

# DEEP MODELING AND OPTIMIZATION OF MEDICAL IMAGE CLASSIFICATION

Yihang Wu<sup>1</sup>, Muhammad Owais<sup>1</sup>, Reem Kateb<sup>2</sup>, Ahmad Chaddad<sup>1,3,\*</sup>

<sup>1</sup>AIPM, School of Artificial Intelligence, Guilin University of Electronic Technology, China

<sup>2</sup>College of Computer Science and Engineering, Jeddah University, Jeddah, Saudi Arabia

<sup>3</sup>Laboratory for Imagery, Vision and Artificial Intelligence, École de Technologie Supérieure, Canada

Correspondence:ahmad8chaddad@gmail.com

## ABSTRACT

Deep models, such as convolutional neural networks (CNNs) and vision transformer (ViT), demonstrate remarkable performance in image classification. However, those deep models require large data to fine-tune, which is impractical in the medical domain due to the data privacy issue. Furthermore, despite the feasible performance of contrastive language image pre-training (CLIP) in the natural domain, the potential of CLIP has not been fully investigated in the medical field. To face these challenges, we considered three scenarios: 1) we introduce a novel CLIP variant using four CNNs and eight ViTs as image encoders for the classification of brain cancer and skin cancer, 2) we combine 12 deep models with two federated learning techniques to protect data privacy, and 3) we involve traditional machine learning (ML) methods to improve the generalization ability of those deep models in unseen domain data. The experimental results indicate that maxvit shows the highest averaged (AVG) test metrics (AVG = 87.03%) in HAM10000 dataset with multimodal learning, while convnext.l demonstrates remarkable test with an F1-score of 83.98% compared to swin\_b with 81.33% in FL model. Furthermore, the use of support vector machine (SVM) can improve the overall test metrics with AVG of  $\sim 2\%$  for swin transformer series in ISIC2018. Our codes are available at <https://github.com/AIPMLab/SkinCancerSimulation>.

**Index Terms**— Federated learning, Foundation models, Medical imaging.

## 1. INTRODUCTION

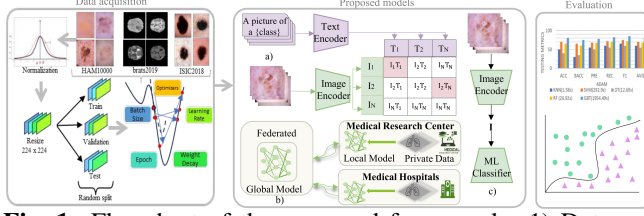
In recent years, deep learning models (DLMs) have significantly advanced medical imaging by using powerful architectures like convnext series, vision transformer (ViT) and maxvit for efficient deployment in realistic scenarios [1]. These models have shown a remarkable ability to learn complex visual patterns and have outperformed traditional machine learning techniques in both accuracy and scalability.

However, advanced DLMs require large amounts of source data to fine-tune, which is impractical in the field

of healthcare due to data privacy concerns [2, 3]. It can be challenging to collect and centralize the data required to effectively train algorithms since medical data is sensitive. Federated learning (FL), as a distributed learning framework that collaboratively trains a robust global model with multiple clients without sharing raw data, focuses on solving this challenge [2]. For example, federated averaging (FedAVG) can simply aggregate each client model with equal weights to produce the global model, demonstrating remarkable performance in many tasks [4]. Furthermore, with the development of CLIP, multi-modal (e.g., image with text) demonstrates competitive performance in natural image classification [5]. However, a comprehensive evaluation using recent CNN or ViT models in the classification of skin cancer is still lacking. In addition, the authors in [6] also claimed that natural foundation models such as CLIP show poor testing recall in the ISIC2019 data set.

Motivated by previous challenges, we introduce two FL approaches to solve the data privacy leakage issue. Furthermore, we propose a CLIP variant (i.e., replace the image encoder with CNN or ViT) to evaluate the generalization ability of those deep models in medical imaging. This approach allows the model to learn comprehensive representations from both images and text, improving its ability to accurately classify skin cancer. Specifically, we perform extensive experiments using three common benchmark data sets including two modalities with 12 deep models (four CNNs, eight ViTs) to provide a comprehensive analysis. We summarize the key contributions of this approach as follows:

1. We propose a novel CLIP-based approach that considers CNN and ViT architectures with CLIP text encoders for multimodal training, effectively integrating image and text data.
2. We introduce two FL approaches (FedAVG and FedProx) into skin cancer classification to solve data privacy leakage.
3. To improve the generalization ability of deep models, we combine traditional ML with deep models. Specifically, we use ML techniques as classifiers while using deep models as feature extractors.



