

# Accurate ECG images classification using Vision Transformer

Fatima Mallak Hanoon, Khawla Hussein Ali\*

Department of Computers, faculty of Education for Pure Sciences, University of Basra, Iraq.

## ARTICLE INFO

Received 06 January 2024  
Accepted 5 March 2024  
Published 30 June 2024

## Keywords :

Electrocardiogram, Classification,  
Vision transformer, Deep Learning

**Citation:** Fatima M. Hanoon, Khawla H. Ali, J. Basrah Res. (Sci.) 50(1), 328 (2024). DOI:<https://doi.org/10.56714/bjrs.50.1.26>

## ABSTRACT

Electrocardiogram (ECG) classification plays a crucial role in the diagnosis and management of cardiovascular diseases. Deep learning-based approaches have shown promising results in automated ECG classification. However, the complexity of ECG signals, including variations in morphology, duration, and amplitude, poses significant challenges for existing deep learning models. In this regard, recent advancements in vision transformer models have shown remarkable performance in images processing and computer vision tasks. In this paper, we propose a deep vision transformer-based approach for ECG classification, which combines the power of convolutional neural networks and self-attention mechanisms. Our proposed model was tuned and enhanced by four hyper-parameters of the proposed model, it can effectively detect internally the main features of ECG images and achieve performance on benchmark ECG datasets. The proposed model can aid in the early detection and diagnosis of cardiovascular diseases, thus improving patient outcomes the final accuracy was 98.23% in dataset <https://data.mendeley.com/datasets/gwbz3fsgp8/2>.

## 1. Introduction

The ElectroCardioGram (ECG) is one of the most important vital signals that is extracted from the human body, as it is the basic diagnostic signal that doctors rely on - in most of their specialties - in the initial and sometimes final diagnosis of a wide range of diseases **Error! Reference source not found.**, Fig.1 shows a normal heartbeat.

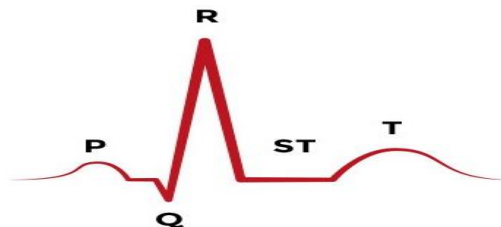


Fig.1: normal heartbeat

\*Corresponding author email: [alsalihfatimah4@gmail.com](mailto:alsalihfatimah4@gmail.com)



©2022 College of Education for Pure Science, University of Basrah. This is an Open Access Article Under the CC by License the [CC BY 4.0](https://creativecommons.org/licenses/by/4.0/) license.

N: 1817-2695 (Print); 2411-524X (Online)  
line at: <https://jou.jobrs.edu.iq>

The electrical signal (ECG) changes in many heart diseases or when a malfunction in one of the other organs or systems in the human body causes a change in the heart's electrical function. Which makes the doctor notice this change and suspect certain diseases, or even sometimes diagnose this disease based on the quoted reference. Fig. 2 shows some cases of heartbeats **Error! Reference source not found.**

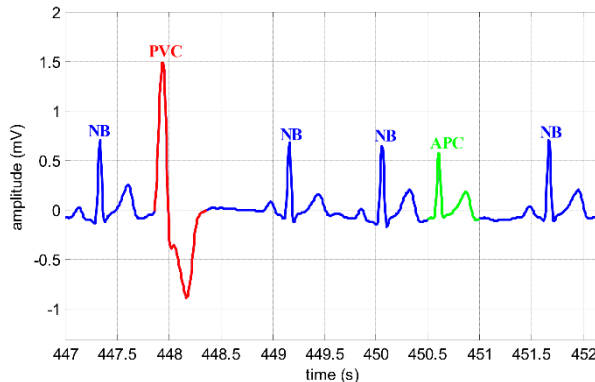


Fig. 2: Different types of heartbeats

Vision Transformers, or ViT, are A new advancement in the field of computer vision that have shown remarkable performance in various visual recognition tasks. Unlike traditional Convolutional neural networks (CNNs) that utilize convolutional operations layers to extract visual features, ViTs employ a self-attention mechanism to directly model the relationships between different image regions. This allows ViTs are designed to capture long-range dependencies and global context information, which can be crucial for accurate image recognition **Error! Reference source not found..**

The architecture of ViTs is based on the Transformer, a powerful sequence modeling architecture originally introduced for natural language processing. The Transformer consists of self-attention layers and feedforward layers, which allow it to process input sequences of variable length and capture long-term dependencies. ViTs adapt the Transformer architecture for image recognition by dividing the image into a **Page Dimensions**

All material on all pages should fit within an area of A4 (21 x 29.7 cm), 2.8 cm from the top of the page and ending with 2.4 cm from the bottom. The left and right margins should both be 2.4 cm.

## 2. Main Text

This section provides details for typesetting your manuscript according to the formatting guidelines set for JOBRs Journal. Use 11-point Times New Roman regular font for typesetting of the main text in the document.

