hackers is the creation of deepfakes, which involves manipulating visual or audio content to impersonate individuals and spread deceptive information.

Deepfakes can be combined with social engineering, extortion, and other schemes to cause confusion and fear among those who consume the manipulated content. Furthermore, hackers can employ data poisoning techniques to alter the training data of AI algorithms, leading to biased or incorrect decisions. Data poisoning can be difficult to detect and can result in severe consequences by the time it is discovered.

In this changing AI environment, individuals and businesses need to review their cybersecurity practices and ensure they follow best practices, especially in areas such as passwords, data privacy, personal cybersecurity, and protection against social engineering. Regularly updating the security measures and always staying informed about the latest cyber-security tips is crucial. While AI offers many benefits in improving cybersecurity, it is important to remain vigilant and adapt security practices to mitigate the risks associated with AI-powered attacks.

One challenge in using AI for cybersecurity is the need for substantial resources and financial investments to build and maintain AI systems effectively. Acquiring diverse and reliable datasets for training AI systems can be time-consuming and costly, making it difficult for many organizations to afford. Inaccurate or incomplete datasets can also lead to incorrect results and false positives, highlighting the great importance of quality data for AI systems to function effectively.

Furthermore, the same AI technologies used for defense can also be leveraged by cybercriminals to analyze their malware and launch more advanced attacks. This highlights the ongoing cat-and-mouse game between cybersecurity professionals and hackers, where advancements in AI technology are utilized on both sides. In other words, while AI has the potential to enhance cybersecurity, it is important to be aware of the ways in which hackers can abuse AI for their malicious purposes.

Implementing robust cybersecurity measures, staying informed about the evolving AI landscape, and adapting security practices accordingly are very crucial for individuals and organizations to protect themselves in this changing environment. Malware and phishing attacks are significant cybersecurity threats that can cause substantial harm to individuals and organizations.

However, the advancements in artificial intelligence (AI) have brought new

possibilities for detecting and mitigating these threats. AI-based cybersecurity systems have shown promising results in malware detection [31–35]. Traditional signature-based approaches can only detect known malware, while AI-powered systems can identify dynamically changing malicious agents more effectively. By utilizing techniques like computer vision and neural networks, researchers have achieved high accuracy in detecting malware across various file formats.

AI systems can analyze the inherent characteristics of malware to identify potential threats, improving the overall security efficiency compared to legacy detection systems. Phishing attacks, which often lead to the activation of malware, can also be combated using AI. Machine learning-based techniques can analyze the structure of emails and classify them as legitimate or phishing emails, achieving high accuracy rates. AI-enabled tools, such as Mimecast's Cyber Graph, employ machine learning to block trackers, detect phishing emails, and alert users about potential threats.

AI's role within cybersecurity goes beyond malware and phishing detection. It helps in knowledge consolidation by leveraging machine learning models to retain and utilize vast amounts of historical data to detect security breaches effectively. AI can keep track of global and industry-specific vulnerabilities, constantly updating its knowledge to defend against new threat actors and prevent upcoming attacks.

Tech giants like Google, IBM, and Microsoft have invested significant resources in developing advanced AI systems for threat identification and mitigation, making substantial progress in protecting users and enterprises [31–35]. Additionally, AI tools can predict breach risks, prioritize security measures, and automate threat detection and mitigation processes. By reducing the time taken to detect and respond to cyber threats, AI contributes to minimizing the damage caused by attacks. It enables organizations to allocate resources more effectively and develop cyber resilience to withstand future attacks.

While AI offers tremendous potential for improving cybersecurity, it also poses certain risks and challenges. Data manipulation, where hackers alter training data or introduce biases, can impact the efficiency of AI models. Hackers themselves can exploit AI techniques to develop intelligent malware that evades detection. Insufficient or biased training data can result in false positives or a false sense of security.

Privacy concerns arise when user data is used to train AI models without adequate protection. Moreover, AI systems themselves can become targets of cyberattacks, with hackers feeding poisonous data to manipulate their behavior.

To address these challenges, it is crucial to build robust infrastructures that counter the risks associated with AI in cybersecurity. Data integrity and privacy protection measures, continuous model updating, and proactive security measures are essential for ensuring the safe and secure operation of AI-powered cybersecurity systems. AI brings significant advancements to malware and phishing detection, knowledge consolidation, threat prediction, and automation in cybersecurity. While there are risks and challenges to overcome, organizations must leverage AI's potential while implementing robust security measures to create a safe digital environment. By combining human expertise with AI capabilities, the cybersecurity landscape can be strengthened to defend against evolving threats and ensure the protection of

individuals and businesses.

Artificial intelligence (AI) has been adopted by several tech giants and cybersecurity companies to enhance their capabilities in the field. Google has been utilizing machine learning techniques in Gmail and various other services for years, with deep learning algorithms allowing for independent adjustments and self-regulation [31–35]. IBM heavily relies on its Watson cognitive learning platform for tasks like knowledge consolidation and threat detection, aiming to automate routine processes in security operations centralized areas.

Juniper Networks envisions a future with autonomous networks, leveraging AI, machine learning, and intent-driven networking. Balbix Security Cloud also uses AI-powered risk predictions and vulnerability management to bolster cyber-security efforts. However, the rise of AI in cybersecurity also presents risks. Adversaries can employ AI and ML techniques to evade defenses and launch more sophisticated attacks. They can target the data used to train security algorithms, manipulate information, or develop mutating malware to avoid detection. It is crucial for organizations to be aware of these downsides and implement safeguards to protect against potential threats.

3. Cybersecurity vulnerabilities: Case studies analysis

Server-Side Template Injection (SSTI) and Client-Side Template Injection (CSTI) are significant security vulnerabilities that occur when attackers are able to inject and execute malicious code within template engines used by web applications. SSTI happens when user-provided input is improperly sanitized and subsequently incorporated into server-side templates, which are then executed by the server.

Common server-side templating engines vulnerable to such attacks include Twig, Jinja2, Django, ExpressJS, and Razor. Conversely, CSTI occurs when user input is unsanitized and injected into client-side templates, which are executed by the victim's browser. Popular client-side templating engines susceptible to CSTI include AngularJS, Vue, Handlebars, and Mustache. An illustrative case of an SSTI attack involved an application allowing users to create email templates using the Twig templating engine.

By inserting the test string `{{7×7}}` into the template, the attackers confirmed the vulnerability when the test email returned the value "49", indicating that the input was executed by the template engine. This discovery allowed the attackers to exploit Twig's 'filter' function to execute arbitrary system commands, leading to remote code execution under the context of the `www-data` user. Such an attack can have severe implications, including potential privilege escalation and unauthorized access to internal services.

To prevent SSTI attacks, it is crucial to sanitize user inputs rigorously, ensuring that no malicious code is processed by the template engine. Using template engines with built-in security measures, such as automatic input escape, strict input validation, and sandboxing, can also mitigate these risks. For CSTI, preventing attacks involves similar measures: properly validating and sanitizing user inputs, adhering to secure coding practices, conducting regular vulnerability assessments, and keeping software updated with the latest security patches. By implementing these protective strategies,

developers can significantly reduce the likelihood of both SSTI and CSTI attacks, thereby enhancing the overall security of web applications.

Next, physical social engineering tests involve a team of experts attempting to gain access to buildings and offices to evaluate the security of the infrastructure and employees. These tests are usually conducted with mature cybersecurity clients, with only a few staff members aware of the ongoing test to ensure genuine responses from employees. In a recent engagement, the challenge was to access two different sites, remain unchallenged, and gather additional information by interacting with employees.

Reconnaissance and Pretext: Effective reconnaissance, or Open-Source Intelligence (OSINT), is crucial for these tests, especially for physical engagements. Tools like Google Maps and LinkedIn were used to gather information about the building layouts, potential entry methods, and develop a convincing pretext for why they should be allowed entry. In this case, posing as contractors or consultants due to ongoing building work proved to be an effective pretext.

Initial Access: The team targeted a satellite site first, considering it easier to breach. Upon arrival, wearing Hi-Viz vests, they were easily allowed inside by a staff member. Inside, they encountered key card access restrictions but managed to find an unoccupied conference room. Here, they connected laptops to Ethernet ports and conducted network scans, even obtaining the corporate Wi-Fi password from staff members who did not question their presence.

Main Target: With the initial success boosting their confidence, the team targeted the head office next. Despite initial resistance at the reception, they eventually gained entry by convincing an employee of their legitimate presence. Inside, the lack of a proper sign-in process and the hot-desking environment allowed them to move freely and conduct further network attacks. They managed to collect user hashes and crack one belonging to a security team member. This led to discovering a Domain Admin account vulnerable to kerberoasting, eventually giving them Domain Administrator credentials and full access to the network.

Recommendations include the various following engagements; the team provided several recommendations to enhance security.

Enforce Visitor Sign-In Processes: Ensuring that all visitors follow a strict sign-in process can prevent unauthorized access.

Staff Training on Social Engineering Risks: Educating staff about the dangers of social engineering can mitigate the risk of such attacks.

Badge Security: Avoid revealing badge details on social media to prevent attackers from creating fake badges.

Monitor All Entrances: Tailgating from non-monitored entrances like smoking areas can be prevented by accounting for all entry points.

Network Security Measures: Implementing MAC filtering for Ethernet connections and securing Wi-Fi access points can prevent unauthorized network access.

These measures can significantly bolster the physical and cybersecurity posture of an organization, making it harder for social engineering attacks to succeed.

4. Results and findings

The integration of artificial intelligence (AI) into cybersecurity is increasingly critical, particularly in addressing the challenges associated with privacy and security in the digital age. As internet technologies rapidly advance and networking capabilities between devices expand, there is a growing need for robust cybersecurity measures that can handle vast and dynamic data flows. This research presents detailed visual analytical illustrations, including results and findings within **Figures 2–7**, which provide a comprehensive overview of AI’s impact on cybersecurity. These figures include conceptual frameworks for AI applications in cybersecurity, detailed security analytics, and insights into specific vulnerabilities and threats identified in the explorations.

The visualizations serve as key tools for understanding the complex interactions between AI and cybersecurity. Each figure is methodically designed to highlight various aspects of AI integration, from its potential to enhance threat detection and response capabilities to the new vulnerabilities it introduces.

For example, **Figures 2–4** demonstrate how AI can be leveraged to predict and mitigate cyber threats more effectively. However, they also reveal areas where AI systems may be susceptible to exploitation, emphasizing the need for continuous monitoring and improvement. The findings of this research underscore the significant influence of AI on digital systems, which are becoming increasingly reliant on advanced computing technologies (**Figures 5–7**). As AI continues to evolve, it is reshaping our approach to cybersecurity, driving the development of new strategies and tools to address emerging threats. The research highlights that while AI offers enhanced capabilities for threat detection, such as real-time anomaly detection and advanced behavioral analysis, it also introduces new challenges, particularly concerning the protection of sensitive information and the management of AI-specific vulnerabilities.

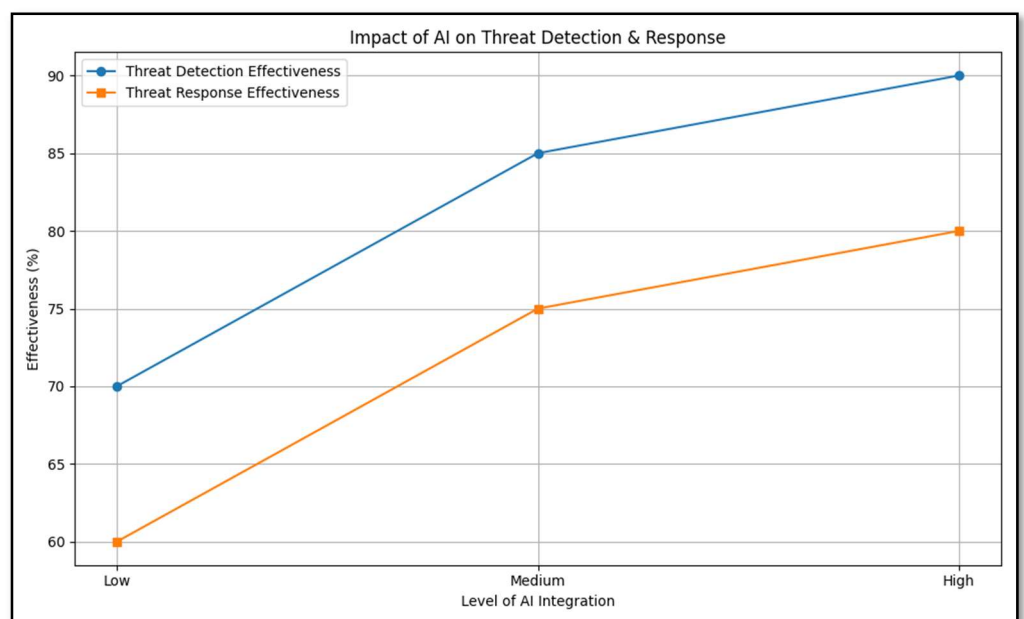


Figure 2. A visualization for impacts of AI on threat detection and response.

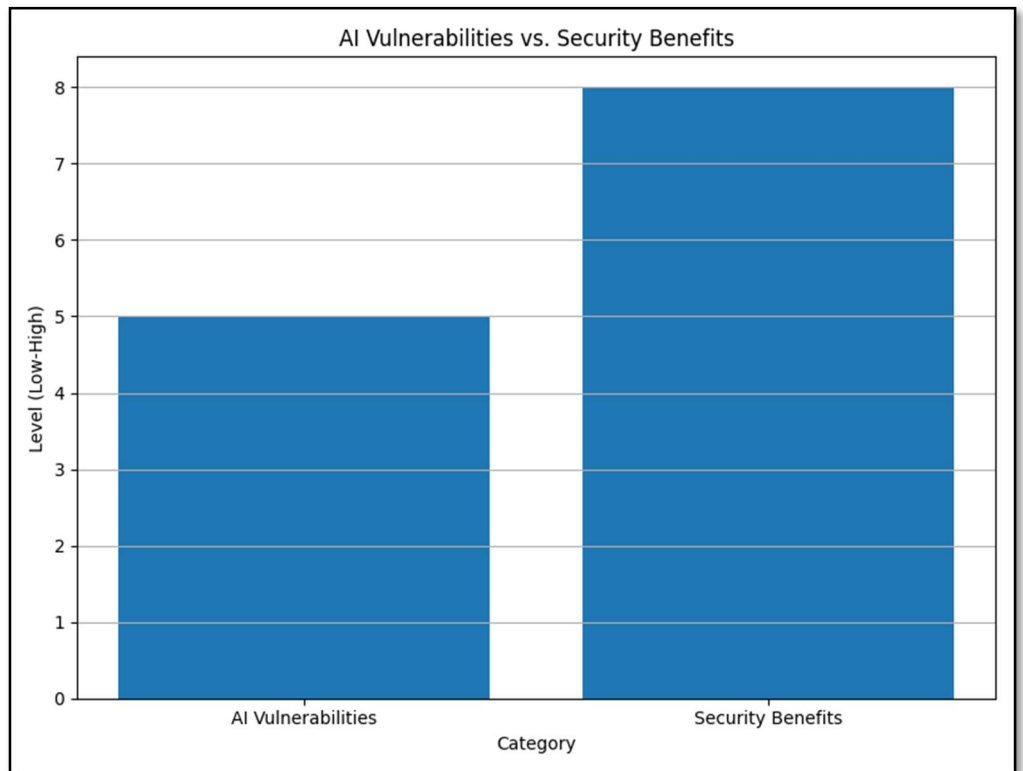


Figure 3. An overview of AI vulnerabilities vs. security benefits.