**Fig. 1:** Flowchart of the proposed framework. 1) Data acquisition: Image data are preprocessed. 2) Proposed models: This involves three key components: multimodal learning, federated learning, and the combination of traditional ML and deep learning models. 3) Evaluation: we evaluate those models performance using standard classification metrics.

## 2. METHODOLOGY

Figure 1 illustrates the flowchart of the proposed model. Specifically, we propose a CLIP variant based on CNN and ViT to extract the image feature and then combine it with the text encoder in CLIP to perform training and inference. In addition, we introduce two FL techniques with four CNNs and eight ViTs as network backbones to solve data privacy issue. Furthermore, we combine ML approaches with deep models to enhance the generalizability of deep classifiers.

**Multimodal learning** The key idea of CLIP is to maximize the similarities between paired image and text (e.g., class label) features while minimize the similarities between unpaired image and text features. However, a comprehensive analysis of using CLIP based deep models in skin cancer classification is still lacking. Motivated by this challenge, we propose to use these deep models to extract image features while using transformer to obtain text features to perform CLIP. Let  $e_I(\cdot)$  be the image encoder while  $e_T(\cdot)$  be the text encoder. For a training example  $\mathbf{x}_j$ , we denote  $\mathbf{I}_j = e_I(\mathbf{x}_j) \in \mathbb{R}^D$  as the  $D$ -dimensional vector of image features. For text features, we use the standard prompt “a picture of a {class}” as input to the text encoder to obtain features  $\mathbf{T}_j = e_T(\mathbf{x}_j) \in \mathbb{R}^D$ . We calculate cosine similarity between the image and text features to measure the probability that  $\mathbf{x}_j$  belongs to a specific class  $c$  as follows.

$$p(Y=c|\mathbf{x}_j) = \frac{\exp(s_{j,c}/\tau)}{\sum_{c'=1}^K \exp(s_{j,c'}/\tau)}, \text{ with } s_{j,c} = \frac{\langle \mathbf{I}_j^*, \mathbf{T}_c \rangle}{\|\mathbf{I}_j^*\| \cdot \|\mathbf{T}_c\|} \quad (1)$$

Following [5], the contrastive loss is used in to train the image encoder. Furthermore, we freeze the text encoder, only optimizing the image encoder. Let  $B$  be the batch size, we compute the contrastive loss over batches of size  $B$ . Let  $\mathbf{S}$  be the  $B \times B$  matrix where  $s_{j,j'}$  is the cosine similarity between the image features  $\mathbf{I}_j^*$  and  $\mathbf{T}_{j'}$  as measured in Eq (1). We compute an image probability matrix  $\mathbf{P} = \text{softmax}(\mathbf{S}/\tau) \in [0, 1]^{B \times B}$  and a text probability matrix  $\mathbf{Q} = \text{softmax}(\mathbf{S}^T/\tau) \in [0, 1]^{B \times B}$ . The contrastive loss is then formulated as follows:

$$\mathcal{L}_{contr} = -\frac{1}{B} \sum_{j=1}^B \frac{1}{2} (\log p_{j,j} + \log q_{j,j}). \quad (2)$$

where  $p_{j,j}$  ( $q_{j,j}$ ) is the vector in matrix  $\mathbf{P}$  ( $\mathbf{Q}$ ).

**Federated learning** In this study, we use two common FL techniques, FedAVG [4] and FedProx [7]. FedAVG is a FL methodology that using averaging aggregation to obtain the global model. FedProx adds a small proximal term to measure the discrepancy between the local and global model to enhance the generalizability of local models. After each local epoch, the clients will send the local models parameters (i.e., weights)  $\omega_i^t$  to the global server, while the global server will aggregate those weights by the aggregation techniques defined as follows.

$$\omega_{glo}^t = \frac{1}{M} \sum_{i=1}^M \omega_i^t \quad (3)$$

where  $M$  is the number of clients.

Let  $f_\theta$  be the global model, our goal is to train a robust global model that can perform well on each client test data, i.e.,

$$\min_f \frac{1}{N} \sum_{i=1}^N \frac{1}{n_i^{test}} \sum_{j=1}^{n_i^{test}} \mathcal{L}(f_\theta(\mathbf{x}_{i,j}^{test}), y_{i,j}^{test}), \quad (4)$$

where  $\mathcal{L}$  is a given loss function,  $\mathbf{x}_{i,j}^{test}$  is the  $j$ -th test input from client  $i$ , and  $y_{i,j}^{test}$  is the corresponding label.

