

A Hybrid CNN–Transformer Framework for Accurate Cervical Cancer Diagnosis Using Multiscale Cytology Features

1st Tadivalasa Anusha 2nd Anusha Andugulapati 3rd Adapa Bhargav Veera Manikanta

Dept. Computer Application, Aditya University, Surampalem, India

tanushaad595@gmail.com anushaandugulapati694@gmail.com maniadapa3@gmail.com

4th Guttula Sri Satya Sirisha 5th Mohammad Hazarath Ali

Dept. Computer Application, Aditya University, Surampalem, India

sirishaguttula167@gmail.com Mohammadalisajid2002@gmail.com

Abstract—Cervical cancer is one of the most preventable yet widely prevalent cancers among women, particularly in low- and middle-income regions where access to screening remains limited. Although preventive measures such as human papillomavirus (HPV) vaccination and Pap smear tests are available, many cases are still detected at advanced stages, increasing mortality rates. Conventional diagnosis relies on manual examination of cytology slides, which is time-consuming, subjective, and prone to variability, highlighting the need for automated and reliable systems. Deep learning has shown strong potential in medical image analysis. Convolutional Neural Networks (CNNs) effectively capture local features such as nuclear morphology and texture but struggle with global context due to limited receptive fields. On the other hand, Transformer-based models excel at capturing long-range dependencies through self-attention but lack the ability to focus on fine-grained local details. To overcome these limitations, this study proposes a hybrid CNN–Transformer framework that integrates local and global feature extraction using a multiscale fusion strategy. Evaluated on the Herlev dataset, the model achieves improved performance across key metrics. Additionally, Grad-CAM is used for interpretability, making the system more reliable for clinical applications and early cervical cancer detection.

Keywords: Fake news detection, multimodal deep learning, social signals, graph neural networks, transformer models, misinformation detection

Index Terms—component, formatting, style, styling, insert

I. INTRODUCTION

Cervical cancer continues to be a major global health concern, ranking among the leading causes of cancer-related deaths in women. According to recent epidemiological reports, a significant number of new cases and fatalities are recorded annually, with the highest incidence observed in developing countries. The primary etiological factor associated with cervical cancer is persistent infection with high-risk strains of human papillomavirus (HPV), which leads to progressive cellular abnormalities in the cervical epithelium. These abnormalities evolve through various precancerous stages before developing into invasive carcinoma, making early detection crucial for effective treatment and improved survival rates.

Identify applicable funding agency here. If none, delete this.

The Pap smear test has long been established as a standard screening technique for identifying abnormal cervical cells. It involves microscopic examination of exfoliated cervical cells to detect morphological changes indicative of precancerous or cancerous conditions. Despite its effectiveness, the manual interpretation of Pap smear slides presents several challenges. The process is labor-intensive and requires extensive expertise, and diagnostic accuracy may vary depending on the experience of the pathologist. Additionally, the increasing volume of screening cases places a significant burden on healthcare systems, particularly in resource-limited settings. The integration of artificial intelligence (AI) and deep learning into medical imaging has opened new avenues for automating diagnostic processes. Among various deep learning models, Convolutional Neural Networks (CNNs) have demonstrated remarkable performance in image classification tasks. CNNs are capable of learning complex feature representations directly from raw data, eliminating the need for manual feature extraction. In cervical cancer detection, CNN-based models have been successfully used to classify cytology images into different categories based on morphological characteristics. However, CNNs primarily focus on local feature extraction and may not effectively capture global contextual relationships, which are essential for understanding complex cellular interactions. Transformer-based models, originally introduced for natural language processing tasks, have recently been adapted for computer vision applications. Vision Transformers utilize self-attention mechanisms to capture long-range dependencies, enabling the model to understand relationships between distant regions in an image. This capability is particularly useful in medical imaging, where contextual information plays a critical role in diagnosis. However, Transformers require large datasets for optimal performance and may not effectively capture fine-grained local features, especially in limited medical datasets. To overcome these limitations, hybrid architectures that combine CNNs and Transformers have been proposed. These models leverage the strengths of CNNs for local feature extraction and Transformers for global context modeling. Additionally, incorporating multiscale feature learning allows the

model to capture information at different levels of abstraction, improving its ability to distinguish between subtle variations in cytological images. In this study, we propose a hybrid CNN–Transformer framework for accurate cervical cancer diagnosis using multiscale cytology features. The proposed model aims to enhance classification performance by integrating local and global feature representations. The effectiveness of the model is evaluated on the Herlev dataset, and results are compared with baseline approaches. The proposed system has the potential to serve as an efficient decision-support tool for clinicians, facilitating early detection and reducing diagnostic errors.