The main text starts at the top of the page and continues in a one-column format. Place an indentation for each paragraph starting from the first in all sections or subsections. There is no space between paragraphs within the text. Add an 11-point space after the text in each section or subsection.

sequence of non-overlapping patches, Treating every patch as a token in the input sequence. The self-attention mechanism in ViTs then operates on these patches to capture global relationships and produce a feature representation for the entire image **Error! Reference source not found..**

ViTs have achieved State-of-the-art achievement in various benchmarks, such as the ImageNet classification task and the COCO object detection task. They have also shown promise in other computer vision tasks such as semantic segmentation and video classification. ViTs are particularly attractive because they require fewer parameters than traditional CNNs, which can result in faster training times and reduced memory requirements. Vision Transformers represent an innovative method for image recognition, harnessing the strength of self-attention to capture global context and long-range connections.

They have demonstrated impressive results in various computer vision tasks and are a promising direction for future research in the field **Error! Reference source not found.**

The choice between Vision Transformers and Convolutional Neural Networks (CNNs) hinges on various factors, such as the dataset size and the complexity of the visual recognition task and the available computing resources.

CNNs have been the go-to method for image recognition for many years and are highly effective at capturing local patterns and features in images. They are well-suited for datasets with a large number of images and a moderate number of classes. CNNs are also highly efficient in terms of computational requirements, making them ideal for real-time applications **Error! Reference source not found.**[7][8].

On the contrary, Vision Transformers are a relatively recent development and have demonstrated outstanding performance across a range of computer vision tasks, such as image classification, object detection, and semantic segmentation. They are particularly effective at capturing long-range dependencies and global context information in images, which can be crucial for accurate recognition of complex visual patterns. However, ViTs require significantly more computational resources than CNNs, making them more suitable for larger datasets and more complex visual recognition tasks **Error! Reference source not found.**

In summary, CNNs are an excellent choice for image recognition tasks that involve large datasets and moderate complexity[10][11], while Vision Transformers are ideal for tasks that require capturing long-range dependencies and global context information in images. The decision between the two approaches ultimately relies on the precise requirements of the current task and the available computing resources **Error! Reference source not found.**

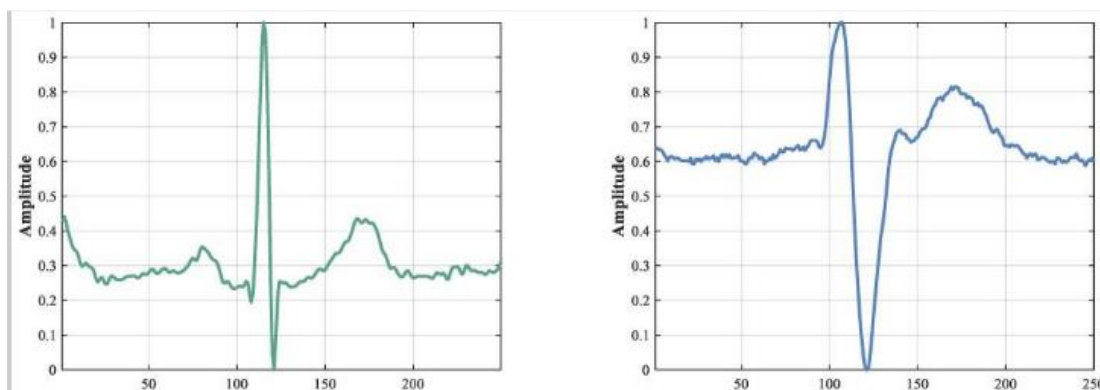
### State of Art

The authors in **Error! Reference source not found.** present a method for detecting congestive heart failure (CHF) in patients using a combination of electrocardiogram (ECG) data and deep learning techniques. Specifically, the authors propose using a convolutional neural network (CNN) to extract features from the ECG signals, and a vision transformer The proposed approach aims to train a network to discover associations between the identified features and the presence or absence of CH is shown to outperform existing approaches to CHF detection in terms of accuracy and generalization to new patients.

In **Error! Reference source not found.**, authors propose a new framework for detecting valvular heart diseases (VHD) using phonocardiography (PCG) signals and vision transformer (ViT) models. The proposed framework includes pre-processing steps to extract the PCG signal and feature extraction using a modified ViT model. The modified ViT model is trained on a large dataset of PCG recordings and can effectively learn the characteristic patterns of different types of VHDs. The proposed framework achieves high accuracy and efficiency in VHD detection, outperforming other state-of-the-art methods. The article concludes that the proposed framework has the potential to be an effective tool for early VHD detection in clinical settings.

Authors in **Error! Reference source not found.** joined ViT with deep learning algorithm with minimal adjustments to current ViT models, the suggested strategy enables the training of deeper ViT models with steady performance gains. An interesting finding is that the Top-1 classification accuracy on ImageNet may be increased by 1.6% by training a deep ViT model with 32 transformer blocks. The dataset was the **ImageNet**, it is general dataset.