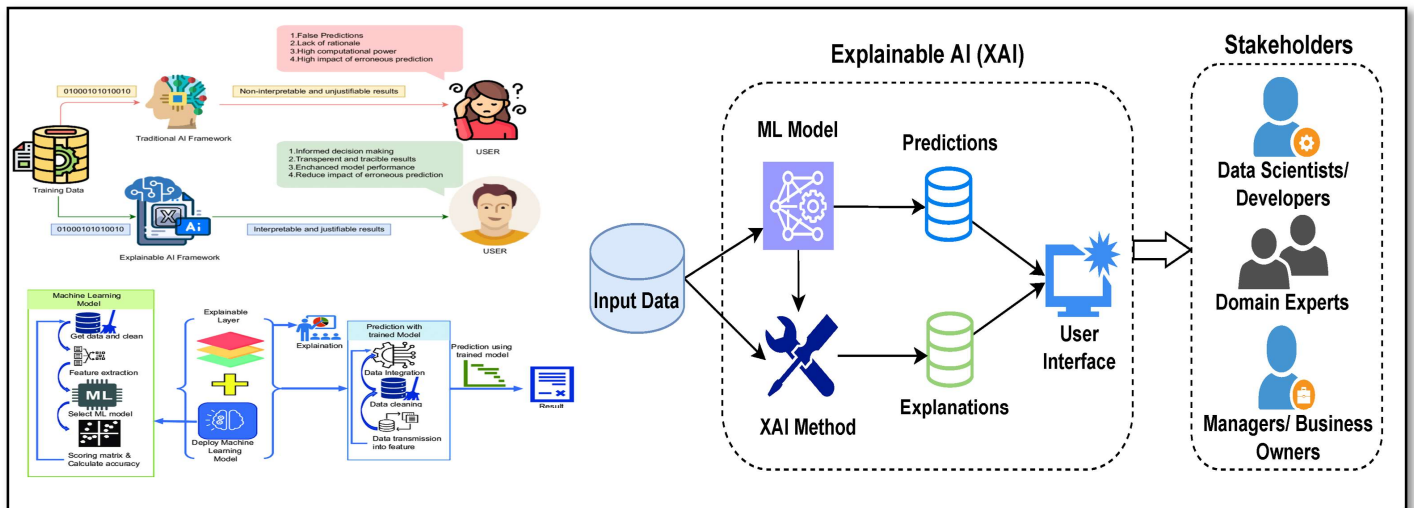


Figure 4. An overview of the AI and XAI-user's perspective context.



Figure 5. A visualization of use cases for AI in cybersecurity with the most dangerous threats.

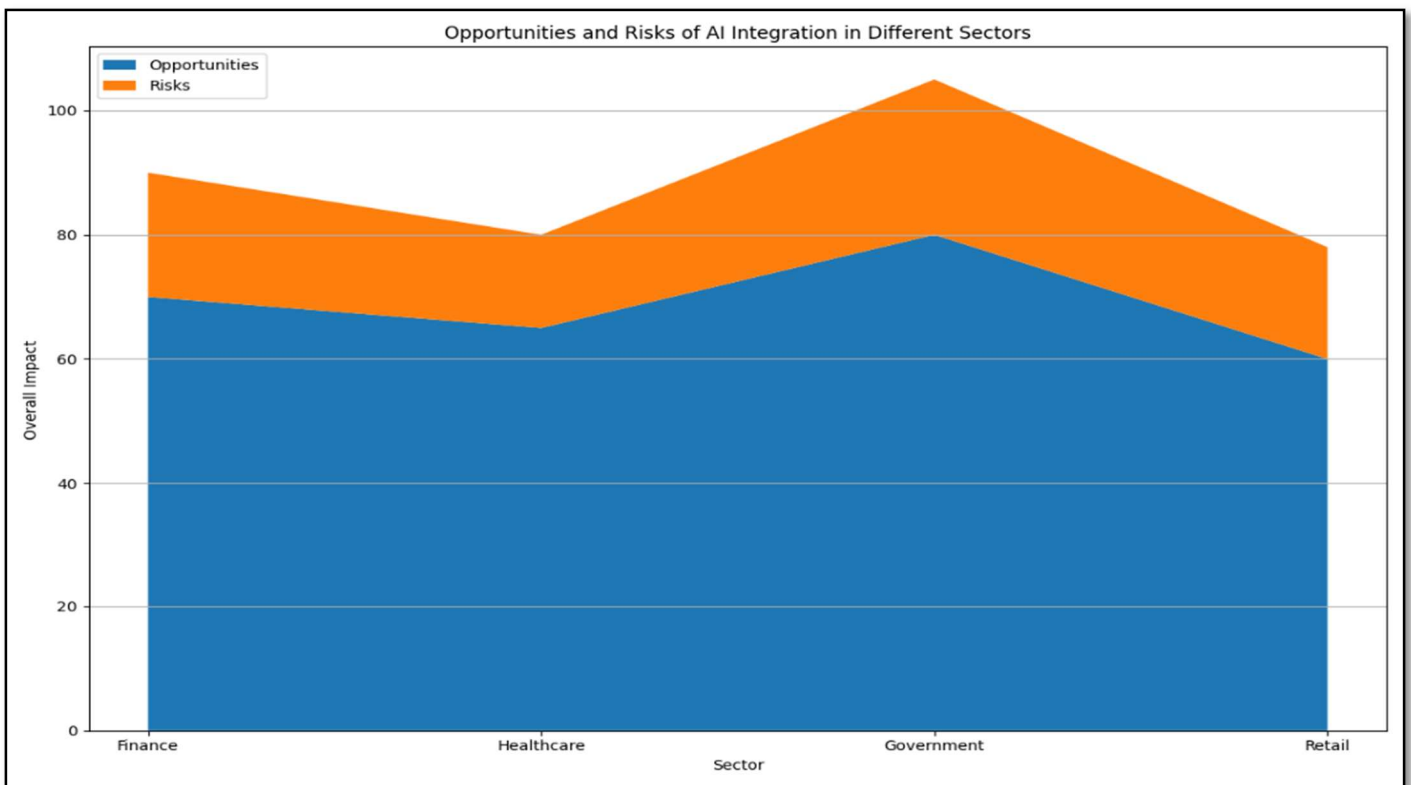


Figure 6. An overview of opportunities and risks of AI integrations in different sectors.

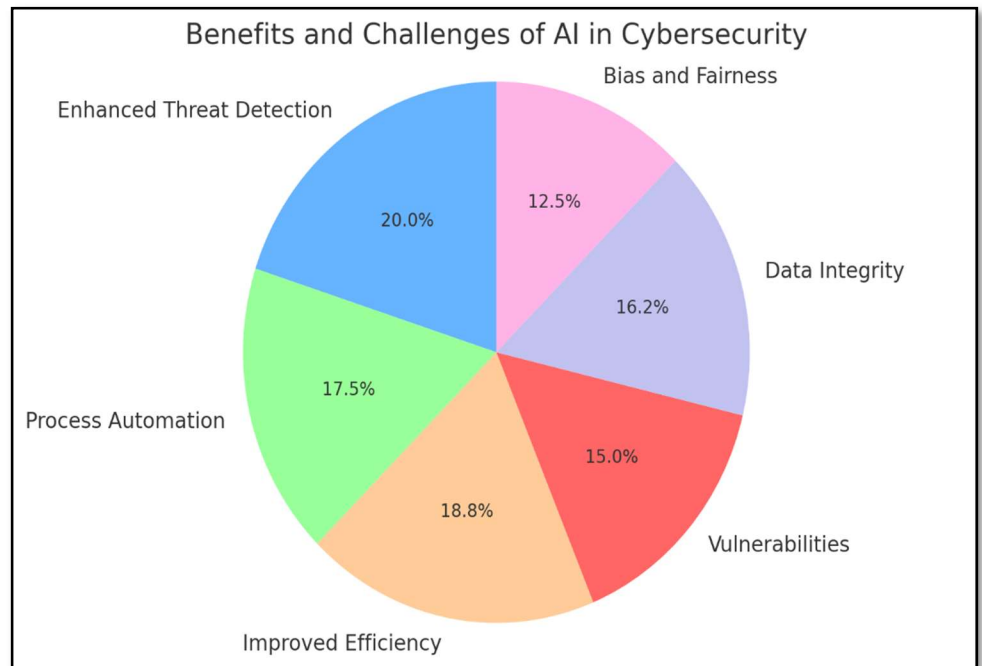


Figure 7. AI in cybersecurity benefits and challenges an overall visualization.

The expansion of AI's role in cybersecurity is expected to bring about transformative changes across various sectors, with cybersecurity remaining a critical area of focus. The research identifies both opportunities and risks associated with AI integration.

On the one hand, AI-driven systems can significantly improve the efficiency and accuracy of cybersecurity measures, offering real-time threat detection and response capabilities that were previously unattainable. On the other hand, the increasing complexity and sophistication of AI systems introduce new vulnerabilities that must be carefully managed to prevent potential exploitation by malicious actors.

The dual-edged nature of AI in cybersecurity is a recurring theme in the findings. While AI offers substantial benefits, such as enhanced data processing capabilities and the ability to adapt to evolving threats, it also poses significant risks, particularly if not implemented with careful consideration of security implications. The visualizations provided in this research offer a clear representation of these dynamics, illustrating both the potential benefits and the associated risks of AI integration in cybersecurity.

This research contributes to a deeper understanding of AI's role in cybersecurity, highlighting the need for ongoing research and development to balance the advantages of AI with the imperative to protect privacy and security [31,32]. The findings emphasize the importance of developing more secure and resilient digital systems that can harness the power of AI while mitigating its risks (**Figures 2–7**). By addressing these critical issues, the research aims to support the development of AI-driven cybersecurity frameworks that are both effective and secure, ensuring that the benefits of AI can be fully realized in the digital age.

5. Discussions and future directions

Artificial intelligence (AI) has become an indispensable tool in the cybersecurity

landscape, offering significant benefits such as enhanced threat detection, process automation, and improved efficiency in managing security operations. However, the adoption of AI in cybersecurity is not without challenges. While AI can significantly bolster defenses, it also introduces new risks that require careful consideration and management. As organizations increasingly rely on AI to safeguard their digital assets, they must remain vigilant about potential vulnerabilities that could be exploited by malicious actors. Ensuring data integrity, preventing data manipulation, and securing high-quality data are critical to maintaining the accuracy and effectiveness of AI systems.

To effectively harness the power of AI while mitigating associated risks, organizations need to implement best practices tailored to their specific security challenges. Developing a well-defined AI strategy that aligns with organizational goals and integrates seamlessly into existing security frameworks is essential. This strategy should prioritize data quality and privacy, ensuring that AI systems are fed with accurate, reliable data while adhering to stringent privacy protections. Furthermore, the ethical implications of AI must be addressed, particularly in terms of bias and fairness in decision-making processes that impact individuals. Building an ethical framework is crucial to mitigate these concerns and ensure that AI-driven cybersecurity solutions are both effective and equitable.

The dynamic nature of cybersecurity threats necessitates continuous adaptation and innovation. Regular testing and updating of AI models are vital to keeping pace with the evolving threat landscape and maintaining optimal system performance. As AI technologies advance, their role in cybersecurity is expected to expand, with emerging technologies such as 5G and the Internet of Things (IoT) offering new opportunities for enhanced security capabilities. The integration of AI with these technologies promises to revolutionize the cybersecurity field, providing more sophisticated tools for threat detection, risk management, and overall digital security.

Looking ahead, the impact of AI on the cybersecurity industry and the job market is likely to be profound. While AI can automate repetitive tasks, reducing the burden on human operators, it also opens up new opportunities for human-machine collaboration [31–35]. Cybersecurity professionals will increasingly partner with AI systems to enhance security at scale, allowing them to focus on more strategic and complex tasks that require human expertise. This shift underscores the need for ongoing education and training to equip cybersecurity professionals with the skills necessary to work alongside AI technologies effectively.

As AI continues to play a more prominent role in cybersecurity, it is crucial to understand and mitigate the associated risks. The cat-and-mouse dynamic between hackers and cybersecurity experts will only intensify as AI-driven tools become more prevalent. Advanced AI technologies must be leveraged to stay ahead of malicious actors, but this requires a proactive approach that anticipates potential threats and responds swiftly to emerging vulnerabilities. In today's rapidly evolving digital landscape, staying updated with the latest technological advancements is not just beneficial but necessary. Our daily lives are increasingly influenced by machine-driven processes and continuous data flows, which, while offering convenience and efficiency, also pose significant security risks if not properly managed. AI has the

potential to transform the cybersecurity landscape into a safer, more secure environment. However, without adequate oversight and control, the same technologies that protect us could be used to create new threats. Ensuring the responsible use of AI in cybersecurity is a collective responsibility that involves researchers, developers, policymakers, and end-users. Establishing proper guidelines and leveraging expert insights are essential steps in maintaining a balance between innovation and security. As AI continues to evolve, stakeholders must remain actively engaged in ensuring that these technologies are developed and deployed responsibly. The future of cybersecurity depends on our ability to harness the positive aspects of AI while preventing its potential misuse, ensuring a secure and balanced digital future for all. While AI holds immense promise for enhancing cybersecurity, it is imperative to maintain vigilant oversight and control. By doing so, organizations can fully realize the benefits of AI while minimizing its risks, ultimately contributing to a more secure digital world.

6. Conclusions

As we stand on the cusp of a new era defined by accelerated computing and technological innovation, artificial intelligence (AI) is set to fundamentally transform various aspects of our world, including the cybersecurity landscape. The integration of AI into cybersecurity systems presents both unprecedented opportunities and significant challenges. AI has the potential to revolutionize threat detection, automate responses to security incidents, and enhance predictive analytics, making our digital environments more secure and resilient. However, these benefits are accompanied by substantial risks, particularly regarding privacy and security vulnerabilities that may arise from AI misuse. The dual-edged nature of AI in cybersecurity underscores the importance of a balanced approach.

While AI can greatly improve our ability to protect sensitive information and detect cyber threats, it also introduces new potential threats if not properly managed. The possibility of AI being exploited by malicious actors cannot be overlooked, and it is crucial that society remain vigilant and adaptive in response to these evolving challenges. To ensure the responsible deployment of AI in cybersecurity, it is imperative to establish robust guidelines and regulatory frameworks that address both ethical and security concerns.

These frameworks must be designed to mitigate the risks associated with AI abuse, safeguarding the privacy and security of individuals and organizations. Moreover, continuous research and development are necessary to stay ahead of emerging threats and to refine AI-driven security measures.

Collaboration among stakeholders, including researchers, developers, policymakers, and industry leaders, is essential in creating a comprehensive approach to AI in cybersecurity. This collaborative effort should focus on developing best practices that promote the ethical use of AI, ensuring that its benefits are maximized while its risks are minimized. The establishment of a culture of responsible AI usage is vital in achieving this balance, fostering an environment where innovation can thrive without compromising security.

As AI continues to evolve and integrate more deeply into the cybersecurity

landscape, it is clear that its impact will be profound. However, with careful management and strategic oversight, we can harness the power of AI to create a safer, more secure digital future. By proactively addressing the ethical and security implications of AI, we can ensure that its advancements contribute positively to society, benefiting all of humanity. The future of AI in cybersecurity holds immense promise, but it also demands careful consideration and responsible action. By embracing innovation while rigorously managing the associated risks, we can ensure that AI serves as a powerful tool for enhancing cybersecurity, ultimately leading to a more secure and resilient digital environment for all.