**Deep models with traditional machine learning** Deep models have remarkable ability of feature extraction, thereby using a simple linear layer can produce feasible performance in classification tasks. However, this may omit the merits of advanced traditional ML techniques such as Random Forest (RF). Inspired by [8], we use traditional ML techniques as the classifier while using deep models as the feature extractor. These classifier models take the features extracted by deep models without the last classifier layer as inputs and predict based on these image features.

**Table 1:** Summary of hyper-parameter settings in optimizer.

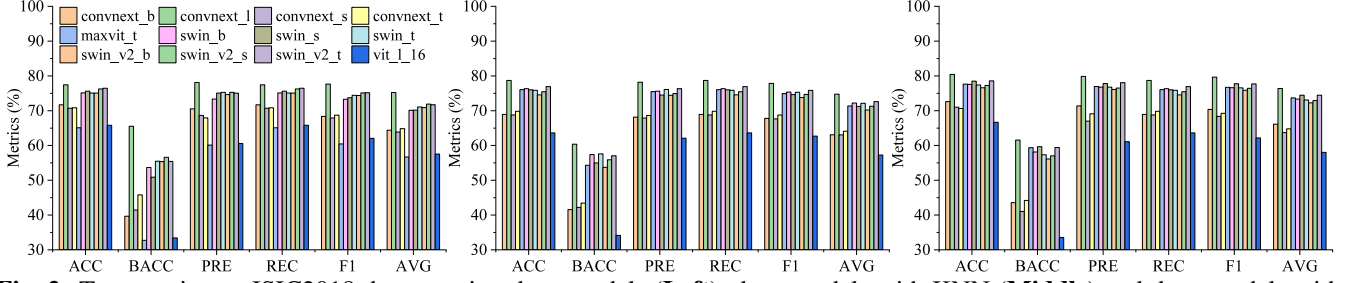
	LR	WD	Betas
SGD	0.01	0.0005	-
Adam	0.001	0.02	(0.9,0.98)
AdamW	0.001	0.02	(0.9,0.98)
Adagrad	0.001	0.0005	-
Adadelta	0.001	0.0005	-

## 3. EXPERIMENTS

### 3.1. Datasets

**Federated HAM10000 (FHAM).** We modified the HAM10000 dataset to build the FL dataset [9]. We split the original training set of HAM10000 into three clients with randomly selected samples, while the original testing set is used for global evaluation. In each client, the data are randomly partitioned into three non-overlapping parts, namely a training set (60%), a validation set (20%) and a testing set (20%).

**ISIC2018.** ISIC2018 is a medium scale skin cancer dataset ( $\sim 10000$  images) with seven classes [9, 10]. For this dataset,



**Fig. 2:** Test metrics on ISIC2018 dataset using deep models (Left), deep models with KNN (Middle) and deep models with SVM (Right).

we only considered its test set (1512 images) for the analysis of generalizability.

**BraTS2019.** We use a public kaggle dataset [11] with pre-divided data (7:1:2 for train, validation and test). It has two classes, namely high grade glioblastoma (HGG) and low grade glioma (LGG). The training set has 231 patients (178 HGG and 53 LGG), the validation set contains 32 patients (25 HGG and 7 LGG), while the test set holds 68 patients (52 HGG and 16 LGG).

### 3.2. Tasks

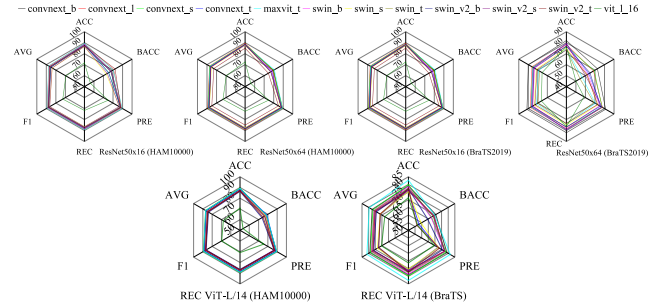
1. Multimodal learning. HAM10000 and BraTS2019 datasets are used to perform multimodal learning using 12 deep models (four CNNs and eight ViTs) with CLIP text encoders pretrained by three image encoders (i.e., ViT-L/14, ResNet50x16 and ResNet50x64). The Adagrad optimizer is considered for optimization.
2. Federated skin cancer classification. We use FHAM dataset as an example to evaluate the usefulness of 12 deep models using two FL techniques (FedAVG and FedProx) with five optimizers.
3. Generalization analysis. The ISIC2018 test set is used to demonstrate the generalizability of four CNNs and eight ViTs using k-nearest neighbours (KNN) and support vector machine (SVM) [12]. Note that the deep models are pretrained on HAM10000 training set using SGD optimizer (no overlap with ISIC2018 test set).