Authors in **Error! Reference source not found.** studied Congestive heart failure (CHF) attacks that can result in symptoms including breathing problems, lightheadedness, or weariness. To diagnose CHF, an (ECG) is a quick and affordable procedure. Misdiagnosis occurs often as a result of the ECGs' intrinsic complexity and the waveforms' minute variations. They provided a novel technique to diagnose CHF utilizing the ECG-Convolution-Vision Transformer (ECVT-Net). They combined the traits of a CNN and a Vision Transformer to extract high-dimensional features from ECGs. The results accuracy was 98.88%. Their model has resistance to noise.



**Fig. 3:** Two images from **Error! Reference source not found.**

The learning of multi-scale feature representations in transformer models for image classification was explored by the authors in **Error! Reference source not found.**. To integrate picture patches of different sizes and create stronger image characteristics, they suggested a dual-branch transformer. Using two different branches with different computational complexities, their approach handles both small-patch and large-patch tokens. Through cross-attention, they develop a simple yet effective token fusion module that utilizes a single token for each branch as a query to communicate with other branches. Their method fared better than a number of works on vision transformer. The dataset was the ImageNet, it is general dataset.

The authors **Error! Reference source not found.** in employed ViT in video classification, they offered models for categorizing videos based just on transforms. Their methodology takes the input video and extracts spatiotemporal tokens, The input is processed by multiple transformer layers, which enable the management of long token sequences in video by factoring the spatial and temporal dimensions of the input. Various effective model versions have been proposed for this purpose. The approach successfully regularizes the model during training and uses pretrained image models to be able to train on very little datasets, even though transformer-based models are known to only be effective when huge training datasets are available. They conducted extensive ablation investigations and outperformed approaches based on deep 3D models. The datasets were evaluated on several video classification benchmarks, such as Kinetics 400 and 600, Epic Kitchens, and Something-Something v2

#### **Objective of this paper**

Problem of classification heartbeats depending on images might have two ways: processing images and get features as table or processing images using pixels values and this done in deep learning. In this approach an optimization step is applied to best selection of some parameters in ViT. ViT has many parameters that can be optimized as will be explained in next paragraph. The classification is binary classification (Normal and Abnormal)

### **3. Methodology**

The ViT model consists of the following components **Error! Reference source not found.****Error! Reference source not found.:**

**Patch Embeddings:** In order to feed images In the ViT model, the input is initially divided into a grid of non-overlapping patches, with each patch being considered as an individual token. These patches are subsequently flattened into a 1D sequence and passed through a linear layer to convert them into a sequence of embeddings.

**Positional Embeddings:** Since the ViT model doesn't use convolutional layers like traditional image classification models, it needs a way to incorporate spatial information into its representations. To do this, the model adds learnable positional

**Transformer Encoder:** The core of the ViT model is a stack of Transformer encoder layers, which process the sequence of patch embeddings and positional embeddings. The encoder layer comprises a multi-head self-attention mechanism and a feedforward neural network .

**Multi-head self-attention:** This sublayer is responsible for processing the input embeddings to capture the most relevant contextual information. In this sublayer, each input embedding attends to all other embeddings in the same sequence, calculating a weighted sum of them based on their relevance to

the current embedding. This is done using a mechanism called self-attention, which is computed in parallel for multiple "heads" of attention. Each head is responsible for learning a different representation of the input sequence, which can capture different types of dependencies and relationships.

**MLP: Multi-Layer Perception or Position-wise feedforward network:** After the self-attention mechanism, the input embeddings are passed through a position-wise feedforward network. This sublayer applies a fully connected feedforward neural network to each position in the sequence independently. This network consists of two linear transformations, separated by a non-linear activation function (typically ReLU), and is applied to each position in the sequence independently

**Classification Head:** At the top of the ViT model is a classification head, which takes the output of the last Transformer encoder layer and maps it to class probabilities using a fully-connected layer and a softmax activation function.

**Last Layer is MLP:** The Multilayer Perceptron (MLP) serves as a robust tool for addressing pattern recognition, classification, and regression tasks. It has found widespread use across various applications, including image recognition, speech recognition, and natural language processing. Notably, MLP is incorporated as the final layer in the Vision Transformer (ViT). Figure 4 illustrates the general structure of ViT, while Figure 5 depicts deep learning with ViT and the potential incorporation of optimization within the complete model according to our approach