Author contributions: Conceptualization, ZBA; methodology, ZBA; software, ZBA; validation, ZBA; formal analysis, ATR; investigation, ZBA; resources, ATR; data curation, ZBA; writing—original draft preparation, ZBA; writing—review and editing, ZBA; visualization, ZBA; supervision, ZBA. All authors have read and agreed to the published version of the manuscript.

Acknowledgments: The authors would like to acknowledge and enthusiastically thank the GOOGLE Deep Mind Research with its associated pre-prints access platforms. This research was deployed and utilized under the various platforms and provided by GOOGLE Deep Mind which is under the support of the GOOGLE Research and the GOOGLE Research Publications under GOOGLE Gemini platform. Using their provided platform of datasets and database files with digital software layouts consisting of free web access to a large collection of recorded models that are found in research access and its related open-source software distributions which is the implementation and simulation of analytics for the proposed research which was undergone and set in motion. There are many datasets, data models which are resourced and retrieved from a wide variety of GOOGLE service domains. All the DATA SOURCES and various domains from which data has been included and retrieved for this research are identified, mentioned and referenced where appropriate. However, various original data sources, some of which are not all publicly available, because they contain various types of private information. The available platform provided data sources that support the findings and information of the research investigations are referenced where appropriate.

Conflict of interest: The authors declare no conflict of interest.

References

1. Schatz D, Bashroush R, Wall J. Towards a More Representative Definition of Cyber Security. *The Journal of Digital Forensics, Security and Law*. 2017; 12(2): 1558-7215. doi: 10.15394/jdfsl.2017.1476
2. Stevens T. Global Cybersecurity: New Directions in Theory and Methods. *Politics and Governance*. 2018; 6(2): 1-4. doi: 10.17645/pag.v6i2.1569
3. Misa TJ. Computer Security Discourse at RAND, SDC, and NSA (1958-1970). *IEEE Annals of the History of Computing*. 2016; 38(4): 12-25. doi: 10.1109/mahc.2016.48
4. Stoneburner G, Hayden C, Feringa A. *Engineering Principles for Information Technology Security (A Baseline for Achieving Security)*, Revision A. National Institute of Standards and Technology; 2004. doi: 10.6028/nist.sp.800-27ra
5. Yost JR. The Origin and Early History of the Computer Security Software Products Industry. *IEEE Annals of the History of Computing*. 2015; 37(2): 46-58. doi: 10.1109/mahc.2015.21

6. Nicole P. How the U.S. Lost to Hackers. The New York Times; 2021.
7. Computer Security and Mobile Security Challenges. Available online: https://www.researchgate.net/publication/298807979_Computer_Security_and_Mobile_Security_Challenges (accessed on 20 June 2024).
8. Multi-Vector Protection Securing users and devices across all stages of a malware attack. Available online: https://www-cdn.webroot.com/4415/0473/1276/WSA_Multi-Vector_Protection_WP_us.pdf (accessed on 20 June 2024).
9. What is a Phishing Attack? Defining and Identifying Different Types of Phishing Attacks. Available online: <https://www.digitalguardian.com/blog/what-phishing-attack-defining-and-identifying-different-types-phishing-attacks> (accessed on 20 June 2024).
10. Bendovschi A. Cyber-Attacks—Trends, Patterns and Security Countermeasures. *Procedia Economics and Finance*. 2015; 28: 24-31. doi:10.1016/S2212-5671(15)01077-1
11. Lebo H. The UCLA Internet Report: Surveying the Digital Future. UCLA Center for Communication Policy; 2000. pp. 1-55.
12. Buchanan BG. A (Very) Brief History of Artificial Intelligence. *AI Magazine*. 2005; 26(4): 53-60. doi: 10.1609/aimag.v26i4.1848
13. Kurzweil R. AI Set to Exceed Human Brain Power. CNN; 2006.
14. Feigenbaum EA, McCorduck P. The Fifth Generation: Artificial Intelligence and Japan's Computer Challenge to the World. Addison-Wesley; 1983.
15. Haugeland J. Artificial Intelligence: The Very Idea. MIT Press; 1989.
16. NRC. Developments in Artificial Intelligence. National Academy Press; 1999.
17. Newell A, Simon HA. GPS: A Program that Simulates Human Thought. In: Feigenbaum EA, Feldman J (editors). *Computers and Thought*. McGraw-Hill; 1963.
18. Newquist HP. The Brain Makers: Genius, Ego, And Greed in the Quest for Machines That Think. Mac-millan/SAMS; 1994.
19. Tversky A, Kahneman D. Judgment under uncertainty: Heuristics and biases. In: *Science*. Cambridge University Press; 1982. pp. 1124-1131.
20. Kaplan A, Haenlein M. Siri, Siri, in my hand: Who's the fairest in the land? On the interpretations, illustrations, and implications of artificial intelligence. *Business Horizons*. 2018; 62(1): 15-25. doi: 10.1016/j.bushor.2018.08.004
21. Poole D, Mackworth A, Goebel R. *Computational Intelligence: A Logical Approach*. Oxford University Press; 1998.
22. Samuel AL. Some studies in machine learning using the game of checkers. *IBM Journal of Research and Development*. 1959; 3(3): 210-219, doi:10.1147/rd.33.0210
23. Luger G, Stubblefield W. *Artificial Intelligence: Structures and Strategies for Complex Problem Solving*, 5th ed. Benjamin/Cummings; 2004.
24. Turing AM. Computing Machinery and Intelligence. *Mind*. 1950; LIX(236): 433-460. doi:10.1093/mind/LIX.236.433
25. Pearl J. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann; 1988.
26. British Standard Institute. Part 1: Concepts and models for information and communications technology security management. In: *Information Technology— Security Techniques—Management of the Information and Communications Technology Security*. British Standard Institute; 2004.
27. Kiountouzis EA, Kokolakis SA. *Information Systems Security: Facing the Information Society of the 21st Century*. Chapman & Hall, Ltd; 1996.
28. Vijayan J. *New Vulnerability Database Catalogs Cloud Security Issues*. Dark Reading; 2022.
29. David H. *Operating System Vulnerabilities. Exploits and Insecurity*; 2015.
30. Most laptops vulnerable to attack via peripheral devices. Available online: <http://www.sciencedaily.com/releases/2019/02/190225192119.htm> (accessed on 20 June 2024).
31. Akhtar Z, Rawol A. Uncovering Cybersecurity Vulnerabilities: A Kali Linux Investigative Exploration Perspective. *International Journal of Advanced Network, Monitoring and Controls*. 2024; 9(2): 11-22. doi:10.2478/ijanmc-2024-0012
32. Akhtar Z. Securing Operating Systems (OS): A Comprehensive Approach to Security with Best Practices and Techniques. *International Journal of Advanced Network, Monitoring and Controls*. 2024; 9(1): 100-111. doi: 10.2478/ijanmc-2024-0010
33. Akhtar ZB. Unveiling the evolution of generative AI (GAI): a comprehensive and investigative analysis toward LLM models (2021-2024) and beyond. *Journal of Electrical Systems and Information Technology*. 2024; 11(1): 22. doi: 10.1186/s43067-024-00145-1
34. Akhtar ZB. The design approach of an artificial intelligent (AI) medical system based on electronical health records (EHR)

- and priority segmentations. *The Journal of Engineering*. 2024; 2024(4): e12381. doi: 10.1049/tje2.12381
35. Bin AZ. Artificial intelligence (AI) within manufacturing: An investigative exploration for opportunities, challenges, future directions. *Metaverse*. 2024; 5(2): 2731. doi: 10.54517/m.v5i2.2731

Review

Applications of reinforcement learning, machine learning, and virtual screening in SARS-CoV-2-related proteins

Yasunari Matsuzaka^{1,2,*}, Ryu Yashiro^{2,3}

¹ Division of Molecular and Medical Genetics, Center for Gene and Cell Therapy, The Institute of Medical Science, The University of Tokyo, Minato-ku, Tokyo 108-8639, Japan

² Administrative Section of Radiation Protection, National Institute of Neuroscience, National Center of Neurology and Psychiatry, Kodaira, Tokyo 187-8551, Japan

³ Department of Mycobacteriology, Leprosy Research Center, National Institute of Infectious Diseases, Tokyo 162-8640, Japan

* **Corresponding author:** Yasunari Matsuzaka, yasunari80808@ims.u-tokyo.ac.jp

CITATION

Matsuzaka Y, Yashiro R.
Applications of reinforcement learning, machine learning, and virtual screening in SARS-CoV-2-related proteins. *Computing and Artificial Intelligence*. 2024; 2(2): 1279.
<https://doi.org/10.59400/cai.v2i2.1279>

ARTICLE INFO

Received: 9 June 2024
Accepted: 27 August 2024
Available online: 10 September 2024

COPYRIGHT



Copyright © 2024 by author(s).
Computing and Artificial Intelligence is published by Academic Publishing Pte. Ltd. This work is licensed under the Creative Commons Attribution (CC BY) license.
<https://creativecommons.org/licenses/by/4.0/>

Abstract: Similarly, to all coronaviruses, SARS-CoV-2 uses the S glycoprotein to enter host cells, which contains two functional domains: S1 and S2 receptor binding domain (RBD). Angiotensin-converting enzyme 2 (ACE2) is recognizable by the S proteins on the surface of the SARS-CoV-2 virus. The SARS-CoV-2 virus causes SARS, but some mutations in the RBD of the S protein markedly enhance their binding affinity to ACE2. Searching for new compounds in COVID-19 is an important initial step in drug discovery and materials design. Still, the problem is that this search requires trial-and-error experiments, which are costly and time-consuming. In the automatic molecular design method based on deep reinforcement learning, it is possible to design molecules with optimized physical properties by combining a newly devised coarse-grained representation of molecules with deep reinforcement learning. Also, structured-based virtual screening uses protein 3D structure information to evaluate the binding affinity between proteins and compounds based on physicochemical interactions such as van der Waals forces, Coulomb forces, and hydrogen bonds, and select drug candidate compounds. In addition, AlphaFold can predict 3D protein structures, given the amino acid sequence, and the protein building blocks. Ensemble docking, in which multiple protein structures are generated using the molecular dynamics method and docking calculations are performed for each, is often performed independently of docking calculations. In the future, the AlphaFold algorithm can be used to predict various protein structures related to COVID-19.

Keywords: angiotensin-converting enzyme 2; AlphaFold; Deep Q Network; molecular dynamics; SARS-CoV-2; reinforcement learning; virtual screening

1. Introduction

The novel coronavirus, severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) emerged as a human pathogen in Wuhan, China at the end of 2019 and has since spread around the world, resulting in a pandemic [1]. Symptoms appear about four to five days after being infected with the virus but can take as long as two weeks. On the other hand, asymptomatic infections have also been reported [2]. The main symptoms include fever, cough, difficulty breathing, body malaise, chills, muscle pain, headache, sore throat, and loss of smell and taste. Elderly people and people with underlying health conditions such as heart disease or diabetes are more likely to develop severe pneumonia [3]. Respiratory symptoms, high fever, diarrhea, and taste disorders have also been reported in other generations. When infected during childhood, the symptoms are mild or asymptomatic, but viral infection itself occurs,

and transmission to the elderly due to asymptomatic infection has also been reported. The host range of SARS-CoV-2 is wide, and this virus infects not only humans and wild animals, but also livestock, pets, laboratory animals, and many other animals, causing various diseases. Genetic sequence analysis has shown that this virus is like the coronavirus found in bats and pangolins [4], and it has been pointed out that these viruses may have undergone genetic recombination. Understanding the structure and function of this virus is essential to developing vaccines and treatments for coronavirus infectious disease, which emerged in 2019 (COVID-19).

Artificial Intelligence (AI) methods are being increasingly utilized to predict various aspects related to SARS-CoV-2. Here are some key findings from the search results:

- I) Prediction of COVID-19 severity based on blood protein profiling:
 - A study aimed to classify COVID-19 patients into mild, severe, critical, and control groups based on blood protein profiling using deep learning, random forest, and gradient-boosted trees [5].
 - The ensemble classifier GBTs produced the highest accuracy in predicting disease severity (96.98%) [5].
 - This approach identified specific proteins associated with COVID-19 severity, highlighting the potential for early diagnosis and treatment strategies [5].
- II) Prediction of SARS-CoV-2 epitopes:
 - Machine learning technologies have been used to predict target human proteins of the SARS-CoV-2 virus based on protein sequences and amino acid composition [6].
 - Studies have focused on epitope prediction for SARS-CoV-2 S protein using machine learning models and immunological data from SARS-CoV [7].
 - The aim is to identify nonallergenic, highly antigenic, and nontoxic epitopes that can be used in vaccine design against SARS-CoV-2 [7].
- III) AI-based mutation prediction in SARS-CoV-2:
 - Research is ongoing to develop AI models that predict the next variants of the SARS-CoV-2 virus based on genomic data.

These studies demonstrated the potential of AI-based methods in predicting COVID-19 severity, identifying epitopes for vaccine design, and forecasting mutations in the SARS-CoV-2 virus.

In this review, we focused our attention on the relevant new fields, such as the prediction of the SARS-CoV-2-related protein with AI, such as reinforcement learning and AlphaFold.

2. Classification and structure of SARS-CoV-2

Coronaviruses that infect birds and mammals belong to the order Nidovirales, family Coronaviridae, subfamily Orthocoronaviridae, which includes four genera: alphacoronavirus, betacoronavirus, gammacoronavirus, and deltacoronavirus. Currently, seven types of coronaviruses are known to infect humans; HCoV-229E, HCoV-NL63, HCoV-OC43, and HCoV-HKU1, which are human coronaviruses (HCoV) that routinely infect humans, SARS coronavirus (SARS-CoV-1), which

caused Severe Acute Respiratory Syndrome (SARS) in 2003, Middle East Respiratory Syndrome (MERS) coronavirus (MERS-CoV), which emerged in 2012, and the new coronaviruses (SARS-CoV-2) that is currently causing a pandemic [8]. Among the seven viruses mentioned above, HCoV-229E and HCoV-NL63 belong to the alphacoronavirus genus, and the remaining five viruses (HCoV-OC43, HCoV-HKU1, SARS-CoV-1, MERS-CoV, and SARS-CoV-2) is classified into the beta coronavirus genus, which is divided into four lineages (A, B, C, and D lineages) (**Figure 1**) [9]. Phylogenetic analysis indicates that all the coronaviruses that infect humans are derived from wild animals including bats and rodents. It is thought that coronaviruses originally carried by natural hosts including bats and rodents first infected intermediate hosts, and then eventually infected humans, causing disease. Regarding SARS-CoV-2, the sequence of a coronavirus closely related to this virus has been found in bats, so, likely, the natural host of SARS-CoV-2 is also a bat. Additionally, a coronavirus closely related to SARS-CoV-2 has been detected in Malayan pangolins, so there is a theory that Malayan pangolins are an intermediate host, but the details are unknown.