### 3.3. Implementation details

We used convnext\_b, convnext\_l, convnext\_s, convnext\_t, maxvit\_t, swin\_b, swin\_s, swin\_t, swin\_v2\_b, swin\_v2\_s, swin\_v2\_t and vit\_l\_16 as the deep network backbones for classification tasks [13, 14]. We choose SGD, Adam, AdamW, Adagrad, and Adadelat optimizers for simulations [15]. Table 1 reports a detailed hyper-parameter settings for those optimizers. The training epoch is set to 50. For **Task 1**, the batch-size is set to 16, while for **Task 2**, the batchsize is set to 32 for HAM10000 and 16 for BraTS2019. The batchsize of **Task 3** is set to 32. The experiment environment is based on the Windows 11 operating system and features an Intel 13900KF

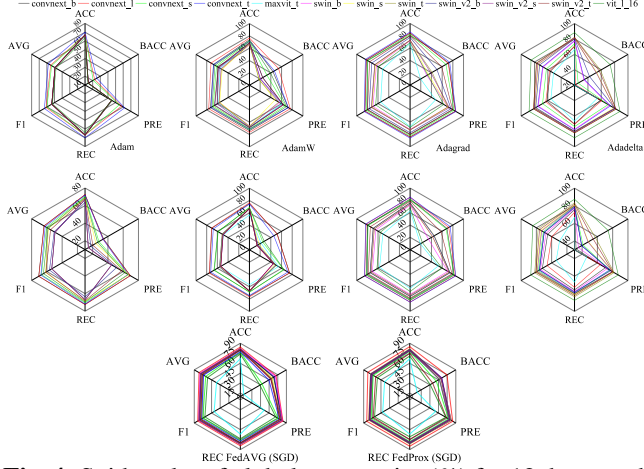
CPU with 128 GB of RAM and an RTX 4090 GPU. We use Pytorch 1.13.1 and Python 3.8. For classification metrics, this study considered accuracy (ACC), balanced accuracy (BACC), weighted-precision (PRE), weighted-recall (REC), weighted-F1 score and AVG (i.e., unweighted mean value of ACC, BACC, PRE, REC and F1).

**Task 1** Figure 3 shows the test metrics in the HAM10000 and BraTS datasets. As illustrated, for Adagrad optimizer, maxvit\_t demonstrates the best overall testing metrics (87.03% with ViT-L/14, 87.58% with ResNet50x16) while convnext-large shows the highest overall testing metrics (87.36% with ResNet50x64). This suggests that maxvit is more suitable to perform multimodal in skin cancer classification. Furthermore, unlike HAM10000, maxvit\_t fits well with the text encoder pretrained by ViT-L/14 (e.g., 79.6% AVG) while performs poorly with the text encoder pretrained by ResNet50x64 (69.41% AVG).



**Fig. 3:** Spider-plot of multimodal test metrics (%) for 12 deep models in HAM10000 and BraTS2019 using text encoder pretrained by three CLIP image encoders (ViT-L/14, ResNet50x16 and ResNet50x64).

**Task 2** Figure 4 shows the test metrics in HAM10000. As illustrated, convnext\_l demonstrates remarkable overall testing metrics (81.66% using FedAVG and 82.72 using FedProx) with the SGD optimizer. Similarly, deep models with Adagrad optimizer indicate better performance compared to Adam, AdamW and Adadelat-based models. These findings suggest that the use of CLIP can achieve feasible performance, however, it introduces large communication costs in the FL context as suggested in [16]. In future work, we will explore efficient training techniques such as Adapter [3].



**Fig. 4:** Spider-plot of global test metrics (%) for 12 deep models in HAM10000 using FedAVG (First and third row) and FedProx (Second and third Row).

**Task 3** Figure 2 shows the metrics in the ISIC2018 test set. We observe that 1) introducing ML techniques such as KNN and SVM can improve the overall testing metrics of maxvit.t  $\sim 15\%$ . 2) Similarly, the use of SVM indicates the best overall results compared to the original deep models. Those findings suggest that using ML techniques can further improve the performance of deep models in unseen domain.