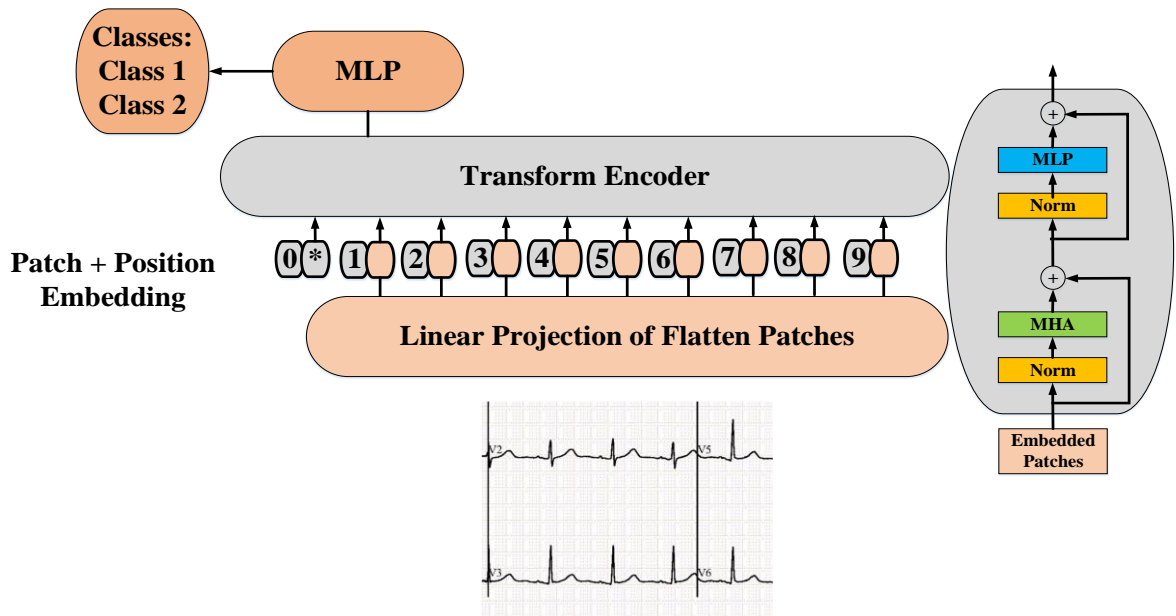
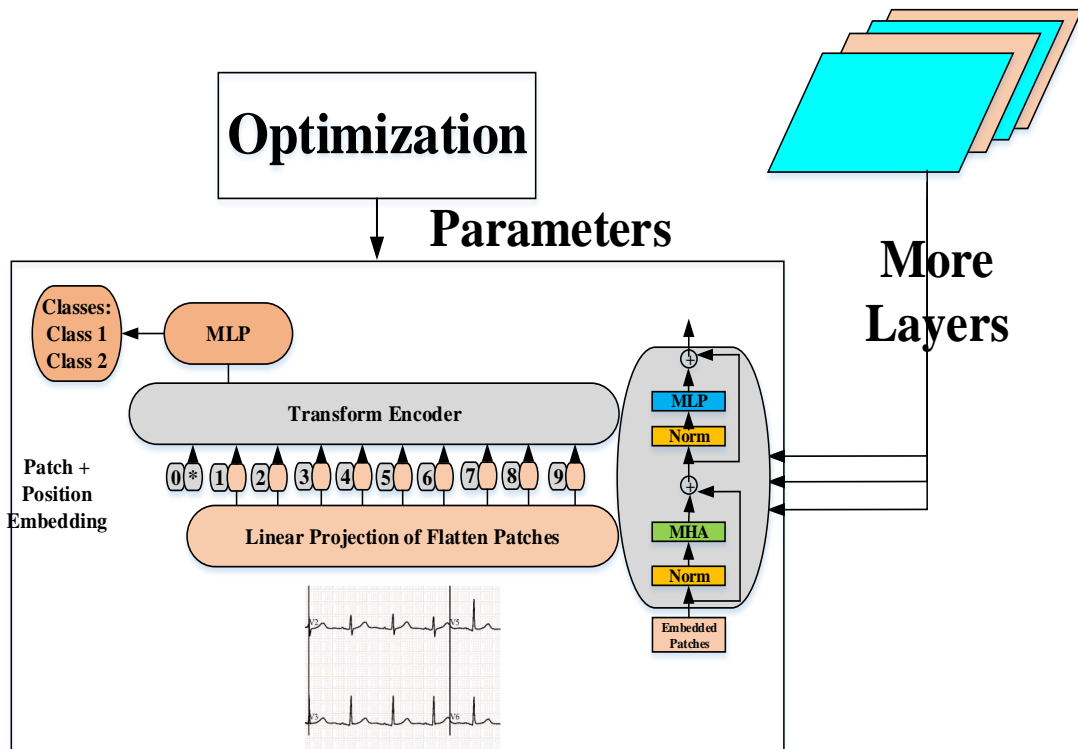


Fig. 4: General structure of ViT



**Fig.5:** Deep ViT with optimization

Algorithm :proposed Deep ViT

Input: Training and learning imaging :  $\{x_i, y_i\}_{i=1}^k$

Out put: predivted labels of ECG dataset images

1. Set batchsize to 100, Optimizer Adam (learning rate: 0.00015), number of iterations to 20000, image dimensions to (535,941, 3)
2. Set the number of mini-batches as:  $nb = n/batchsize$
3. For iteration = 1: Number of iterations
  - 3.1 For batch = 1 : nb
    - Pick a batch from the training set,
    - Incrises another batch of augmented images using a particular augmentation method,
    - Train the model on the original and augmented images by minimizing the SparseCategoricalCrossentropy loss.
    - Backpropagate the loss.Update the model parameters.
4. Classify test images

**Deep learning injection**

The parts where deep learning can be noticed are MLP basically. MLP is used in transformer encoder, default ViT has two layers of MLP. Each layer has number of nodes; thus if MLP has three layers and first layer has 64 nodes, second layer has 128 layers and the last layer has 256 nodes, so the total number of nodes is  $64+128+256=448$  nodes.