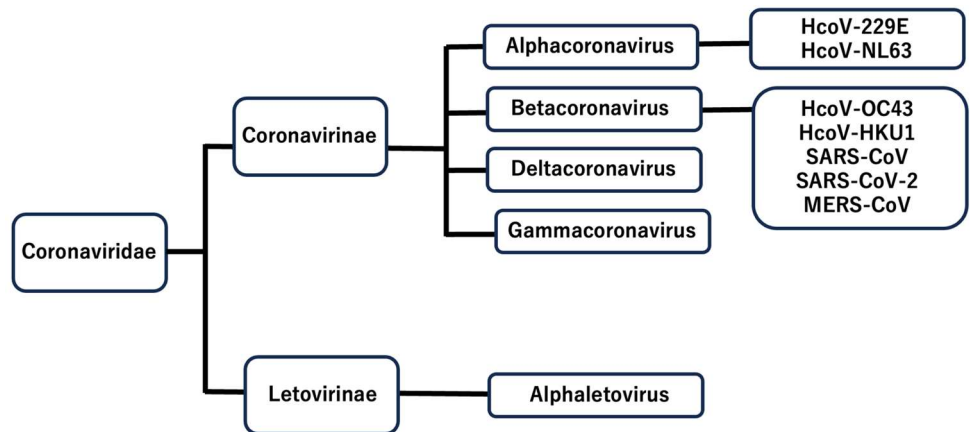


Figure 1. Taxonomy of coronaviridae family [10].

The homology of SARS-CoV-2 genomic RNA and viral proteins with SARS-CoV-1 is 79.0% for genomic RNA, 76.2% for S protein, 94.7% for E protein, 90.1% for M protein, and 90.3% for N protein [11]. Betacoronavirus lineage B, which is included by SARS-CoV-1 and SARS-CoV-2, and an enveloped, single-stranded RNA virus characterized by spikes protruding from its surface and an unusually large RNA genome whose size is approximately 27 to 32 kb that is the largest among currently known RNA viruses [12]. The SARS-CoV-2 genome, whose size is approximately 30 kb, encodes four structural proteins; spike (S) protein, nucleocapsid (N) protein, membrane (M) protein, and envelope protein, each of which is essential for constructing the virus particle (**Figure 2**) [13]. The genomic RNA has a cap structure and a poly (A) sequence, at the 5' end at the 3' end, respectively, so it can infect host cells and function directly as mRNA. There are two open reading frames (ORF1a and ORF1b) in approximately 20 kb at the 5' end of the viral RNA, and the start codon of ORF1b is located slightly upstream of the stop codon of ORF1a. Two proteins are translated from ORF1a and ORF1b:1a and 1a + 1b, which is synthesized by frameshifting of ribosomes. These proteins are cleaved by their proteases into more

than a dozen types of nonstructural proteins, including RNA-dependent RNA polymerase.

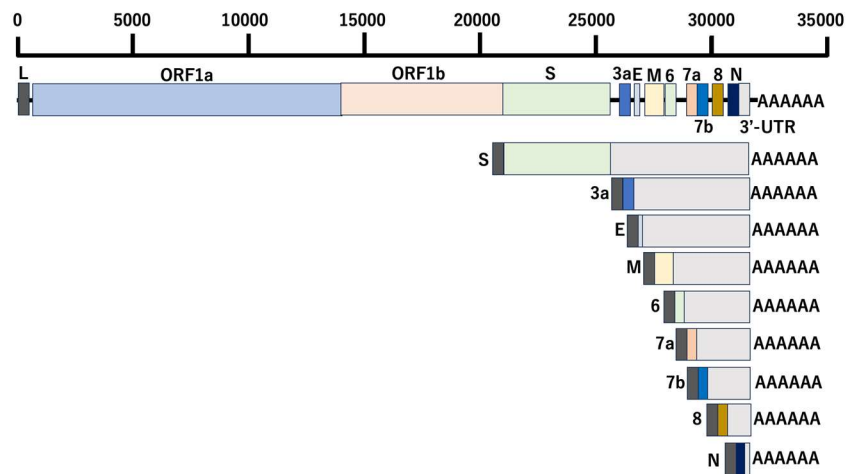


Figure 2. SARS-CoV-2 genome organization and the canonical subgenomic mRNAs. The full-length genomic RNA (29,903 nt) which also serves as an mRNA, ORF1a and ORF1b are translated [14].

3. Interaction of SARS-CoV-2 S protein with a receptor on the cell

Angiotensin-converting enzyme 2 (ACE2) is recognizable by the S protein on the surface of the SARS-CoV-2 or SARS-CoV virus [15]. On the other hand, it has also been shown that the S protein of SARS-CoV-2 does not recognize Dipeptidyl peptidase 4 (DPP4), the receptor for MERS-CoV, and APV, the receptor for HCoV-229E [16]. ACE2 and S protein combine like a lock and key, allowing the virus to enter human cells [17]. The SARS-CoV-2 virus is very similar to the SARS-CoV virus that causes severe acute respiratory syndrome (SARS), but the receptor binding region of the S protein significantly enhances the binding affinity of the SARS-CoV-2 virus to ACE2 via several mutations [18]. Like all coronaviruses, SARS-CoV-2 uses the S glycoprotein to enter host cells, which contains two functional domains: S1 and S2 RBDs [19]. The two subunits, S1 and S2 are cleaved from the S protein by host cell proteases [20]. S1 plays a role in receptor binding via RBD, and S2 plays a role in membrane fusion between the viral envelope and the cell [21].

The SARS-CoV-2 S protein first binds to the host cell's ACE2 receptor, which is a membrane protein with an enzyme domain in the cell membrane of human cells, via the S1 RBD [22]. The receptor specificity of the S protein is a major factor determining the host range and tissue tropism which is the ability to selectively infect specific tissues or organs of coronavirus [23]. It has been identified that ACE2 is the receptor for SARS-CoV-1, DPP4 is the receptor for MERS-CoV, aminopeptidase N (APN) is the receptor for HCoV-229E, and 9-O-acetylated sialic acid is a receptor for HCoV-OC43 and HCoV-HKU1, respectively [24]. The cell entry mechanisms of coronaviruses have two routes after binding to the receptor, 1) entry into the cell from the cell surface, and 2) entry into the cell via endosomes after the virus particle is taken into the cell by endocytosis [25]. When an enveloped virus invades a cell, the viral envelope needs to fuse with the cell's lipid bilayer membrane [26]. In the case of

coronaviruses, the S protein subunit S2 contains a fusion peptide, which plays an important role in membrane fusion. In SARS-CoV-1, the second route is the main pathway, in which the viral S protein taken up by endocytosis is activated by host proteases and causes membrane fusion between the endosomes and the viral envelope [27].

Host proteases that can activate the S protein of SARS-CoV-1 include cathepsin, trypsin, elastase, and TMPRSS2 [28]. Additionally, the S protein of MERS-CoV is cleaved into S1 and S2 by Furin [29]. One of the major differences between SARS-CoV-1 and SARS-CoV-2 is that SARS-CoV-2 has a characteristic sequence of consecutive basic amino acids (RRAR) in the S1/S2 cleavage site of the S protein, called “Furin cleavage site” which is absent in SARS-CoV-1, but is present in the S protein of MERS-CoV and HCoV-OC43, and efficiently cleaved by Furin and other proteases. During the virus replication cycle, the S protein is cleaved into S1 and S2, but the location and timing of cleavage differs depending on the types of coronaviruses, that is 1) S protein is synthesized in infected cells and then cleaved by host protease, and 2) when a virus invades a target cell, the S protein binds to a receptor and is then cleaved by host protease [30].

The mechanism in SARS-CoV-1 is the latter, so the S protein exists in an uncleaved state on the surface of the virus particle, and when the virus invades cells, it is cleaved by host proteases (trypsin, elastase, cathepsin, TMPRSS2) [31]. In contrast, in the case of SARS-CoV-2, cleavage occurs within the cell after S protein synthesis due to the first mechanism [32]. Experiments using pseudotyped viruses have shown that the S protein of virus particles exists as cleaved forms of S1 and S2 [33]. In addition, it has been suggested that the S protein cleavage site with the Furin needs SARS-CoV-2 to efficiently infect the human respiratory tract and that the S protein activation by TMPRSS2 is important [34]. The RBD in the SARS-CoV-1 S protein is composed of a core structure and a receptor binding motif (RBM), and the RBM directly binds to the ACE2 surface [35]. The six amino acids Y442, L472, N479, D480, T487, and Y491 in the RBM of SARS-CoV-1 are critical for binding to ACE2 and are involved in determining the host range of SARS-related coronaviruses [36]. In SARS-CoV-2, those corresponding to these six amino acids are L455, F486, Q493, S494, N501, and Y505, but except for Y505 (Y491 in SARS-CoV-1), different from amino acids [37]. Regarding the binding affinity between the RBM of SARS-CoV-2 and ACE2 in various animal species, such as humans using computer analysis of the protein structure, the RBM of SARS-CoV-2 has a high binding affinity for ACE2 in humans, civets, pigs, ferrets, cats, orangutans, monkeys (green monkeys), and bats (acetone), and the high binding affinity of mouse and rat ACE2 was predicted to be low [38].

After membrane fusion, the virus unsheds and the virus genome is released into the cell, whereupon virus replication begins within the cytoplasm [39]. Since coronaviruses positive-strand genomic RNA can function as mRNA, it binds to host cell ribosomes and synthesizes RNA-dependent RNA polymerase and other substances necessary for virus replication [40]. Using the positive-strand genomic RNA as a template, the mRNA encoding each viral protein is transcribed based on the synthesized complementary negative-strand RNA, and the viral protein is produced

[41]. Replication of positive-strand genomic RNA for progeny viruses also takes place. Newly synthesized viral structural proteins (S, E, and M proteins) are transported to the endoplasmic reticulum-Golgi apparatus intermediate (ERGIC) [42]. The nucleocapsid formed by the N protein and viral RNA, together with other structural proteins, forms the virus particle and buds into the ERGIC [43]. M and E proteins play important roles in the virus budding step [44]. Progeny viruses budded within ERGIC are released outside the cell by exocytosis [45].

ACE is an enzyme that catalyzes the conversion of the peptide hormone angiotensin I (Ang I) to angiotensin II (Ang II) and is well-known as a vasoconstrictor that promotes muscle contraction of blood vessel walls and narrows the lumen of blood vessels [46]. ACE2, a viral receptor, also plays a role as a vasodilator. This is because it balances ACE and relaxes blood vessel walls [47]. Both ACE and ACE2 play pivotal roles in the renin-angiotensin system (RAS), which regulates blood pressure and blood flow in multiple organs, including the lung, heart, and kidney, and conjugates a complex network of enzymes, peptide hormones, and receptors [48]. Angiotensinogen, a precursor of Ang secreted by the liver, is cleaved by the kidney enzyme renin to produce Ang I, which is converted to Ang II, an eight amino acid hormone peptide by ACE [49]. Ang II binds to the type 1 angiotensin receptor (AT1R) on the surface of microvascular muscle cells, causing vasoconstriction and promoting salt reabsorption in the kidney [50]. Vasoconstriction and salt reabsorption both contribute to increased blood pressure [51]. Therefore, when ACE activity becomes abnormally high, the amount of Ang II increases, causing hypertension.

On the other hand, ACE2 catalyzes the eight amino acid peptide of Ang II to a seven amino acid peptide (Ang 1–7) [52]. Though its action on a different receptor the Mas-1 receptor (MasR), it has the opposite effect on Ang II [53]. Although the detailed role of Ang 1–7 in blood pressure regulation is not completely understood, it is believed that Ang 1–7 decreases blood pressure and induces vasodilation [18]. Further, ACE2 splits Ang I into Ang 1–9, thereby balancing the effects of ACE by removing the substrate [54]. By converting Ang II to Ang (1–7) and Ang I to Ang 1–9, ACE2 plays an important role in maintaining the balance between vasoconstriction and vasodilation to sustain blood pressure within an appropriate range [55].

4. Reinforcement learning

4.1. Q-learning in a finite Markov decision process

Reinforcement learning is a type of machine learning that is a “mechanism for AI to automatically learn” and is a technology for machines to automatically identify and predict based on learned data [56]. It refers to a technology in which the system learns appropriate control methods through repeated trials and error. Its main feature is that it can analyze data without human intervention. In conventional machine learning, humans had to extract and adjust “feature values”, which are indices for learning the data to be analyzed [57]. However, deep learning does not require human intervention to extract feature values, so machine learning can be easily performed [58]. Machine learning is mainly composed of the following three types: supervised learning with correct data, unsupervised learning with no correct data, and

reinforcement learning. The machine learns by recognizing many images as correct data. This method is called “supervised learning” [59].

On the other hand, “unsupervised learning “is a method of learning without giving correct data [60]. Machines analyze the characteristics of data, making it possible to identify and classify data. In reinforcement learning, an agent placed in a certain environment act on the environment and seeks a policy that maximizes the reward obtained [61]. The learning progresses through a cycle in which the agent acts on the environment, the environment updates the state evaluates the action, and informs the agent of the state and reward. The action-value function and policy are optimized through learning so that the total reward obtained is maximized. The reinforcement learning repeats the following steps: 1) the agent acts on the environment, 2) the environment gives agents updated states and rewards, and 3) the agent modifies its behavioral strategy based on the reward and returns to 1).

Q-learning, a kind of reinforcement learning, is one of the policy-off Temporal Difference (TD) learning of machine learning methods [62]. *Q*-learning converges to the optimal evaluation value when it tries an infinite number of episodes in which all states are sufficiently sampled in a finite Markov decision process [63]. In *Q*-learning, each rule to be executed has a value called *Q*-value, which indicates the effectiveness of the rule, which is a pair of a state and an agent’s possible actions under that state, and the value is updated each time the agent acts. For example, assume that the agent’s current state is *St*, and there are four possible actions *a*, *b*, *c*, and *d* in this state. At this time, the agent decides the action to take based on the four *Q* values, *Q* (*St*, *a*), *Q* (*St*, *b*), *Q* (*St*, *c*), and *Q* (*St*, *d*).

Theoretically, the *Q* value convergence has been proven even if the trial is performed an infinite number of times. Still, to expedite the convergence, actions with a large *Q* value are selected with a high probability. As a selection method, select randomly with a small probability *e*, otherwise select the action with the maximum *Q* value, *e*-Greedy method, and roulette selection used in genetic algorithm, Boltzmann distribution as below softmax Equation (1) is used.

$$\pi (St, a) = \exp (Q (st, a)/T) / \sum_{p \in A} \exp (Q (St, p)/T) \quad (1)$$

where *T* is a positive constant and *A* is the set of possible actions of the agent in state *St*. If the action is decided, then update the state and the *Q* value of the action. As an example, the state *St* agent chooses action *a* and the state transitions to *st + 1*. The updated formula for the action-value function, *Q* function in *Q*-learning (2) is as follows.

$$Q (st, at) = Q (st, at) + \alpha [rt + 1 + \gamma \max_{a \in A} (st + 1) Q (st + 1, at + 1) - Q (st, at)] \quad (2)$$

here, alpha is called a learning rate, which is a numerical value that satisfies the conditions described later, and gamma is called a discount rate, which is a constant between 0 and 1 inclusive. Also, *rt + 1* is the reward the agent got when it transitioned to *St + 1*. The above update formula means that when the current state moves to the next state, the *Q* value is brought closer to the value of the state with the highest *Q* value in the next state. This means that if a state has a high reward, that reward will propagate to states that can reach that state with each update. As a result, optimal state transition learning is performed. When the learning rate satisfies the following conditions, in *Q*-learning, all *Q* values converge to the optimal value with probability.

$$\sum_{t=0}^{\infty} \alpha(t) \rightarrow \infty \quad (3)$$

$$\sum_{t=0}^{\infty} \alpha(t)^2 < \infty \quad (4)$$

Due to this good convergence, many studies have been done on Q -learning, but some problems have been pointed out.