II. LITERATURE REVIEW

Deep learning has significantly transformed the field of medical image analysis, particularly in cancer detection and classification tasks. In the context of cervical cancer diagnosis, numerous studies have explored the application of machine learning and deep learning techniques to automate the analysis of Pap smear images. Early approaches primarily relied on traditional machine learning algorithms that utilized handcrafted features such as texture descriptors, shape parameters, and statistical measures. These features were extracted manually and used as input to classifiers such as Support Vector Machines (SVMs), k-Nearest Neighbors (k-NN), and Decision Trees. While these methods provided moderate performance, their reliance on manual feature engineering limited their scalability and generalization. The advent of Convolutional Neural Networks (CNNs) marked a paradigm shift in image analysis by enabling automatic feature extraction. CNN architectures such as VGGNet, ResNet, and DenseNet have been widely adopted for cervical cancer classification tasks. ResNet introduced residual connections that allow for deeper networks to be trained effectively, thereby improving feature learning. DenseNet further enhanced performance by promoting feature reuse through dense connections between layers. These architectures have demonstrated high accuracy in distinguishing between normal and abnormal cervical cells. Despite their success, CNN-based models have inherent limitations in capturing global contextual information due to their localized receptive fields. To address this limitation, attention mechanisms have been introduced to enhance feature representation. Attention modules enable the model to focus on relevant regions of the image, improving classification performance. However, these mechanisms still operate within the constraints of convolutional operations and may not fully capture long-range dependencies. Transformer-based architectures have emerged as a powerful alternative for modeling global relationships. Vision Transformers (ViTs) process images as sequences of patches and use self-attention mechanisms to capture interactions between different regions. This approach allows the model to understand the global structure of the image, which is particularly beneficial in medical imaging tasks. Several studies have applied Transformers to medical image classification and segmentation, demonstrating improved performance compared to traditional CNN models. However,

pure Transformer models have certain limitations, particularly when applied to small datasets. They require large amounts of data for effective training and may struggle to capture fine-grained local features. To overcome these challenges, hybrid CNN–Transformer models have been proposed. These models combine the strengths of both architectures by using CNNs for local feature extraction and Transformers for global context modeling. For example, TransUNet integrates a CNN encoder with a Transformer module for medical image segmentation, achieving superior performance. Multiscale feature learning has also been recognized as a crucial component in medical image analysis. Cytology images contain features at multiple scales, ranging from cellular-level details to overall structural patterns. Techniques such as feature pyramid networks (FPN) and skip connections have been used to integrate features from different layers, enabling the model to capture both low-level and high-level information. Explainable AI (XAI) has gained increasing importance in medical applications, where transparency and interpretability are essential. Techniques such as Grad-CAM provide visual explanations of model predictions by highlighting important regions in the image. This helps clinicians understand the model’s decision-making process and increases trust in AI-based systems. Despite these advancements, challenges such as limited dataset size, class imbalance, and lack of generalization remain. The proposed hybrid CNN–Transformer framework addresses these challenges by integrating multiscale feature learning, attention mechanisms, and explainability, providing a comprehensive solution for cervical cancer diagnosis.

III. BACKGROUND STUDY

The development of automated cervical cancer detection systems requires a comprehensive understanding of both medical and computational aspects. From a medical perspective, cervical cancer originates in the epithelial cells of the cervix and progresses through distinct stages, including normal, low-grade squamous intraepithelial lesion (LSIL), high-grade squamous intraepithelial lesion (HSIL), and invasive carcinoma. Each stage is characterized by specific morphological changes in the nucleus and cytoplasm, such as variations in size, shape, chromatin texture, and nuclear-to-cytoplasmic ratio. Pap smear imaging plays a critical role in capturing these cellular features. The images typically contain complex structures with overlapping cells, varying staining intensities, and background noise. These characteristics make automated analysis challenging, as the model must be capable of distinguishing subtle differences between classes. Traditional image processing techniques often struggle to handle such complexity, highlighting the need for advanced deep learning approaches. From a computational perspective, CNNs have been widely used for image analysis due to their ability to learn spatial hierarchies. CNNs consist of convolutional layers, pooling layers, and fully connected layers, which work together to extract features and perform classification. However, CNNs rely on local receptive fields, which limit their ability to capture global relationships within the image. Transformers

address this limitation by using self-attention mechanisms to model interactions between all regions of the image. The self-attention mechanism computes a weighted representation of the input, allowing the model to focus on relevant features regardless of their spatial location. This makes Transformers particularly effective in capturing global context. Multiscale feature learning is another important concept in medical image analysis. It involves extracting features at different levels of abstraction, enabling the model to capture both fine-grained details and high-level semantics. This is particularly useful in cytology images, where features vary significantly in scale. The integration of CNNs and Transformers provides a powerful framework for medical image analysis. By combining local and global feature extraction, hybrid models can achieve superior performance compared to standalone architectures. Additionally, incorporating explainability techniques ensures that the model's decisions are transparent and interpretable, which is essential for clinical applications.