The parts where deep learning can be noticed are: MLP: the number of nodes was increased in transform encoder by two ways, increasing the number of layers and the number of nodes in each layer. Each MLP in ViT was increased by this way.

**Tuning hyper parameters**

There are several methods for tuning any hyper-parameter; suppose the hyper-parameter is learning rate:

**Grid search:** This involves trying out different learning rates over a grid of values and evaluating the performance of the model on a validation set. The learning rate that gives the best performance is then selected.

**Random search:** This involves randomly sampling learning rates from a distribution and evaluating the performance of the model on a validation set. The learning rate that gives the best performance is then selected.

**Adaptive methods:** These are optimization algorithms that adapt the learning rate during training based on the gradients. Examples include Adam, Adagrad, and RMSprop.

**Learning rate schedules:** These are predetermined schedules for decreasing the learning rate over time during training. Examples include step decay, exponential decay, and cosine annealing. In our approach, let's assume the hyper-parameter to be tuned is denoted by, and the function to be maximized is the test accuracy of the training data. Therefore, the fitness function is defined as the tuning shown in the following equations :

$$L = 1 - acc$$

$$grad = L_{new} - L_{old}$$

$$H_{new} = H_{old} + \alpha * grad$$

Where  $\alpha$  is constant term; it is search step; in this approach it is  $10^{-4}$ .

The tuned parameters in our case are:

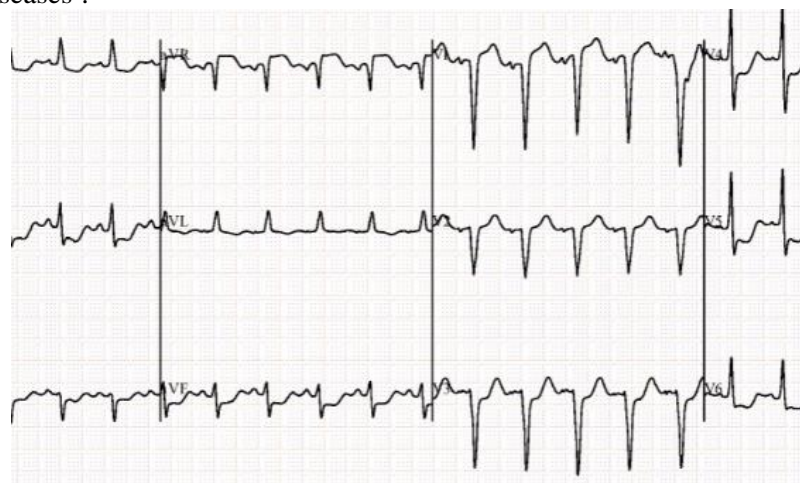
- Learning rate: was tuned using the above equations.
- Number of layers in transformer encoder: was tuned by a table of values.
- The number of patches which is applied the first of ViT: was tuned according the image size: suppose the image with size (H, W) so the number of patches was selected by supposing a kernel with size 5\*5 pixels:

$$patches = \frac{H * W}{5 * 5}$$

Scheduling the values as 1

#### 4. Implementations and Results Dataset

Dataset of Cardiac Patients created under the auspices of Ch. Pervaiz Elahi Institute of Cardiology Multan, Pakistan to help the scientific community for conducting the research for cardiovascular diseases<sup>1</sup>.



**Fig. 6:** Abnormal heartbeat from dataset

<sup>1</sup> <https://data.mendeley.com/datasets/gwbz3fsgp8/2>



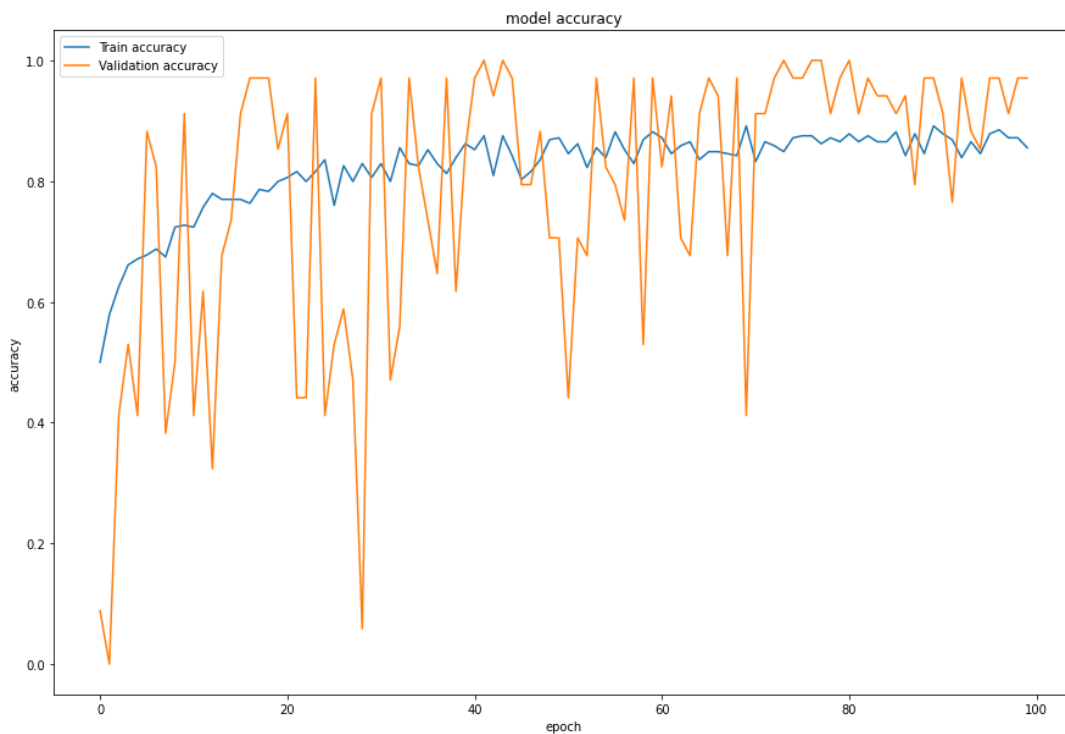
**Fig.7:** Normal heartbeat from dataset