In Q -learning, the Q -function was updated by updating the number of states $s \times$ the number of actions a . However, as the number of states grows, it becomes impractical to represent the Q function with a table. To solve this problem, the Deep Q Network (DQN) takes the approach of expressing the Q function with a convolutional neural network and devises ways to converge learning [64]. However, simply replacing the Q function with a convolutional neural network (CNN) does not result in successful learning, so efforts have been made to converge the learning. Deep reinforcement learning is a combination of reinforcement learning and deep learning methods, and the representative method is this DQN, which is an approximation that replaces the action value function, Q function in Q learning with a CNN.

In the automatic molecular design method based on deep reinforcement learning, it is possible to design molecules with optimized physical properties by combining a newly devised coarse-grained representation of molecules with deep reinforcement learning.

Reinforcement learning and virtual screening in drug discovery and materials design.

Reinforcement learning is a type of machine learning that uses two factors: agent and environment. The agent acts and learns the feedback (reward) from the environment regarding that action, thereby deriving a behavioral guideline (strategy) to maximize the reward. The main feature is that it is less dependent on datasets. In reinforcement learning, based on feedback from the environment, it does not require a static dataset unlike unsupervised and supervised learning because the agent learns from the experience it collects. In other words, there is no need for data collection, preprocessing, or labeling before learning. The reinforcement learning workflow generally is as follows; (i) creating the environment: define the environment in which the agent operates, including the interface between the agent and the environment—introducing simulation from the standpoint of safety and experiment ability. (ii) Definition of remuneration: define rewards for goals and decide how to calculate rewards. Rewards guide the agent's behavioral choices. (iii) Creating an agent: define an agent consisting of a policy and a reinforcement learning algorithm. Specifically, the selection of the method of representing the policy: neural networks, lookup tables, etc. Choosing an appropriate learning algorithm: neural networks are commonly used because they are well suited for learning in large state and action spaces. (iv) Agent learning and verification: set conditions for learning, such as stopping conditions, and perform agent learning. After learning, verify the policy derived by the agent. Reconsider the design of reward signals and strategies, and rerun learning if necessary. Reinforcement learning is sample-inefficient, especially for model-free, on-policy algorithms, and can take minutes to days to train. Therefore, learning is often

parallelized on multiple central processing *units* (CPUs), Graphics Processing Units (GPUs), or clusters. (v) Development of measures: Investigate the learned strategies. Based on this result, the process may return to the initial stage of the workflow. Specifically, if the learning process and policy derivation do not converge within the calculation time, the following items need to be updated before relearning.

Searching for new compounds is an important initial step in drug discovery and materials design, but the problem is that this search requires trial-and-error experiments, which are costly and time-consuming. On the other hand, in recent years, *in silico* drug discovery and materials search, in which chemical compounds are searched for in a computer, have been attracting attention. However, it is generally difficult to search the space of discrete chemical structures of compounds, and an efficient method is required. Therefore, in recent years, new search methods using deep neural networks (DNNs), such as methods using generative models such as variational autoencoders (VAE), have been proposed [65]. These techniques attempt to circumvent this problem by learning the mapping between the discrete compound space and the continuous latent space by a generative model approximated by DNNs, and by allowing compound optimization to be performed in this continuous latent space.

However, this method had some problems. For example, there is no objective metric to evaluate whether the learned mapping is suitable for efficient optimization. In addition, the learning process of the generative model is separated from the optimization process of the compound concerning the score function of the optimization target molecule. On the other hand, in methods based on reinforcement learning; by thinking of molecular design as a Markov decision process, the agent learns the optimal policy through the rewards provided by the surrounding environment.

Virtual screening (VS) in drug discovery is a method of selecting drug candidate compounds from many compounds using computers. Naturally, it cannot be used as a drug unless it shows medicinal efficacy, so VS mainly focuses on medical efficacy and evaluates the presence or absence of activity against drug target proteins. Such VS can be broadly divided into methods based on known active compounds (ligand-based VS; LBVS) and methods based on protein three-dimensional (3D) structure (structure-based VS; SBVS). LBVS is a method that mainly uses similarity evaluation of compounds and machine learning and uses known experimental results to construct regression prediction models and classification prediction models and uses these to select compounds [66]. While drug-candidate compounds can be selected with relatively high precision. Because it learns based on the few compounds that have been tested against the target protein, it is difficult to develop guidelines for how to optimize the chemical structure of selected drug candidate compounds due to the lack of novelty in the chemical structure of predicted active compounds and the lack of 3D structural knowledge, having been pointed out as major problems.

On the other hand, SBVS uses protein 3D structure information to evaluate the binding affinity between proteins and compounds based on physicochemical interactions such as van der Waals forces, Coulomb forces, and hydrogen bonds, and select drug candidate compounds [67]. Although this method is less accurate than

LBVS because it does not use known experimental information about the target protein, it can discover highly novel drug candidates. Furthermore, the estimated binding mode between the protein and the compound can be obtained, which can provide guidelines for subsequent structural optimization of the compound. Due to the above two advantages, it is attracting a lot of attention just like LBVS. SBVS and LBVS are used together, and compounds commonly selected by both are sent to in vitro experiments. By introducing LBVS methods such as machine learning using information on known active compounds, drug candidate compounds are often narrowed down to a certain extent before the SBVS process is performed. On the other hand, many methods combine SBVS and machine learning, such as machine learning methods for protein 3D structures and for predicted binding structures obtained by docking calculations [68]. These methods hardly expose the weakness of LBVS, which is that the predicted active compound has little novelty in its chemical structure, due to the use of the structures of various proteins or the binding structures of various proteins and compounds rather than simply using knowledge of known compounds for specific proteins.

In the investigation and selection of target proteins for drug discovery, first, the target protein is selected from among the proteins involved in SARS-CoV-2, which is the target of drug discovery [69]. In addition to selecting target proteins simply based on known infection mechanisms, protein selection is performed using bioinformatics methods, such as selecting target proteins using omics analysis. However, inhibiting proteins that play an important role in the human body can lead to side effects, so it must be avoided as much as possible. If there is an essential protein that has the same function as the target protein, even if it is not a perfect match, it is necessary to show selectivity despite slight structural differences, which increases the difficulty of drug discovery. It is also important to be able to conduct experiments using gene knockout rats and mice during non-clinical trials. In addition to these conditions, to perform SBVS, it is also necessary that the 3D structure of the protein is known, or that a reliable 3D structure can be estimated by homology modeling, etc. Furthermore, the final target protein should be determined by considering the difficulty of the drug binding site.

In SBVS, even changes inside chains can greatly affect the results of docking calculations, so it is necessary to carefully prepare the protein 3D structure [70]. Various protein 3D structure is registered in the Protein Data Bank (PDB), but if the complex structure with a compound is known, the local structure is likely to allow the compound to bind easily, and highly accurate compound selection can be expected [71]. However, the required resolution is strict, and docking calculations require a resolution of at least 2.2Å to 2.5Å. On the other hand, if the protein 3D structure is unknown, it is necessary to predict the 3D structure. In SBVS, the *ab initio* method is rarely used due to the resolution mentioned above, and homology modeling is used to predict the 3D structure using homologous proteins whose structures are known. Examples of homology modeling tools include MODELLER and SWISS-MODEL, etc., [72].

However, the situation in protein 3D structure prediction changed significantly with the release of AlphaFold2 [73]. Furthermore, the ionization state of some protein

residues changes depending on the environment. Since interactions due to Coulomb forces are stronger than der Waals forces and hydrogen bonds, it is extremely important to consider the ionization state. However, changes in the ionization state cannot occur during docking calculations or molecular dynamics (MD) simulations. Therefore, it is necessary to generate an ionized state in advance, and PROPKA is most widely used for this purpose. In most cases, the human body has a nearly neutral environment, so an ionized state of pH 7.0 is often generated and used for docking calculations. However, it should be noted that proteins present in the stomach, for example, must produce an ionization state under acidic conditions.

In drug binding site prediction and selection, identifying protein surface sites (druggable sites) where drugs can be expected to bind is essential for estimating more detailed binding structures [74]. The conditions for a druggable site include having a concave region called a “pocket” when a compound binds, the concave region being of appropriate size and deep enough and having a hydrophobic surface [75]. Among these features, widely used methods include POVME, which predicts drug binding sites based on the protein surface shape, Fpocket and SiteMap, which make estimations by considering the properties of the protein surface, and FTMap, which locates small probe molecules and finds energetically stable spaces [76]. When a clear active site exists, such as in an enzyme, drug design is often aimed at that active site, binding site estimation methods are especially important if a clear concave region appears only after a compound bind [77].

Examples of such cases include when a protein binds to a compound while changing its structure (induced fit), and when designed inhibitors of protein-protein interactions [78]. Another aspect of considering whether a site is a druggable site is the degree of conservation of the amino acid residues that make up the binding site. Significant differences in target protein sequences between experimental animals such as rats and mice and humans can lead to differences in drug efficacy, leading to the suspension of drug development during clinical trials. In addition, with antiviral drugs, it is possible to suppress the acquisition of drug resistance by designing drugs that target highly conserved sites that are essential for protein function [79]. For drug binding site prediction, binding site prediction using 3D convolutional neural networks (3D-CNN) has been actively proposed, such as DeepSite, Kalasanty, and DeepSurf, and a method for predicting peptide binding sites rather than compound binding sites [80].

Evaluation of compounds based on protein 3D structure usually involves computational difficulties. Furthermore, even if it is possible to estimate a drug candidate compound that promotes or inhibits protein function, many compounds are unsuitable as drugs due to problems such as compound solubility and side effects. Based on the above, compound filtering is performed from various perspectives. The most widely used rule for designed oral drugs is Lipinski’s rule of five. This is a rule that Lipinski et al. summarize the chemical properties of drugs approved for oral use. It lists four conditions: molecular weight of 500 or less, hydrogen bond accepting groups of ten or less, hydrogen bond donating groups of five or less, and water-octanol partition coefficient logP of five or less (It’s called the rule of five because everything is a multiple of five). QED (quantitative estimate of drug-likeness) is also widely used

as a method to evaluate this “oral drug-likeness” using real numbers [81]. Additionally, indicators related to side effects and toxicity have been proposed, such as PAINS, which summarizes the characteristics of compound substructures that frequently cause off-target effects that bind to and inhibit or activate other proteins. In addition, LBVS-like methods are often used to select compounds to reduce the amount of docking calculations. However, this should not be done too much, as the result approaches “the discovery rate of binding compounds is high, but the novelty of the compounds is low.”

Like proteins, the ionization state of compounds also changes depending on the environment. Compounds often have a range of ionization states, and an ionization state of approximately pH 7.0 ± 2.0 is generated and used for docking calculations [82]. Tools that generate the ionization state of compounds include Schrodinger’s Epik, ChemAxon’s JChem Protonation Plugin, and the open-source software Dimorphite-DL [83]. Additionally, some compounds may have tautomers or optical isomers may not be separated and may be grouped in one compound entry. Such isomers often have significant effects, such as changes in the interaction mode with proteins and the occurrence of collisions with proteins due to changes in the 3D structure of the compound [84]. Therefore, it is necessary to generate each isomer for these as well. Regarding this, there are JChem Protonation Plugin from ChemAxon, LigPrep from Schrodinger, and open-source software Gypsum-DL.

4.2. Protein-compound docking calculation

Docking calculation is a method for predicting the binding affinity and binding mode of a certain compound to a drug-binding site of a protein. DUD-E is a benchmark data set for SBVS, and the enrichment factor (EF) is a ratio that indicates how much the proportion of active compounds has increased after selection compared to before selection [85]. For example, the EF (EF_{x%}) when selecting the top x% is calculated as follows.

$$EF_x\% = (\text{Pos}_x\% / \text{All}_x\%) / (\text{Pos}_{100}\% / \text{All}_{100}\%)$$

The denominator is the proportion of active compounds included in the benchmark data set, and the numerator is the proportion of active compounds after selection, which based on docking calculations often narrows down the evaluation target to 1/100 or less, so it is often set to a small value such as EF1%. Commercial software such as Glid and Surflex have high prediction accuracy, whereas open-source software AutoDock and AutoDock Vina tend to have lower prediction accuracy. Also, docking calculation takes about ten seconds per compound in Glide SP mode when using one CPU core. In addition, GPU-based docking software such as Quantum. Ligand. Dock and BUDE have been developed, and AutoDock has been implemented with GPU, achieving 250 times faster speed than one CPU core when using NVIDIA Titan V [86].

While docking calculations consider structural changes in compounds, structural changes in proteins are generally not considered. The structure of a protein changes to a greater or lesser degree due to the binding of a compound (induced fit), so taking protein structural changes into account is important for improving prediction accuracy [87]. However, although there are methods that consider structural changes in protein

side chains during docking calculations, they have not become common due to the computational complexity problem [88]. Ensemble docking, in which multiple protein structures are generated using the molecular dynamics (MD) method and docking calculations are performed for each, is often performed independently of docking calculations [89]. However, since the amount of calculation is doubled by the number of protein structures used in the docking calculation, the number of applications is limited to a small number of cases.

In protein-compound docking calculations, a reranking method has been proposed that outputs multiple predicted bond structures in the docking calculation and predicts the interaction mode or interaction energy of the bond structures [90]. Therefore, accuracy is improved compared to ranking based on scores obtained by docking calculations. 3D-CNNs that use the connection structure as an input are being proposed for these as well, but interestingly, there is no significant performance difference between methods that use the interaction mode as a feature and deep learning methods [91]. This suggests a lack of data for deep learning and sufficient maturity of domain knowledge regarding interactions. An example of the application of SBVS to COVID-19 is that SBVS was performed on approximately 2100 approved drugs and active compounds with $IC_{50} < 10$ microns were identified [92]. As a result of binding energy estimation using the MM-PBSA method for each compound, since they showed good binding energies of -8.73 kcal/mol or less, in vitro assays were performed on all of them, and a good hit compound with $IC_{50} < 10$ microns was obtained.

4.3. Compound selection using the MD method

MD methods, which simulate the temporal changes in the coordinates of each atom in environments where solutions such as proteins and solvents such as water exist, are used in a variety of analyses [93]. Programs that perform MD simulations include AMBER, GROMACS, NAMD, CHARMM, and Desmond. In addition, in MD simulation, the speedup rate by using accelerators such as GPU is extremely high. From the perspective of SBVS, MD simulation makes it possible to evaluate the binding strength between a protein and a compound while explicitly considering protein structural changes, solvation, entropic effects, etc., making it possible to select compounds with higher precision [94]. In MD calculations for SBVS, simulations are performed using the predicted binding structure from docking calculations as the starting points. For example, several methods of conducting multiple short-term simulations and evaluating how stable the predicted bond structure is and of highly accurate estimation for binding energy using MM-GBSA, MM-PBSA, or MP-CAFEE, etc., have been proposed. Since the orientation of even a single side chain is important for protein structures in drug discovery, there is a possibility that attention will be focused on estimating protein structures to which compounds can easily bind using MD simulations [95].

5. Prevention of asymptomatic infections of COVID-19

To control the SARS-CoV-2 pandemic, many countries have placed restrictions on non-essential travel, and have subsequently implemented travel restrictions using a

combination of the following four strategies to lift restrictions: whitelist, unrestricted travel permission; gray list, travelers providing proof of a negative PCR and reverse transcription before arrival; red list, travelers quarantined on arrival; blacklist, ban on non-essential travel. Decisions about which list to assign to this vary by country and are often based on publicly available population-level epidemiological indicators: cases per capita, deaths per population, and positivity rate [96]. However, it has been pointed out that these indicators are incomplete, with problems such as underreporting, bias in symptomatic populations, and reporting delays.