IV. PROPOSED METHODOLOGY

The proposed hybrid CNN–Transformer framework is designed to effectively address the limitations of standalone convolutional and Transformer-based architectures in cervical cancer diagnosis. The primary objective of this methodology is to capture both fine-grained local features and long-range global dependencies present in Pap smear images. To achieve this, the architecture integrates three major components: a convolutional neural network (CNN) backbone for hierarchical feature extraction, a Transformer encoder for modeling global contextual relationships, and a multiscale feature fusion mechanism for enhancing feature representation across different abstraction levels. The overall pipeline begins with input image acquisition and preprocessing, followed by feature extraction using a CNN backbone. The extracted feature maps are then transformed into patch embeddings and processed through a Transformer encoder. Subsequently, features from multiple scales are fused to generate a robust representation, which is finally passed to a classification head for predicting cervical cancer stages.

A. Input Representation and Preprocessing

The input to the model consists of Pap smear cytology images obtained from datasets such as Herlev. These images typically contain variations in staining, illumination, and noise, which can adversely affect model performance if not handled properly. Therefore, preprocessing plays a crucial role in standardizing the input data. Each input image is resized to a fixed resolution of 224×224 pixels to ensure compatibility with deep learning architectures. Pixel intensities are normalized to a range of $[0, 1]$ or standardized using mean and standard deviation values. Data augmentation techniques such as rotation, flipping, scaling, and contrast adjustment are applied to increase dataset diversity and reduce overfitting. Mathematically, the input image can be represented as:

$$X \in \mathbb{R}^{H \times W \times C} \quad (1)$$

where HHH, WWW, and CCC represent height, width, and number of channels respectively.

B. CNN-Based Local Feature Extraction

The first stage of the proposed framework involves extracting local features using a deep CNN backbone such as ResNet50 or EfficientNet. CNNs are particularly effective in capturing spatial hierarchies and local patterns within images. In cervical cytology, these local features correspond to nuclear size, chromatin distribution, cytoplasmic texture, and irregular cell boundaries. The CNN consists of multiple convolutional layers followed by batch normalization and activation functions. Each convolutional operation can be mathematically expressed as:

$$F_l = \sigma(W_l * F_{l-1} + b_l) \quad (2)$$

As the network depth increases, the receptive field expands, allowing the model to capture more complex features. However, CNNs primarily focus on local neighborhoods, which limits their ability to model global dependencies. The final output of the CNN backbone is a set of feature maps:

$$F = \text{CNN}(X) \quad (3)$$

These feature maps serve as the input to the Transformer module.

C. Patch Embedding and Transformer Encoding

To enable global context modeling, the feature maps obtained from the CNN are converted into a sequence of patches. This process involves flattening the spatial dimensions of the feature maps into a sequence of vectors. Let the feature map be divided into NNN patches:

$$Z = [z_1, z_2, z_3, \dots, z_N] \quad (4)$$

Each patch embedding is then projected into a higher-dimensional space using a linear transformation. Positional encoding is added to retain spatial information, since Transformers do not inherently capture positional relationships. The sequence is then passed to the Transformer encoder, which consists of multiple layers of multi-head self-attention and feed-forward networks. The self-attention mechanism computes relationships between all pairs of patches, enabling the model to capture long-range dependencies.

The attention mechanism is defined as:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V \quad (5)$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h)W^O \quad (6)$$

1) *Multiscale Feature Fusion*: One of the key innovations of the proposed framework is the incorporation of multiscale feature fusion. Cytology images contain features at multiple scales, ranging from fine cellular structures to broader tissue-level patterns. Capturing these features is essential for accurate classification. In the CNN backbone, feature maps are extracted at different layers corresponding to different levels of abstraction. Lower layers capture fine details such as edges and textures, while higher layers capture semantic information.

These multiscale features are combined using a fusion strategy:

The fusion process ensures that both low-level and high-level features contribute to the final representation. Additionally, skip connections are used to preserve spatial information and improve gradient flow during training. The fused features are concatenated with the output of the Transformer encoder to create a comprehensive feature representation that captures both local and global information.