#### 4. CONCLUSION

This study proposed three models covering multimodal, FL and traditional ML with deep models in medical image classification tasks. The findings suggest that maxvit.t shows potential for multimodal, convnext.l indicates remarkable overall test metrics using FedAVG and FedProx with SGD optimizer, while introducing SVM and KNN can improve the overall performance of maxvit.t, vit.l16, convnext.b, convnext.l and swin transformer series in unseen domain. In future work, we will introduce domain adaptation [17] to minimize the data distribution shifts among different datasets to further improve deep models performance.

#### 5. COMPLIANCE WITH ETHICAL STANDARDS

This is a numerical simulation study for which no ethical approval was required.

#### 6. ACKNOWLEDGEMENTS

This research was funded by the National Natural Science Foundation of China #82260360, the Guilin Innovation Platform and Talent Program #2022C264164, and the Guangxi Science and Technology Base and Talent Project (#2022AC18004 and #2022AC21040).

#### 7. REFERENCES

- [1] I. Pacal, "Maxcervix: A novel lightweight vision transformer-based approach for precise cervical cancer detection," *Knowledge-Based Systems*, vol. 289, p. 111482, 2024.
- [2] A. Chaddad, Y. Wu, and C. Desrosiers, "Federated learning for health-care applications," *IEEE Internet of Things Journal*, 2023.
- [3] Y. Wu, C. Desrosiers, and A. Chaddad, "Facmic: Federated adaptive clip model for medical image classification," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 531–541, Springer, 2024.
- [4] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*, pp. 1273–1282, PMLR, 2017.
- [5] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, *et al.*, "Learning transferable visual models from natural language supervision," in *International conference on machine learning*, pp. 8748–8763, PMLR, 2021.
- [6] J. P. Huix, A. R. Ganeshan, J. F. Haslum, M. Söderberg, C. Matsoukas, and K. Smith, "Are natural domain foundation models useful for medical image classification?," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 7634–7643, 2024.
- [7] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated optimization in heterogeneous networks," *Proceedings of Machine learning and systems*, vol. 2, pp. 429–450, 2020.
- [8] A. Almomani, K. Nahar, M. Alauthman, M. A. Al-Betar, Q. Yaseen, and B. B. Gupta, "Image cyberbullying detection and recognition using transfer deep machine learning," *International Journal of Cognitive Computing in Engineering*, vol. 5, pp. 14–26, 2024.
- [9] P. Tschandl, C. Rosendahl, and H. Kittler, "The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions," *Scientific data*, vol. 5, no. 1, pp. 1–9, 2018.
- [10] N. Codella, V. Rotemberg, P. Tschandl, M. E. Celebi, S. Dusza, D. Gutman, B. Helba, A. Kalloo, K. Liopyris, M. Marchetti, *et al.*, "Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic)," *arXiv preprint arXiv:1902.03368*, 2019.
- [11] M. 2019, "Brain tumor classification dataset..," <https://www.kaggle.com/datasets/shirtgm/brats2019-classification-divided-by-patients/data>.
- [12] M. Bansal, A. Goyal, and A. Choudhary, "A comparative analysis of k-nearest neighbor, genetic, support vector machine, decision tree, and long short term memory algorithms in machine learning," *Decision Analytics Journal*, vol. 3, p. 100071, 2022.
- [13] S. Woo, S. Debnath, R. Hu, X. Chen, Z. Liu, I. S. Kweon, and S. Xie, "Convnext v2: Co-designing and scaling convnets with masked autoencoders," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16133–16142, 2023.
- [14] S. Bangalore Vijayakumar, K. T. Chitty-Venkata, K. Arya, and A. K. Somani, "Convion benchmark: A contemporary framework to benchmark cnn and vit models," *AI*, vol. 5, no. 3, pp. 1132–1171, 2024.
- [15] R. Abdulkadirov, P. Lyakhov, and N. Nagornov, "Survey of optimization algorithms in modern neural networks," *Mathematics*, vol. 11, no. 11, p. 2466, 2023.
- [16] J. Hu, D. Wang, Z. Wang, X. Pang, H. Xu, J. Ren, and K. Ren, "Federated large language model: Solutions, challenges and future directions," *IEEE Wireless Communications*, 2024.
- [17] A. Chaddad and Y. Wu, "Enhancing classification tasks through domain adaptation strategies," in *2023 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pp. 1832–1835, IEEE, 2023.