To address these issues, it will be possible to derive optimal border policies by using real-time estimates of COVID-19 prevalence and estimating the number of asymptomatic infected people with high accuracy. Unlike normal restriction protocols, allocations can be made from limited information, based on demographic information and past test results of the incoming population. This system estimates the prevalence of COVID-19 based on test results used in the past; i) Adaptively extract a minimum set of traveler types based on demographic characteristics, country, region, age, and gender, using the least absolute shrinkage and selection operator (LASSO) regression from high-dimensional statistic [97]. ii) Estimate the prevalence of each type using the empirical Bayes method, deriving prior probabilities from previous experience. This system environment is such that the prevalence of COVID-19 is low, two in 1,000 people and arrival rates vary widely by country. As a result, testing data is unbalanced (few cases among those eligible for testing) and sparse (few arrivals from specific countries). These data characteristics are sequentially processed using the empirical Bayes method to perform appropriate processing. Utilizing the prevalence estimates described above, a subset of travelers for PCR testing is derived based solely on traveler type. This allocation of tests is done by adjusting the exploration-exploitation trade-off between the two objectives. Specifically, i) maximize the number of infected asymptomatic travelers based on current information (exploitation), and ii) assign tests based on experience to travelers for whom there is no accurate estimate, and accurately understand and update the epidemic status (exploration). For a greedy allocation to this tradeoff, allocating tests to concentrate on types with high prevalence, test data for the types with the highest number of patients and moderate prevalence will not be extracted. As the prevalence of COVID-19 is rapidly increasing in some cases, it is necessary to understand as much as possible of moderate symptoms to carry out appropriate learning. These challenges can be viewed as multi-armed bandit problems in reinforcement learning especially batch bandit problems with non-stationary, contextual, delayed feedback, and constraints. Information from pipeline tests, that are not returning results, must be considered. To solve this exploration-exploitation trade-off, the algorithm is built based on the Gittins index. Each type introduces a deterministic index representing a risk score, incorporating both estimated prevalence and uncertainty, according to which allocations are made.

Reinforcement learning, machine learning, and VS have been utilized in the search for inhibitors against SARS-CoV-2-related proteins. machine learning techniques are commonly employed to identify potential compounds for drug development quickly and accurately. In a study focusing on the SARS-CoV-2 main protease (3CLpro), machine learning-based virtual screening was used to predict new

inhibitors. Algorithms such as K-nearest neighbor (KNN), support vector machine (SVM), and Random Forest (RF) were employed, with RF showing the best performance in classifying phytochemicals as potential inhibitors.

The use of VS combined with molecular docking and molecular dynamics simulations has led to the identification of high-potential therapeutic compounds that could inhibit SARS-CoV-2 pathogenesis. These advanced computational approaches have helped narrow down a list of over 4000 compounds to 26 promising candidates [98].

In another study, deep reinforcement learning was employed after an initial virtual screening to design dual-target inhibitors against SARS-CoV-2 main protease (Mpro) and papain-like protease (PLpro) [99]. Additionally, graph generative models have been explored for designing novel drug candidates targeting SARS-CoV-2 viral proteins. Addressing minor issues, it is important to note that while virtual screening is a powerful tool for drug discovery, it can yield a high proportion of false positive hits. To mitigate this, machine learning-based approaches are increasingly being integrated into virtual screening workflows to enhance accuracy and efficiency.

6. SARS-CoV-2 protein structure prediction by AlphaFold algorithm

With the increasing number of COVID-19 cases, the AlphaFold algorithm, a deep-learning algorithm developed by DeepMind, was utilized to predict various protein structures related to COVID-19 [100]. Given the amino acid sequence, and the building blocks of a protein, AlphaFold can predict 3D protein structures. The analysis of amino acid sequences into 3D structures is typically a long-term and intensive process, involving visualization techniques for a variety of protein and structural analyses, including nuclear magnetic resonance, cryo-electron microscopy, and X-ray crystallography, and is costly. However, AlphaFold, which is an AI system predicting the 3D structure of proteins from amino acid sequence information and won the CASP13 (Critical Assessment of Structure Prediction) competition, an international competition for protein 3D structure prediction, eschews these techniques and uses a DNN that predicts distances and angles between amino acids scored with gradient descent, resulted in achieving a dramatic high score [72]. Proteins have a variety of functions due to the folding of linear chains of amino acids linked by peptide bonds to form 3D structures. By elucidating this structure, it will be possible to elucidate the proteins involved in most diseases involving proteins, especially those related to SARS-CoV-2. However, the method by which proteins fold into their final 3D structure remains a black box. Because the theoretical number is astronomical, it has been pointed out that enumerating all possible configurations of a typical protein by brute force calculations takes a long time and is known as the “protein folding problem.” By using free modeling, AlphaFold can ignore similar structures in predictions, which is particularly useful for COVID-19.

AlphaFold consists of three different layers of DNNs [101]. The first layer consists of a variational autoencoder stacked with an attention model to generate realistic fragments based on a single amino acid sequence. In the second layer, it is divided into two sublayers. The first sublayer uses a 1D Convolutional Neural

Network (CNN) on the contact map to optimize inter-residue distances. This is a 2D amino acid residue distance representation by projection of the contact map into one dimension for input into the CNN. In the second sublayer, it optimizes the scoring network and the degree to the generated substructures observed like proteins using CNN with 3D structure. After normalization, a third neural network layer is added that scores the generated proteins against the actual model. AlphaFold's structure module takes as input the features of the amino acid sequence corresponding to the input sequence and the pair representation features of the MSA (Multiple Sequence Alignment) extracted by the Evoformer part and outputs the coordinates of all atoms and the prediction reliability score pLDDT for each residue [70]. AlphaFold2 consists of four modules. i) Data preparation module: using the amino acid sequence (input sequence) of the predicted 3D structure as a query, create MSA from the database and search template 3D structure (template structure) from the database using bioinformatics tools. However, the use of a template 3D structure is optional. ii) Embedding module: the creation of an MSA representation that links raw MSA with target sequence information and a pair representation that records the relative positional relationship between residues [102]. Dense vector transformation with embedding, which fully connected layer without activation for sparse input values. iii) Evoformer (Transformer for molecular evolution) module: feature extraction from MSA and pair representation [103]. Information exchange between MSA and pair representation. Axial attention and triangular attention are performed keeping in mind the characteristics of MSA and the physical constraints of spatial graphs (proteins). iv) Structure module: integration of MSA representation, residue pair representation, and current 3D structure using IPA (Invariant point attention) module. Prediction of relative movement instructions for each residue (= (3, 3) rotation matrix and (x, y, z) translation vector for the number of residues) and side chain torsion angle. The structure module consists of eight layers with shared weights [104]. Each layer updates the features S of the amino acid sequence and the 3D representation plotted in the coordinate system T_i (corresponding to object coordinates) defined for each residue. T_i is a pair of a rotation matrix R_i , which represents a rotation that superimposes the coordinate system defined for each residue on the global coordinate system, and a vector t_i , which represents a translation.

$$T_i = (R_i, t_i) \quad (5)$$

This model was trained on Protein Data Bank, a freely accessible database containing 3D structures of larger biomolecules, including proteins and nucleic acids. The output is a distribution map containing the secondary structure and accessible surfaces predicted. After cross-validation of the results for the COVID-19 spike protein using the experimentally determined structure, they submitted predictions for proteins whose structure is not readily determined. These proteins have membrane proteins, proteins 3a, nsp2, nsp4, nsp6, and C-terminal domains such as papain. The structures of these proteins may represent docking sites for new drugs and therapeutics and could aid drug development in efforts to contain COVID-19. Utilizing a protein structure prediction AI program for the unique structure of the "mutant strain" of SARS-CoV-2 has the potential to change the way research is done in the field of

biology, allowing researchers to search for potential targets for new treatments before samples physically arrive.

AlphaFold2 has been released, making highly accurate protein structure prediction results available [105]. The premise of SBVS is that there is a reliable 3D protein structure, and as the 3D structures of proteins have been known so far, the targets for SBVS are naturally limited. In contrast, with the advent of AlphaFold, it has become possible to perform SBVS on proteins whose structures are unknown.

AlphaFold has significantly advanced the prediction of protein structures, including those of SARS-CoV-2, the virus responsible for COVID-19. The algorithm predicts the three-dimensional structures of proteins from their amino acid sequences, a process that traditionally requires extensive experimental techniques such as cryo-electron microscopy, nuclear magnetic resonance, and X-ray crystallography.

I. S protein

The SARS-CoV-2 spike (S) glycoprotein, which is the main target of antibodies, has been a primary focus for AlphaFold predictions. These predictions have helped elucidate the structural features of the spike protein, including its interaction with the angiotensin-converting enzyme 2 (ACE2) receptor, which is critical for the virus's entry into human cells [106]. AlphaFold's predictions have also been used to study the structural changes in different variants of the spike protein, such as those in the Omicron variant, to understand their impact on vaccine efficacy and viral transmission [107].

II. Other SARS-CoV-2 proteins

AlphaFold has also been used to predict the structures of several other SARS-CoV-2 proteins that are less well-studied but are essential for the virus's lifecycle. These include the membrane protein, Nsp2, Nsp4, Nsp6, and the papain-like proteinase (C-terminal domain) [106].

III. Methodology and Validation

AlphaFold employs a neural network architecture that integrates evolutionary, physical, and geometric constraints of protein structures. The algorithm uses multi-sequence alignments and a deep neural network to predict distances and angles between amino acids, achieving high accuracy even when no homologous structures are available [108]. The accuracy of AlphaFold's predictions has been validated by comparing them with experimentally determined structures, showing close agreement in many cases [106]. Thus, AlphaFold has revolutionized the field of protein structure prediction, particularly for SARS-CoV-2, by providing high-accuracy models that facilitate drug development and enhance our understanding of viral biology.

7. The application of AI in SARS-CoV-2-related proteins

AI has been extensively applied in various aspects related to SARS-CoV-2 proteins, particularly in COVID-19 drug discovery and vaccine design. Here are some examples of AI applications in this field:

- I) In prediction of vaccine candidates, AI tools like XGBoost have been used to predict vaccine candidates from non-structural proteins of SARS-CoV-2 [109].

- II) In prediction of HLA-binding peptides, feed-forward neural networks have been employed to predict HLA-binding peptides from the SARS-CoV-2 virus based on binding stability [109]
- III) As for the design of multiple-epitope vaccines, deep neural networks have been utilized for the prediction and design of multi-epitope vaccines that can manage the mutation of the virus [109].

These applications demonstrate how AI and machine learning play a crucial role in accelerating the discovery of effective drugs, vaccines, and treatment strategies for combating COVID-19 by leveraging the understanding of SARS-CoV-2-related proteins.

The computational techniques have been instrumental in various aspects related to SARS-CoV-2 research. These techniques have been applied in computational protein design for COVID-19 research, including the rapid design of peptides for detecting SARS-CoV-2 proteins [110]. *In silico* methods, computational tools, and bioinformatics resources have been utilized to annotate SARS-CoV-2 genomes and understand viral proteins [110]. Additionally, a computational study focused on cooperative binding to multiple SARS-CoV-2 proteins has been conducted, aiming to identify compounds with potential therapeutic effects through systems computational analysis [111]. These computational approaches play a crucial role in advancing our understanding of the virus and developing strategies for diagnosis and treatment.

The application of reinforcement learning, machine learning, and virtual screening in SARS-CoV-2-related proteins has shown promising results in identifying potential inhibitors for the virus. Studies have utilized machine learning-based virtual screening, molecular docking, and molecular dynamics simulations to identify novel compounds with the potential to inhibit key proteins like the main protease (Mpro) and papain-like protease (PLpro) of SARS-CoV-2 [97]. These approaches have led to discovering inhibitors effectively targeting these proteins, offering new avenues for developing antiviral agents. Additionally, *in silico* reinforcement learning has been employed to design spike/ACE2 inhibitory macrocycles, showcasing the use of AI in drug discovery for COVID-19. The combination of deep reinforcement learning, and virtual screening has been instrumental in optimizing hit molecules and developing effective non-covalent inhibitors for SARS-CoV-2 proteins. Furthermore, a novel protein design framework using reinforcement learning has been proposed to design a variant of the human ACE2 that binds more tightly to the SARS-CoV-2 S protein, potentially aiding in developing therapeutic solutions for COVID-19.

The application of AI in the study of SARS-CoV-2-related proteins has been a significant area of research, particularly in the context of the COVID-19 pandemic [112]. AI has been utilized in various domains including drug repurposing, structural biology, diagnostics, and vaccine development [113]. AI has played a crucial role in determining the structure of SARS-CoV-2 proteins. By predicting the structures of viral proteins, AI helps researchers understand the virus's mechanisms and identify potential targets for drug development.

I. Protein structure prediction and analysis

AI techniques have been used to predict and analyze the structure of SARS-CoV-2 proteins, which is crucial for understanding the virus and developing targeted

therapies. For example, deep learning models have been applied to predict protein structures and interactions [72].

II. Epitope prediction for vaccine design

AI algorithms have been employed to identify potential epitopes on SARS-CoV-2 proteins that could be targets for vaccine development. One study used computational analysis to compare SARS-CoV-2 nucleocapsid protein epitopes with those of related coronaviruses [114].

III. Drug target identification

AI-powered approaches have been used to identify potential drug targets among SARS-CoV-2 proteins. For instance, graph convolutional neural networks have been applied to predict drug-target interactions [114].

IV. Vaccine candidate ranking

Machine learning tools like Vaxign-ML have been developed to rank non-structural proteins as potential SARS-CoV-2 vaccine candidates using network-based algorithms [115].

V. Inhibitor discovery

AI has been utilized to rapidly screen large compound libraries to identify potential inhibitors of SARS-CoV-2 proteins. One study used deep docking to screen 1.3 billion compounds for potential inhibitors of the SARS-CoV-2 main protease [114].

VI. Antigenicity prediction

AI models have been used to predict the protective antigenicity of SARS-CoV-2 proteins. The spike (S) protein was found to have the highest protective antigenicity score [116].

VII. Immunogenic landscape prediction

AI techniques have been applied to predict the immunogenic landscape of SARS-CoV-2, which can guide universal vaccine design strategies [116]. These applications demonstrate how AI is being leveraged to accelerate research on SARS-CoV-2 proteins, potentially leading to faster development of effective drugs and vaccines against COVID-19. The integration of AI with biological and structural data has enabled researchers to rapidly analyze vast amounts of information and generate insights that can guide experimental work in the fight against the pandemic.

8. Conclusions

The SARS-CoV-2, like other coronaviruses, utilizes the S glycoprotein with S1 and S2 domains to enter host cells by binding to ACE2 receptors. Mutations in the S protein's RBD can enhance its affinity to ACE2. Searching for new compounds in COVID-19 involves trial-and-error experiments, but methods like deep reinforcement learning and structure-based virtual screening aid in drug discovery. AlphaFold, an AI system by DeepMind, predicts protein structures accurately by combining physical and biological approaches. It uses deep learning to predict 3D protein structures from amino acid sequences, achieving atomic accuracy even without homologous structures available.