**Implementation**

The implementation was on Google-Colab, the dataset was split into train set (338 images) and test set (78 images). Images augmentation was also used to increase the dataset size. The check point is being saved to be used later. The validation metric is the metric used by the model training.

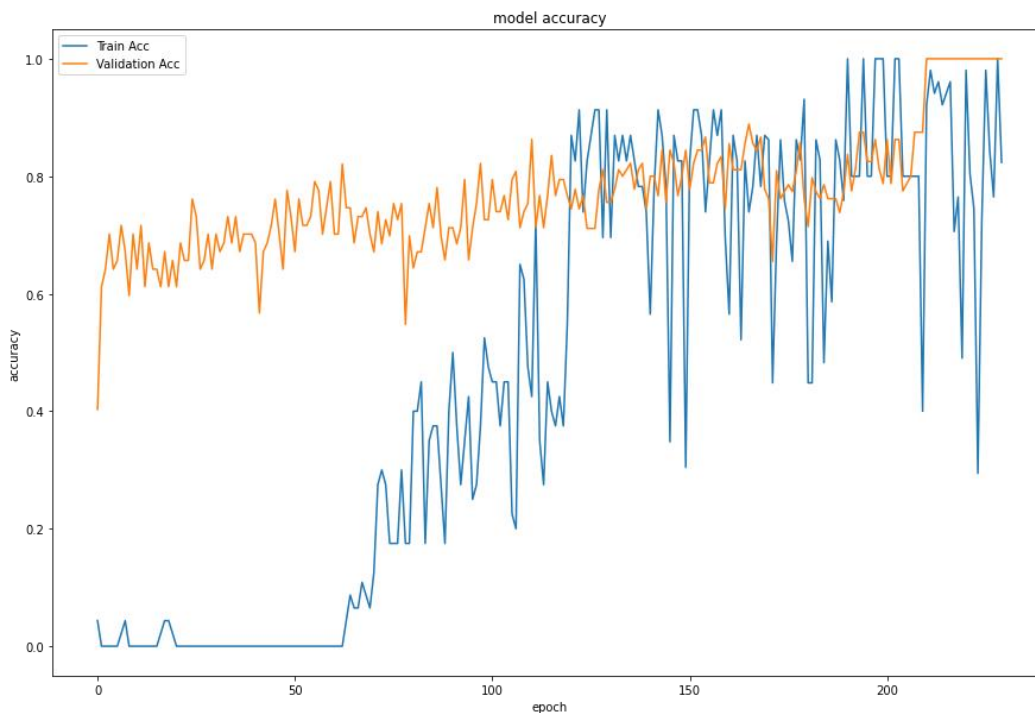
**Results and discussion**

The first ran of program with tuning of learning rate, just optimizing number of patches and number of layers in transformer encoder. With five layers in MLP and layers has 64, 128, 192, 256, 320 nodes (in transformer encoder). Each transformer layer has these five layers of MLP. **Fig.** shows the final results with tuning learning rate, **Fig.** shows the final results with tuning learning



**Fig. 8:** First results





**Fig. 9:** Final results

After many iterations (about 20000 iterations) the final accuracy of classification was 98.23%. The table below summarize all results (number of patches=NP, number of transformer layers in transformer encoder = NL), the accuracy when tuning Learning rate LR is reached only 97.43%. Then applying other tuning parameters for NL, NP and increasing the number of nodes in each MLP part in Deep-ViT model to reach final accuracy:

**Table 1:** learning rate tuning results

<b>Learning rate (NP=6, NL=6)</b>	<b>Accuracy</b>
<b>0.038</b>	92.21
<b>0.035</b>	92.45
<b>0.031</b>	93.01
<b>0.026</b>	94.21
<b>0.023</b>	94.31
<b>0.018</b>	94.85
<b>0.012</b>	95.15
<b>0.005</b>	96.12
<b>0.008</b>	97.18
<b>0.001</b>	97.25
<b>0.00015</b>	<b>97.43</b>
<b>0.00008</b>	97.28
<b>0.00003</b>	96.41
<b>0.00001</b>	96.03
<b>0.000008</b>	95.11

Comparison with state of art is not effective since the no papers worked on the same dataset in this approach. After applying tuning on NL and NP. And increasing number of nodes to 256, the results are shown in table below with comparison with state if arts result.