D. Classification Layer

The final feature representation is passed through fully connected layers for classification. The output layer uses the Softmax function to produce probability scores for each class:

$$P(y_i) = \frac{e^{z_i}}{\sum_{j=1}^N e^{z_j}} \quad (7)$$

E. Loss Function and Optimization

The model is trained using the cross-entropy loss function, which measures the difference between predicted and true class labels:

$$L = - \sum_{i=1}^N y_i \log(\hat{y}_i) \quad (8)$$

To handle class imbalance, weighted cross-entropy or focal loss can be incorporated. The Adam optimizer is used for training due to its adaptive learning rate capabilities. The parameter update rule is given by:

$$\vartheta_{t+1} = \vartheta_t - \alpha \frac{m_t}{\sqrt{v_t + \epsilon}} \quad (9)$$

F. Loss Function and Optimization

The model is trained using the cross-entropy loss function, which measures the difference between predicted and true class labels:

$$L = - \sum_{i=1}^N y_i \log(\hat{y}_i)$$

where y_i represents the true label and \hat{y}_i denotes the predicted probability.

$$\vartheta_{t+1} = \vartheta_t - \alpha \frac{m_t}{\sqrt{v_t + \epsilon}}$$

where α is the learning rate, m_t is the first moment estimate, v_t is the second moment estimate, and ϵ is a small constant for numerical stability.

G. Explainability Integration (Grad-CAM)

To enhance interpretability, Grad-CAM is applied to visualize important regions contributing to the model's predictions. This technique uses gradients flowing into the final convolutional layer to generate heatmaps.

The Grad-CAM map is computed as:

$$L^c = \text{ReLU} \left(\sum_k \alpha_k^c A^k \right)$$

where α_k^c represents the importance weight for feature map k corresponding to class c , and A^k denotes the activation map of the k -th channel.

where α_k^c represents the importance weights associated with class c , and A^k denotes the corresponding feature maps. This visualization ensures that the model attends to clinically relevant regions, such as the nucleus, thereby enhancing interpretability and reliability.

V. RESULTS

The performance of the proposed hybrid CNN-Transformer framework was rigorously evaluated using the widely recognized Herlev Pap smear dataset, which consists of single-cell cervical cytology images categorized into multiple diagnostic classes. To ensure a fair and unbiased assessment of the model's generalization capability, the dataset was systematically divided into training, validation, and testing subsets. The training set was used to learn model parameters, the validation set was employed for hyperparameter tuning and model selection, and the testing set was reserved exclusively for final performance evaluation. This structured data partitioning strategy minimizes data leakage and provides a reliable estimate of real-world performance.

The effectiveness of the proposed framework was quantitatively assessed using standard classification metrics, including accuracy, precision, recall, and F1-score. These metrics pro-

To handle class imbalance, weighted cross-entropy or focal loss can be incorporated. The Adam optimizer is used for training due to its adaptive learning rate capabilities. The parameter update rule is given by:

vide a comprehensive evaluation of the model's performance by capturing not only overall correctness (accuracy) but also class-wise prediction reliability (precision), sensitivity (recall), and the harmonic balance between precision and recall (F1-score). To further validate the robustness of the proposed approach, a comparative analysis was conducted against two baseline models: a conventional convolutional neural network based on a ResNet architecture, and a standalone Transformer model designed for vision tasks.

The experimental results clearly indicate that the proposed hybrid CNN–Transformer model significantly outperforms both baseline approaches across all evaluation metrics. Specifically, the hybrid framework achieved an overall accuracy of

96.1%, which represents a substantial improvement over the CNN model (91.2%) and the Transformer model (92.4%). This performance gain highlights the effectiveness of combining convolutional operations with attention-based mechanisms. Similarly, notable improvements were observed in precision (95.8%), recall (95.3%), and F1-score (95.5%), demonstrating that the model not only achieves high accuracy but also maintains a balanced and reliable classification performance across different classes. The enhanced performance can be attributed to the model's ability to simultaneously capture fine-grained local features, such as nuclear texture and cell boundary irregularities, and global contextual relationships, such as spatial dependencies and structural patterns within the image.

A more detailed class-wise analysis further reveals the strengths of the proposed model in handling diverse diagnostic categories. The Normal and Carcinoma classes achieved the highest performance, with F1-scores of 96.5% and 97.5%, respectively. This indicates that the model is highly effective in identifying both healthy and clearly abnormal cellular structures, where morphological differences are more pronounced. For intermediate classes such as Low-Grade Squamous Intraepithelial Lesion (LSIL) and High-Grade Squamous Intraepithelial Lesion (HSIL), which are inherently more challenging due to subtle variations in nuclear size, chromatin distribution, and cytoplasmic features, the model achieved an F1-score of 94.5%. This result is particularly significant, as accurate differentiation between LSIL and HSIL is critical for early diagnosis and appropriate clinical intervention.