The SARS-CoV-2 virus is very similar to the SARS-CoV virus that causes SARS, but several mutations in the RBR of the S protein greatly enhance the binding affinity of the SARS-CoV-2 virus to ACE2. The SARS-CoV-2 uses the S glycoprotein to enter

host cells, which has two functional domains: S1 and S2 RBD. New search methods using DNNs, such as methods using generative models such as VAE can learn the mapping between the discrete compound space and the continuous latent space by a generative model approximated by DNNs, and by allowing compound optimization to be performed in this continuous latent space. The learning process in the generative model is separated from the optimization process of the compound concerning the score function of the optimization target molecule. On the other hand, in methods based on reinforcement learning; by thinking of molecular design as a Markov decision process, the agent learns the optimal policy based on the rewards provided by the surrounding environment. By utilizing the AI program of a protein structure prediction for the unique structure of the “mutant strain” of SARS-CoV-2, it has the potential to search for potential targets for new drugs for SARS-CoV-2.

Disclaimer/Publisher’s note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.

Conflicts of interest: The authors declare no conflict of interest.

References

1. Yashavantha Rao HC, Jayabaskaran C. The emergence of a novel coronavirus (SARS-CoV-2) disease and their neuroinvasive propensity may affect in COVID-19 patients. *Journal of Medical Virology*. 2020; 92(7): 786-790. doi: 10.1002/jmv.25918
2. Ma Y, Deng J, Liu Q, et al. Long-Term Consequences of Asymptomatic SARS-CoV-2 Infection: A Systematic Review and Meta-Analysis. *International Journal of Environmental Research and Public Health*. 2023; 20(2): 1613. doi: 10.3390/ijerph20021613
3. Bongiovanni M, De Lauretis A, Manes G, et al. Clinical characteristics and outcome of COVID-19 pneumonia in elderly subjects. *Journal of Infection*. 2021; 82(2): e33-e34. doi: 10.1016/j.jinf.2020.08.023
4. Gupta SK, Minocha R, Thapa PJ, et al. Role of the Pangolin in Origin of SARS-CoV-2: An Evolutionary Perspective. *International Journal of Molecular Sciences*. 2022; 23(16): 9115. doi: 10.3390/ijms23169115
5. Yaşar Ş, Çolak C, Yoloğlu S. Artificial Intelligence-Based Prediction of Covid-19 Severity on the Results of Protein Profiling. *Computer Methods and Programs in Biomedicine*. 2021; 202: 105996. doi: 10.1016/j.cmpb.2021.105996
6. Dey L, Chakraborty S, Mukhopadhyay A. Machine learning techniques for sequence-based prediction of viral–host interactions between SARS-CoV-2 and human proteins. *Biomedical Journal*. 2020; 43(5): 438-450. doi: 10.1016/j.bj.2020.08.003
7. Cihan P, Ozger ZB. A new approach for determining SARS-CoV-2 epitopes using machine learning-based in silico methods. *Computational Biology and Chemistry*. 2022; 98: 107688. doi: 10.1016/j.compbiolchem.2022.107688
8. Alluwaimi AM, Alshubaith IH, Al-Ali AM, et al. The Coronaviruses of Animals and Birds: Their Zoonosis, Vaccines, and Models for SARS-CoV and SARS-CoV2. *Frontiers in Veterinary Science*. 2020; 7. doi: 10.3389/fvets.2020.582287
9. Kesheh MM, Hosseini P, Soltani S, et al. An overview on the seven pathogenic human coronaviruses. *Reviews in Medical Virology*. 2021; 32(2). doi: 10.1002/rmv.2282
10. Alexandersen S, Chamings A, Bhatta TR. SARS-CoV-2 genomic and subgenomic RNAs in diagnostic samples are not an indicator of active replication. *Nature Communications*. 2020; 11(1). doi: 10.1038/s41467-020-19883-7
11. Naqvi AAT, Fatima K, Mohammad T, et al. Insights into SARS-CoV-2 genome, structure, evolution, pathogenesis and therapies: Structural genomics approach. *Biochimica et Biophysica Acta (BBA) - Molecular Basis of Disease*. 2020; 1866(10): 165878. doi: 10.1016/j.bbadis.2020.165878

12. Chateau A, Van der Verren SE, Remaut H, et al. The Bacillus anthracis Cell Envelope: Composition, Physiological Role, and Clinical Relevance. *Microorganisms*. 2020; 8(12): 1864. doi: 10.3390/microorganisms8121864
13. Bai C, Zhong Q, Gao GF. Overview of SARS-CoV-2 genome-encoded proteins. *Science China Life Sciences*. 2021; 65(2): 280-294. doi: 10.1007/s11427-021-1964-4
14. Dërmaku-Sopjani M, Sopjani M. Interactions between ACE2 and SARS-CoV-2 S Protein: Peptide Inhibitors for Potential Drug Developments Against COVID-19. *Current Protein & Peptide Science*. 2021; 22(10): 729-744. doi: 10.2174/1389203722666210916141924
15. Wartecki A, Rzymiski P. On the Coronaviruses and Their Associations with the Aquatic Environment and Wastewater. *Water*. 2020; 12(6): 1598. doi: 10.3390/w12061598
16. Roy AN, Gupta AM, Banerjee D, et al. Unraveling DPP4 Receptor Interactions with SARS-CoV-2 Variants and MERS-CoV: Insights into Pulmonary Disorders via Immunoinformatics and Molecular Dynamics. *Viruses*. 2023; 15(10): 2056. doi: 10.3390/v15102056
17. Scialo F, Daniele A, Amato F, et al. ACE2: The Major Cell Entry Receptor for SARS-CoV-2. *Lung*. 2020; 198(6): 867-877. doi: 10.1007/s00408-020-00408-4
18. Shirbhate E, Pandey J, Patel VK, et al. Understanding the role of ACE-2 receptor in pathogenesis of COVID-19 disease: a potential approach for therapeutic intervention. *Pharmacological Reports*. 2021; 73(6): 1539-1550. doi: 10.1007/s43440-021-00303-6
19. Lan J, Ge J, Yu J, et al. Structure of the SARS-CoV-2 spike receptor-binding domain bound to the ACE2 receptor. *Nature*. 2020; 581(7807): 215-220. doi: 10.1038/s41586-020-2180-5
20. Huang Y, Yang C, Xu X feng, et al. Structural and functional properties of SARS-CoV-2 spike protein: potential antiviral drug development for COVID-19. *Acta Pharmacologica Sinica*. 2020; 41(9): 1141-1149. doi: 10.1038/s41401-020-0485-4
21. Li X, Yuan H, Li X, et al. Spike protein mediated membrane fusion during SARS-CoV-2 infection. *Journal of Medical Virology*. 2022; 95(1). doi: 10.1002/jmv.28212
22. Raghuvamsi PV, Tulsian NK, Samsudin F, et al. SARS-CoV-2 S protein: ACE2 interaction reveals novel allosteric targets. *eLife*. 2021; 10. doi: 10.7554/elife.63646
23. Belouzard S, Millet JK, Licitra BN, et al. Mechanisms of Coronavirus Cell Entry Mediated by the Viral Spike Protein. *Viruses*. 2012; 4(6): 1011-1033. doi: 10.3390/v4061011
24. Bosch BJ, Smits SL, Haagmans BL. Membrane ectopeptidases targeted by human coronaviruses. *Current Opinion in Virology*. 2014; 6: 55-60. doi: 10.1016/j.coviro.2014.03.011
25. Shang J, Wan Y, Luo C, et al. Cell entry mechanisms of SARS-CoV-2. *Proceedings of the National Academy of Sciences*. 2020; 117(21): 11727-11734. doi: 10.1073/pnas.2003138117
26. Harrison SC. Mechanism of membrane fusion by viral envelope proteins. *Adv Virus Res*. 2005; 64: 231-261. doi: 10.1016/S0065-3527(05)64007-9. PMID: 16139596.
27. Koppiseti RK, Fulcher YG, Van Doren SR. Fusion Peptide of SARS-CoV-2 Spike Rearranges into a Wedge Inserted in Bilayered Micelles. *Journal of the American Chemical Society*. 2021; 143(33): 13205-13211. doi: 10.1021/jacs.1c05435
28. Simmons G, Zmora P, Gierer S, et al. Proteolytic activation of the SARS-coronavirus spike protein: Cutting enzymes at the cutting edge of antiviral research. *Antiviral Research*. 2013; 100(3): 605-614. doi: 10.1016/j.antiviral.2013.09.028
29. Millet JK, Whittaker GR. Host cell entry of Middle East respiratory syndrome coronavirus after two-step, furin-mediated activation of the spike protein. *Proceedings of the National Academy of Sciences*. 2014; 111(42): 15214-15219. doi: 10.1073/pnas.1407087111
30. Takeda M. Proteolytic activation of SARS-CoV-2 spike protein. *Microbiology and Immunology*. 2021; 66(1): 15-23. doi: 10.1111/1348-0421.12945
31. Bertram S, Glowacka I, Müller MA, et al. Cleavage and Activation of the Severe Acute Respiratory Syndrome Coronavirus Spike Protein by Human Airway Trypsin-Like Protease. *Journal of Virology*. 2011; 85(24): 13363-13372. doi: 10.1128/jvi.05300-11
32. Chan YA, Zhan SH. The Emergence of the Spike Furin Cleavage Site in SARS-CoV-2. Kumar S, ed. *Molecular Biology and Evolution*. 2021; 39(1). doi: 10.1093/molbev/msab327
33. Hoffmann M, Kleine-Weber H, Pöhlmann S. A Multibasic Cleavage Site in the Spike Protein of SARS-CoV-2 Is Essential for Infection of Human Lung Cells. *Molecular Cell*. 2020; 78(4): 779-784. doi: 10.1016/j.molcel.2020.04.022

34. Jawad B, Adhikari P, Podgornik R, et al. Key Interacting Residues between RBD of SARS-CoV-2 and ACE2 Receptor: Combination of Molecular Dynamics Simulation and Density Functional Calculation. *Journal of Chemical Information and Modeling*. 2021; 61(9): 4425-4441. doi: 10.1021/acs.jcim.1c00560
35. Carvalho PPD, Alves NA. Featuring ACE2 binding SARS-CoV and SARS-CoV-2 through a conserved evolutionary pattern of amino acid residues. *Journal of Biomolecular Structure and Dynamics*. 2021; 40(22): 11719-11728. doi: 10.1080/07391102.2021.1965028
36. Yerukala Sathipati S, Shukla SK, Ho SY. Tracking the amino acid changes of spike proteins across diverse host species of severe acute respiratory syndrome coronavirus 2. *iScience*. 2022; 25(1): 103560. doi: 10.1016/j.isci.2021.103560
37. Zhai X, Sun J, Yan Z, et al. Comparison of Severe Acute Respiratory Syndrome Coronavirus 2 Spike Protein Binding to ACE2 Receptors from Human, Pets, Farm Animals, and Putative Intermediate Hosts. Gallagher T, ed. *Journal of Virology*. 2020; 94(15). doi: 10.1128/jvi.00831-20
38. Nour AM, Li Y, Wolenski J, et al. Viral Membrane Fusion and Nucleocapsid Delivery into the Cytoplasm are Distinct Events in Some Flaviviruses. Pierson TC, ed. *PLoS Pathogens*. 2013; 9(9): e1003585. doi: 10.1371/journal.ppat.1003585
39. V'kovski P, Kratzel A, Steiner S, et al. Coronavirus biology and replication: implications for SARS-CoV-2. *Nature Reviews Microbiology*. 2020; 19(3): 155-170. doi: 10.1038/s41579-020-00468-6
40. Ahlquist P, Noueiry AO, Lee WM, et al. Host Factors in Positive-Strand RNA Virus Genome Replication. *Journal of Virology*. 2003; 77(15): 8181-8186. doi: 10.1128/jvi.77.15.8181-8186.2003
41. Upadhyay M, Gupta S. Endoplasmic reticulum secretory pathway: Potential target against SARS-CoV-2. *Virus Research*. 2022; 320: 198897. doi: 10.1016/j.virusres.2022.198897
42. Scherer KM, Mascheroni L, Carnell GW, et al. SARS-CoV-2 nucleocapsid protein adheres to replication organelles before viral assembly at the Golgi/ERGIC and lysosome-mediated egress. *Science Advances*. 2022; 8(1). doi: 10.1126/sciadv.abl4895
43. Siu YL, Teoh KT, Lo J, et al. The M, E, and N Structural Proteins of the Severe Acute Respiratory Syndrome Coronavirus Are Required for Efficient Assembly, Trafficking, and Release of Virus-Like Particles. *Journal of Virology*. 2008; 82(22): 11318-11330. doi: 10.1128/jvi.01052-08
44. Villanueva RA, Rouillé Y, Dubuisson J. Interactions between virus proteins and host cell membranes during the viral life cycle. *Int Rev Cytol*. 2005; 245: 171-244. doi: 10.1016/S0074-7696(05)45006-8. PMID: 16125548
45. Seltzer S. Linking ACE2 and angiotensin II to pulmonary immunovascular dysregulation in SARS-CoV-2 infection. *International Journal of Infectious Diseases*. 2020; 101: 42-45. doi: 10.1016/j.ijid.2020.09.041
46. Burrell LM, Johnston CI, Tikellis C, et al. ACE2, a new regulator of the renin-angiotensin system. *Trends in Endocrinology & Metabolism*. 2004; 15(4): 166-169. doi: 10.1016/j.tem.2004.03.001
47. Silhol F, Sarlon G, Deharo JC, et al. Downregulation of ACE2 induces overstimulation of the renin-angiotensin system in COVID-19: should we block the renin-angiotensin system? *Hypertension Research*. 2020; 43(8): 854-856. doi: 10.1038/s41440-020-0476-3
48. Chappell MC. Biochemical evaluation of the renin-angiotensin system: the good, bad, and absolute? *American Journal of Physiology-Heart and Circulatory Physiology*. 2016; 310(2): H137-H152. doi: 10.1152/ajpheart.00618.2015
49. Tamura K, Wakui H, Azushima K, et al. Angiotensin II Type 1 Receptor Binding Molecule ATRAP as a Possible Modulator of Renal Sodium Handling and Blood Pressure in Pathophysiology. *Current Medicinal Chemistry*. 2015; 22(28): 3210-3216. doi: 10.2174/0929867322666150821095036
50. Blaustein MP, Leenen FHH, Chen L, et al. How NaCl raises blood pressure: a new paradigm for the pathogenesis of salt-dependent hypertension. *American Journal of Physiology-Heart and Circulatory Physiology*. 2012; 302(5): H1031-H1049. doi: 10.1152/ajpheart.00899.2011
51. Pratiwi A, Hakim TR, Abidin MZ, et al. Angiotensin-converting enzyme inhibitor activity of peptides derived from Kacang goat skin collagen through thermolysin hydrolysis. *January-2021*. 2021; 14(1): 161-167. doi: 10.14202/vetworld.2021.161-167
52. Karnik SS, Singh KD, Tirupula K, et al. Significance of angiotensin 1-7 coupling with MAS1 receptor and other GPCRs to the renin-angiotensin system: IUPHAR Review 22. *British Journal of Pharmacology*. 2017; 174(9): 737-753. doi: 10.1111/bph.13742
53. Santos RA. Angiotensin-(1-7). *Hypertension*. 2014; 63(6): 1138-1147. doi: 10.1161/hypertensionaha.113.01274