**Table 2:** Final results with comparison

<b>LR = 0.00015, NP=12, NL=9</b>	<b>Accuracy %</b>	<b>Dataset</b>
<b>Ours</b>	98.23	Heartbeats ECG
<b>Results from Error! Reference source not found.</b>	84.1	ImageNet1K dataset
<b>Results from Error! Reference source not found.</b>	85.8	Kinetics dataset <b>Error! Reference source not found.</b>
<b>Results from Error! Reference source not found.</b>	98.88	PhysioNet <b>Error! Reference source not found.</b>

The results from **Error! Reference source not found.** is higher than ours, but the data is different since authors had one image for only one heartbeat, in ours each image has a full ECG heartbeat as stated in **Fig. 3**.

## 5. Conclusion

In this research an enhanced Deep-ViT model was implemented and tested on heartbeats images of ECG for early diagnosis. The enhanced Deep-ViT has mainly contributed in layers in transformer encoder, layers in each MLP part and number of nodes in these layers and tuning parameters such as learning rate and number of patches. The final accuracy was 98.23% of detecting abnormal heartbeats.

## References

1. J. Kellett and S. Rasool, The prediction of the in-hospital mortality of acutely ill medical patients by electrocardiogram (ECG) dispersion mapping compared with established risk factors and predictive scores—A pilot study, *\*European Journal of Internal Medicine\**, vol. 22, no. 4, pp. 394-398, 2011. doi: 10.1016/j.ejim.2011.04.008.
2. S. I. Safie, J. J. Soraghan, and L. Petropoulakis, "Electrocardiogram (ECG) biometric authentication using pulse active ratio (PAR)," *\*IEEE Transactions on Information Forensics and Security\**, vol. 6, no. 4, pp. 1315-1322, 2011. doi: 10.1109/TIFS.2011.2168982.
3. M. M. Naseer, K. Ranasinghe, S. H. Khan, M. Hayat, F. Shahbaz Khan, and M. H. Yang, "Intriguing properties of vision transformers," *\*Advances in Neural Information Processing Systems\**, vol. 34, pp. 23296-23308, 2021. doi: 10.48550/arXiv.2106.00751.
4. K. Islam, "Recent advances in vision transformer: A survey and outlook of recent work," *\*arXiv preprint arXiv:2203.01536\**, 2022. doi: 10.48550/arXiv.2203.01536.
5. A. Parvaiz, M. A. Khalid, R. Zafar, H. Ameer, M. Ali, and M. M. Fraz, "Vision transformers in medical computer vision—A contemplative retrospection," *\*Engineering Applications of Artificial Intelligence\**, vol. 122, p. 106126, 2023. doi: 10.1016/j.engappai.2022.106126.
6. M. Maaz, A. Shaker, H. Cholakkal, S. Khan, S. W. Zamir, R. M. Anwer, and F. Shahbaz Khan, "Edgenext: efficiently amalgamated cnn-transformer architecture for mobile vision applications," in *\*Computer Vision—ECCV 2022 Workshops: Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part VII\**, Cham: Springer Nature Switzerland, 2023, pp. 3-20. doi: 10.1007/978-3-030-66823-5\_1.