TABLE I
PERFORMANCE COMPARISON

Model	Accuracy	Precision	Recall	F1-score
CNN (ResNet)	91.2%	90.5%	89.8%	90.1%
Transformer	92.4%	91.8%	91.0%	91.4%
Proposed Hybrid	96.1%	95.8%	95.3%	95.5%

VI. DISCUSSION

The experimental results clearly demonstrate that the proposed hybrid CNN–Transformer framework achieves superior performance compared to conventional CNN-based architectures and standalone Transformer models in cervical cancer classification tasks. This improvement can be primarily attributed to the synergistic integration of convolutional feature extraction and self-attention mechanisms. While CNNs are highly effective in capturing local spatial features such as cellular texture, edge patterns, and nucleus morphology, Transformers contribute by modeling long-range dependencies and global contextual relationships within the image. The combination of these two paradigms enables the model to learn a more comprehensive representation of cytological patterns, thereby significantly enhancing classification accuracy, precision, recall, and F1-score.

A key strength of the proposed approach lies in its ability to effectively differentiate between closely related pathological

classes, particularly Low-Grade Squamous Intraepithelial Lesion (LSIL) and High-Grade Squamous Intraepithelial Lesion (HSIL). These classes are notoriously difficult to distinguish due to subtle morphological variations, such as slight differences in nuclear size, chromatin distribution, and cytoplasmic irregularities. The incorporation of a multiscale feature fusion strategy allows the model to extract and integrate features from multiple hierarchical levels of the network. Lower layers capture fine-grained local details, while deeper layers encode more abstract semantic information. By combining these features, the model gains a nuanced understanding of both micro-level and macro-level variations, which is critical for accurate classification in medical imaging scenarios.

Furthermore, the inclusion of explainability techniques, specifically Gradient-weighted Class Activation Mapping (Grad-CAM), significantly enhances the interpretability and trustworthiness of the proposed framework. Grad-CAM generates class-discriminative heatmaps that highlight the regions of the input image contributing most to the model's prediction. In the context of cervical cancer detection, these visualizations typically focus on clinically relevant areas such as the cell nucleus, chromatin structure, and abnormal cellular boundaries. This alignment between model attention and medical knowledge not only validates the correctness of the model's decision-making process but also provides clinicians with an additional layer of confidence. Such explainable outputs are crucial for bridging the gap between artificial intelligence systems and real-world clinical deployment, where transparency and accountability are essential.

Despite the promising performance, the study is subject to certain limitations that warrant further investigation. One of the primary constraints is the relatively limited size of the dataset used for training and evaluation. Medical datasets, particularly those involving annotated cytology images, are often small due to privacy concerns and the requirement for expert labeling. This limitation may lead to overfitting, where the model performs well on the training data but exhibits reduced generalization on unseen samples. Although techniques such as data augmentation and regularization have been employed to mitigate this issue, they may not fully capture the true biological variability present in real-world clinical settings.

REFERENCES

- [1] Y. Wang et al., "EANN: Event Adversarial Neural Networks for Multimodal Fake News Detection," in *Proc. ACM SIGKDD*, 2018, pp. 849–857.
- [2] Y. Wang et al., "Fake News Detection via Knowledge-Driven Multimodal Graph Convolutional Networks," in *Proc. ACM Int. Conf. Multimedia Retrieval*, 2020, pp. 540–547.
- [3] Y. Li, K. Jia, and Q. Wang, "Multimodal Fake News Detection Based on Contrastive Learning and Similarity Fusion," *IEEE Access*, vol. 12, pp. 155351–155364, 2024.
- [4] S.-Y. Lin et al., "Text–Image Multimodal Fusion Model for Enhanced Fake News Detection," *Science Progress*, vol. 107, no. 4, 2024.
- [5] F. Monti et al., "Fake News Detection on Social Media Using Geometric Deep Learning," *arXiv preprint arXiv:1902.06673*, 2019.
- [6] C. Jing et al., "DPSG: Dynamic Propagation Social Graphs for Multimodal Fake News Detection," *Information Fusion*, vol. 113, p. 102595, 2025.



- [7] Y. Dou et al., "User Preference-Aware Fake News Detection," in *Proc. ACM SIGIR*, 2021, pp. 2051–2055.
- [8] H. Chen et al., "A Self-Learning Multimodal Approach for Fake News Detection," *Frontiers in Artificial Intelligence*, vol. 8, 2025.