54. Bosso M, Thanaraj TA, Abu-Farha M, et al. The Two Faces of ACE2: The Role of ACE2 Receptor and Its Polymorphisms in Hypertension and COVID-19. *Molecular Therapy - Methods & Clinical Development*. 2020; 18: 321-327. doi: 10.1016/j.omtm.2020.06.017
55. Valente J, António J, Mora C, et al. Developments in Image Processing Using Deep Learning and Reinforcement Learning. *Journal of Imaging*. 2023; 9(10): 207. doi: 10.3390/jimaging9100207
56. Pudjihartono N, Fadason T, Kempa-Liehr AW, et al. A Review of Feature Selection Methods for Machine Learning-Based Disease Risk Prediction. *Frontiers in Bioinformatics*. 2022; 2. doi: 10.3389/fbinf.2022.927312
57. De Teyou GK, Tarabalka Y, Manighetti I, et al. Deep Neural Networks for automatic extraction of features in time series satellite images. Available online: <https://arxiv.org/abs/2008.08432> (accessed on 17 May 2024).
58. Sodhani S, Faramarzi M, Mehta SV, et al. An Introduction to Lifelong Supervised Learning. Available online: <https://arxiv.org/abs/2207.04354> (accessed on 17 May 2024).
59. Yang R. Unsupervised machine learning for physical concepts. Available online: <https://arxiv.org/abs/2205.05279> (accessed on 17 May 2024).
60. Goel D, Neumann A, Neumann F, et al. Evolving Reinforcement Learning Environment to Minimize Learner's Achievable Reward: An Application on Hardening Active Directory Systems. Available online: <https://arxiv.org/abs/2304.03998> (accessed on 17 May 2024).
61. Chitnis R, Xu Y, Hashemi B, et al. IQL-TD-MPC: Implicit Q-Learning for Hierarchical Model Predictive Control. Available online: <https://arxiv.org/abs/2306.00867> (accessed on 17 May 2024).
62. Neufeld A, Sester J. Robust Q-learning Algorithm for Markov Decision Processes under Wasserstein Uncertainty. Available online: <https://arxiv.org/abs/2210.00898> (accessed on 17 May 2024)
63. Ronecker MP, Zhu Y. Deep Q-Network Based Decision Making for Autonomous Driving. Available online: <https://arxiv.org/abs/2303.11634> (accessed on 17 May 2024).
64. Kadurin A, Nikolenko S, Khrabrov K, et al. druGAN: An Advanced Generative Adversarial Autoencoder Model for de Novo Generation of New Molecules with Desired Molecular Properties in Silico. *Molecular Pharmaceutics*. 2017; 14(9): 3098-3104. doi: 10.1021/acs.molpharmaceut.7b00346
65. Dai W, Guo D. A Ligand-Based Virtual Screening Method Using Direct Quantification of Generalization Ability. *Molecules*. 2019; 24(13): 2414. doi: 10.3390/molecules24132414
66. Maia EHB, Assis LC, de Oliveira TA, et al. Structure-Based Virtual Screening: From Classical to Artificial Intelligence. *Frontiers in Chemistry*. 2020; 8. doi: 10.3389/fchem.2020.00343
67. Tran-Nguyen VK, Junaid M, Simeon S, et al. A practical guide to machine-learning scoring for structure-based virtual screening. *Nature Protocols*. 2023; 18(11): 3460-3511. doi: 10.1038/s41596-023-00885-w
68. Wu C, Liu Y, Yang Y, et al. Analysis of therapeutic targets for SARS-CoV-2 and discovery of potential drugs by computational methods. *Acta Pharmaceutica Sinica B*. 2020; 10(5): 766-788. doi: 10.1016/j.apsb.2020.02.008
69. Fassihi A, Hatami S, Sirous H, et al. Preparing a database of corrected protein structures important in cell signaling pathways. *Research in Pharmaceutical Sciences*. 2023; 18(1): 67. doi: 10.4103/1735-5362.363597
70. Revillo Imbernon J, Chiesa L, Kellenberger E. Mining the Protein Data Bank to inspire fragment library design. *Frontiers in Chemistry*. 2023; 11. doi: 10.3389/fchem.2023.1089714
71. Junk P, Kiel C. HOMELETTE: a unified interface to homology modelling software. Valencia A, ed. *Bioinformatics*. 2021; 38(6): 1749-1751. doi: 10.1093/bioinformatics/btab866
72. Yang Z, Zeng X, Zhao Y, et al. AlphaFold2 and its applications in the fields of biology and medicine. *Signal Transduction and Targeted Therapy*. 2023; 8(1). doi: 10.1038/s41392-023-01381-z
73. Konc J, Janežič D. Protein binding sites for drug design. *Biophysical Reviews*. 2022; 14(6): 1413-1421. doi: 10.1007/s12551-022-01028-3
74. Alzyoud L, Bryce RA, Al Sorkhy M, et al. Structure-based assessment and druggability classification of protein-protein interaction sites. *Scientific Reports*. 2022; 12(1). doi: 10.1038/s41598-022-12105-8
75. Piazza I, Beaton N, Bruderer R, et al. A machine learning-based chemoproteomic approach to identify drug targets and binding sites in complex proteomes. *Nature Communications*. 2020; 11(1). doi: 10.1038/s41467-020-18071-x
76. Rufer AC. Drug discovery for enzymes. *Drug Discovery Today*. 2021; 26(4): 875-886. doi: 10.1016/j.drudis.2021.01.006
77. Farooq Q ul A, Shaikat Z, Aiman S, et al. Protein-protein interactions: Methods, databases, and applications in virus-host study. *World Journal of Virology*. 2021; 10(6): 288-300. doi: 10.5501/wjv.v10.i6.288

78. Matthew AN, Leidner F, Lockbaum GJ, et al. Drug Design Strategies to Avoid Resistance in Direct-Acting Antivirals and Beyond. *Chemical Reviews*. 2021; 121(6): 3238-3270. doi: 10.1021/acs.chemrev.0c00648
79. Wang Y, Wei Z, Xi L. Sfcnn: a novel scoring function based on 3D convolutional neural network for accurate and stable protein–ligand affinity prediction. *BMC Bioinformatics*. 2022; 23(1). doi: 10.1186/s12859-022-04762-3
80. Kosugi T, Ohue M. Quantitative Estimate Index for Early-Stage Screening of Compounds Targeting Protein-Protein Interactions. *International Journal of Molecular Sciences*. 2021; 22(20): 10925. doi: 10.3390/ijms222010925
81. Eberhardt J, Santos-Martins D, Tillack AF, et al. AutoDock Vina 1.2.0: New Docking Methods, Expanded Force Field, and Python Bindings. *Journal of Chemical Information and Modeling*. 2021; 61(8): 3891-3898. doi: 10.1021/acs.jcim.1c00203
82. Ropp PJ, Kaminsky JC, Yablonski S, et al. Dimorphite-DL: an open-source program for enumerating the ionization states of drug-like small molecules. *Journal of Cheminformatics*. 2019; 11(1). doi: 10.1186/s13321-019-0336-9
83. Nash S, Vachet RW. Gas-Phase Unfolding of Protein Complexes Distinguishes Conformational Isomers. *Journal of the American Chemical Society*. 2022; 144(48): 22128-22139. doi: 10.1021/jacs.2c09573
84. Cleves AE, Jain AN. Structure- and Ligand-Based Virtual Screening on DUD-E+: Performance Dependence on Approximations to the Binding Pocket. *Journal of Chemical Information and Modeling*. 2020; 60(9): 4296-4310. doi: 10.1021/acs.jcim.0c00115
85. Tang S, Chen R, Lin M, et al. Accelerating AutoDock Vina with GPUs. *Molecules*. 2022; 27(9): 3041. doi: 10.3390/molecules27093041
86. Jumper J, Hassabis D. Protein structure predictions to atomic accuracy with AlphaFold. *Nature Methods*. 2022; 19(1): 11-12. doi: 10.1038/s41592-021-01362-6
87. Chen T, Shu X, Zhou H, et al. Algorithm selection for protein–ligand docking: strategies and analysis on ACE. *Scientific Reports*. 2023; 13(1). doi: 10.1038/s41598-023-35132-5
88. Mohammadi S, Narimani Z, Ashouri M, et al. Ensemble learning from ensemble docking: revisiting the optimum ensemble size problem. *Scientific Reports*. 2022; 12(1). doi: 10.1038/s41598-021-04448-5
89. Verburt J, Kihara D. Benchmarking of structure refinement methods for protein complex models. *Proteins: Structure, Function, and Bioinformatics*. 2021; 90(1): 83-95. doi: 10.1002/prot.26188
90. Peivaste I, Ramezani S, Alahyarizadeh G, et al. Rapid and accurate predictions of perfect and defective material properties in atomistic simulation using the power of 3D CNN-based trained artificial neural networks. *Scientific Reports*. 2024; 14(1). doi: 10.1038/s41598-023-50893-9
91. Aziz S, Waqas M, Mohanta TK, et al. Identifying non-nucleoside inhibitors of RNA-dependent RNA-polymerase of SARS-CoV-2 through per-residue energy decomposition-based pharmacophore modeling, molecular docking, and molecular dynamics simulation. *Journal of Infection and Public Health*. 2023; 16(4): 501-519. doi: 10.1016/j.jiph.2023.02.009
92. Gazi R, Maity S, Jana M. Conformational Features and Hydration Dynamics of Proteins in Cosolvents: A Perspective from Computational Approaches. *ACS Omega*. 2023; 8(3): 2832-2843. doi: 10.1021/acsomega.2c08009
93. Wang X, Chong B, Sun Z, et al. More is simpler: Decomposition of ligand-binding affinity for proteins being disordered. *Protein Science*. 2022; 31(7). doi: 10.1002/pro.4375
94. Kurniawan J, Ishida T. Protein Model Quality Estimation Using Molecular Dynamics Simulation. *ACS Omega*. 2022; 7(28): 24274-24281. doi: 10.1021/acsomega.2c01475
95. Li Y, Hou S, Zhang Y, et al. Effect of Travel Restrictions of Wuhan City Against COVID-19: A Modified SEIR Model Analysis. *Disaster Medicine and Public Health Preparedness*. 2021; 16(4): 1431-1437. doi: 10.1017/dmp.2021.5
96. Chakraborty M, Shakir Mahmud M, Gates TJ, et al. Analysis and Prediction of Human Mobility in the United States during the Early Stages of the COVID-19 Pandemic using Regularized Linear Models. *Transportation Research Record: Journal of the Transportation Research Board*. 2022; 2677(4): 380-395. doi: 10.1177/03611981211067794
97. Samad A, Ajmal A, Mahmood A, et al. Identification of novel inhibitors for SARS-CoV-2 as therapeutic options using machine learning-based virtual screening, molecular docking and MD simulation. *Frontiers in Molecular Biosciences*. 2023; 10. doi: 10.3389/fmolb.2023.1060076
98. Zhang L, Zhao H, Liu J, et al. Design of SARS-CoV-2 Mpro, PLpro Dual-Target Inhibitors Based on Deep Reinforcement Learning and Virtual Screening. *Future Medicinal Chemistry*. 2022; 14(6): 393-405. doi: 10.4155/fmc-2021-0269
99. Higgins MK. Can We AlphaFold Our Way Out of the Next Pandemic? *Journal of Molecular Biology*. 2021; 433(20): 167093. doi: 10.1016/j.jmb.2021.167093

100. Ismi DP, Pulungan R, Afiahayati. Deep learning for protein secondary structure prediction: Pre and post-AlphaFold. *Computational and Structural Biotechnology Journal*. 2022; 20: 6271-6286. doi: 10.1016/j.csbj.2022.11.012
101. Marcu ȘB, Tăbircă S, Tangney M. An Overview of AlphaFold's Breakthrough. *Frontiers in Artificial Intelligence*. 2022; 5. doi: 10.3389/frai.2022.875587
102. Bertoline LMF, Lima AN, Krieger JE, et al. Before and after AlphaFold2: An overview of protein structure prediction. *Frontiers in Bioinformatics*. 2023; 3. doi: 10.3389/fbinf.2023.1120370
103. DeBenedictis EA, Chory EJ, Gretton DW, et al. Systematic molecular evolution enables robust biomolecule discovery. *Nature Methods*. 2021; 19(1): 55-64. doi: 10.1038/s41592-021-01348-4
104. Xu YC, ShangGuan TJ, Ding XM, et al. Accurate prediction of protein torsion angles using evolutionary signatures and recurrent neural network. *Scientific Reports*. 2021; 11(1). doi: 10.1038/s41598-021-00477-2
105. Kilim O, Mentés A, Pál B, et al. SARS-CoV-2 receptor-binding domain deep mutational AlphaFold2 structures. *Scientific Data*. 2023; 10(1). doi: 10.1038/s41597-023-02035-z
106. Gutnik D, Evseev P, Miroshnikov K, Shneider M. Using AlphaFold Predictions in Viral Research. *Curr Issues Mol Biol*. 2023; 45: 3705-3732. doi: 10.3390/cimb45040240
107. Ali MA, Caetano-Anollés G. AlphaFold2 Reveals Structural Patterns of Seasonal Haplotype Diversification in SARS-CoV-2 Spike Protein Variants. *Biology*. 2024; 13(3): 134. doi: 10.3390/biology13030134
108. Jumper J, Evans R, Pritzel A, et al. Highly accurate protein structure prediction with AlphaFold. *Nature*. 2021; 596(7873): 583-589. doi: 10.1038/s41586-021-03819-2
109. Lv H, Shi L, Berkenpas JW, et al. Application of artificial intelligence and machine learning for COVID-19 drug discovery and vaccine design. *Brief Bioinform*. 2021; 22: 320. doi: 10.1093/bib/bbab320
110. Kalita P, Tripathi T, Padhi AK. Computational Protein Design for COVID-19 Research and Emerging Therapeutics. *ACS Central Science*. 2023; 9(4): 602-613. doi: 10.1021/acscentsci.2c01513
111. Li J, McKay KT, Remington JM, et al. A computational study of cooperative binding to multiple SARS-CoV-2 proteins. *Scientific Reports*. 2021; 11(1). doi: 10.1038/s41598-021-95826-6
112. Ashique S, Mishra N, Mohanto S, et al. Application of artificial intelligence (AI) to control COVID-19 pandemic: Current status and future prospects. *Heliyon*. 2024; 10(4): e25754. doi: 10.1016/j.heliyon.2024.e25754
113. Prasad K, Kumar V. Artificial intelligence-driven drug repurposing and structural biology for SARS-CoV-2. *Current Research in Pharmacology and Drug Discovery*. 2021; 2: 100042. doi: 10.1016/j.crphar.2021.100042
114. Keshavarzi Arshadi A, Webb J, Salem M, et al. Artificial Intelligence for COVID-19 Drug Discovery and Vaccine Development. *Frontiers in Artificial Intelligence*. 2020; 3. doi: 10.3389/frai.2020.00065
115. Ghosh A, Larrondo-Petrie MM, Pavlovic M. Revolutionizing Vaccine Development for COVID-19: A Review of AI-Based Approaches. *Information*. 2023; 14(12): 665. doi: 10.3390/info14120665
116. Wang L, Zhang Y, Wang D, et al. Artificial Intelligence for COVID-19: A Systematic Review. *Frontiers in Medicine*. 2021; 8. doi: 10.3389/fmed.2021.704256



Academic Publishing Pte. Ltd.

Add: 73 Upper Paya Lebar Road #07-02B-01 Centro Bianco Singapore 534818

Tel: +65 83184869

E-mail: editorial_office@acad-pub.com

Web: <http://ojs.acad-pub.com/>