7. M. J. Jassim, et al., "Greedy Learning of Deep Boltzmann Machine (GDBM)'s Variance and Search Algorithm for Efficient Image Retrieval," \*IEEE\* DOI [10.1109/ACCESS.2019.2948266](https://doi.org/10.1109/ACCESS.2019.2948266).
8. M. J. Jassim, et al., "An effective image retrieval based on optimized genetic algorithm utilized a novel SVM-based convolutional neural network classifier," 2019. [DOI org/10.1186/s13673-019-0191-8].
9. O. Dalmaz, M. Yurt, and T. Çukur, "ResViT: residual vision transformers for multimodal medical image synthesis," \*IEEE Transactions on Medical Imaging\*, vol. 41, no. 10, pp. 2598-2614, 2022. doi: 10.1109/TMI.2022.3182801.
10. M. S. Khalefa and Z. A. Abduljabar, "FINGERPRINT IMAGE ENHANCEMENT BY DEVELOP MEHTRE TECHNIQUE," 2011. DOI: 10.5121/acij.2011.2616
11. M. J. Jassim, Z. A. Hussien, M. S. Khalefa, Z. A. Abduljabbar, et al., "Fully automated model on breast cancer classification using deep learning classifiers," \*International Journal of Electrical and Computer Engineering\*, vol. 28, no. 1, pp. 183-191, 2022. doi: 10.11591/ijece.v28i1.183-191.
12. Y. Zhao, G. Wang, C. Tang, C. Luo, W. Zeng, and Z. J. Zha, "A battle of network structures: An empirical study of cnn, transformer, and mlp," \*arXiv preprint arXiv:2108.13002\*, 2021. doi: 10.48550/arXiv.2108.13002.
13. T. Liu, Y. Si, W. Yang, J. Huang, Y. Yu, G. Zhang, and R. Zhou, "Inter-Patient Congestive Heart Failure Detection Using ECG-Convolution-Vision Transformer Network," \*Sensors\*, vol. 22, no. 9, p. 3283, 2022. doi: 10.3390/s22093283.
14. S. Jamil and A. M. Roy, "An efficient and robust phonocardiography (pcg)-based valvular heart diseases (vhd) detection framework using vision transformer (vit)," \*Computers in Biology and Medicine\*, vol. 158, p. 106734, 2023. doi: 10.1016/j.combiomed.2023.106734.
15. D. Zhou, B. Kang, X. Jin, L. Yang, X. Lian, Z. Jiang, et al., "Deepvit: Towards deeper vision transformer," \*arXiv preprint arXiv:2103.11886\*, 2021. doi: 10.48550/arXiv.2103.11886.
16. T. Liu, Y. Si, W. Yang, J. Huang, Y. Yu, G. Zhang, and R. Zhou, "Inter-Patient Congestive Heart Failure Detection Using ECG-Convolution-Vision Transformer Network," \*Sensors\*, vol. 22, no. 9, p. 3283, 2022. doi: 10.3390/s22093283.
17. C. F. R. Chen, Q. Fan, and R. Panda, "Crossvit: Cross-attention multi-scale vision transformer for image classification," in \*Proceedings of the IEEE/CVF International Conference on Computer Vision\*, 2021, pp. 357-366. doi: 10.1109/ICCV.2021.00040.
18. A. Arnab, M. Dehghani, G. Heigold, C. Sun, M. Lučić, and C. Schmid, "Vivit: A video vision transformer," in \*Proceedings of the IEEE/CVF International Conference on Computer Vision\*, 2021, pp. 6836-6846. doi: 10.1109/ICCV.2021.00677.
19. K. Han, Y. Wang, H. Chen, X. Chen, J. Guo, Z. Liu, et al., "A survey on vision transformer," \*IEEE Transactions on Pattern Analysis and Machine Intelligence\*, vol. 45, no. 1, pp. 87-110, 2022. doi: 10.1109/TPAMI.2021.3050524.
20. R. Ranftl, A. Bochkovskiy, and V. Koltun, "Vision transformers for dense prediction," in \*Proceedings of the IEEE/CVF International Conference on Computer Vision\*, 2021, pp. 12179-12188. doi: 10.1109/ICCV.2021.01200.

## تصنيف دقيق لصور تخطيط القلب باستخدام Vision Transformer

فاطمة معك حنون<sup>1</sup> ، خولة حسين علي<sup>2\*</sup>

<sup>1</sup> قسم الحاسوب, كلية التربية للعلوم الصرفة, جامعة البصرة.

<sup>2</sup> قسم الحاسوب, كلية التربية للعلوم الصرفة, جامعة البصرة.

### الملخص

### معلومات البحث

يلعب تصنيف تخطيط القلب (ECG) دوراً مهماً في تشخيص وإدارة أمراض القلب والأوعية الدموية. أظهرت الأساليب العميقة القائمة على التعلم نتائج واعدة في تصنيف تخطيط القلب الآلي. ومع ذلك، فإن تعقيد إشارات ECG، بما في ذلك الاختلافات في التشكل، والمدة، والسعة، يشكل تحديات كبيرة لنماذج التعلم العميق الحالية. في هذا الصدد، أظهرت التطورات الحديثة في نماذج محاولات الرؤية أداءً رائعاً في معالجة الصور ومهام رؤية الكمبيوتر. في هذه الورقة، نقترح نهجاً قائماً على محول الرؤية العميقة لتصنيف ECG، والذي يجمع بين قوة الشبكات العصبية التلافيفية وآليات الاهتمام الذاتي. تم ضبط نموذجنا المقترح وتعزيزه من خلال أربع أجهزة برمائية مفرطة في النموذج المقترح، يمكنه اكتشاف داخلياً للميزات الرئيسية لصور ECG وتحقيق الأداء على مجموعات بيانات ECG القياسية. يمكن أن يساعد النموذج المقترح في الكشف المبكر وتشخيص أمراض القلب والأوعية الدموية، وبالتالي تحسين نتائج المريض كانت الدقة النهائية 98.23٪ في مجموعة البيانات <https://data.mendeley.com/datasets/gwbz3fsgp8/2>.

6 كانون الثاني 2024  
5 آذار 2024  
30 حزيران 2024

الاستلام  
القبول  
النشر

### الكلمات المفتاحية

Electrocardiogram,  
Classification, Vision  
transformer, Deep  
.Learning

**Citation:** Fatima M. Hanoon,  
Khawla H. Ali, J. Basrah Res.  
(Sci.) 50(1), 328 (2024).  
DOI:<https://doi.org/10.56714/bjrs.50.1.26>

\*Corresponding author email: alsalihfatimah4@gmail.com

