# A Vision Transformer-Based Approach For Ovarian Cancer Detection And Classification

[1]Apeksha Babu A, [2]Prof. Shruthi B Gowda

[1]Student, [2]Assistant Professor

[1]Department Of Computer Science,

[1]Bangalore Institute Of Technology, Bengaluru, India

***Abstract:*** Deep learning has transformed medical imaging, particularly cancer detection. This paper introduces a Vision Transformer (ViT)-based approach for ovarian cancer classification. Unlike CNN-based models such as ResNet-50, ViTs utilize self-attention mechanisms to enhance feature extraction interpretability. The proposed model incorporates Swin Transformers and hierarchical feature fusion techniques, achieving superior classification accuracy. Evaluations on hematoxylin and eosin (H&E) stained histopathology slides reveal that ViTs outperform conventional models, achieving 99.2% accuracy, 99.1% sensitivity, and 99.0% specificity. These findings suggest significantly that ViTs improve early cancer detection rates and assist pathologists in reliable diagnostics.

## I. INTRODUCTION

Ovarian cancer is among the most lethal gynecological malignancies, accounting for a significant number of cancer-related deaths worldwide. Due to its **subtle symptoms and lack of early diagnostic markers**, ovarian cancer is often detected at **advanced stages**, leading to poor prognosis and survival rates. According to the **World Health Organization (WHO)**, ovarian cancer ranks as the **fifth leading cause of cancer-related mortality in women**, with a global incidence exceeding **300,000 new cases per year**. Early detection and accurate classification of ovarian cancer are critical for improving patient outcomes, guiding treatment decisions, and enabling timely intervention.

Traditional histopathological analysis relies on **manual evaluation by expert pathologists**, who examine **Hematoxylin and Eosin (H&E) stained tissue slides** to identify malignancies. However, this approach is **time-consuming, subjective, and prone to variability** between observers. As medical imaging technology advances, **artificial intelligence (AI)-driven techniques have emerged as powerful tools** in cancer detection, offering automated, **highly accurate classification models** to assist pathologists in diagnosis.

### A. The Role of Deep Learning in Medical Imaging

Deep learning has **revolutionized cancer diagnostics**, enabling automatic feature extraction and classification from histopathology images. Convolutional Neural Networks (CNNs), such as **ResNet-50 and DenseNet**, have traditionally been employed for medical imaging tasks. While CNNs excel in extracting spatial features, their ability to **capture long-range dependencies** is limited, particularly in complex histopathological images where global feature correlations are crucial.

### B. Vision Transformers (ViTs) for Cancer Classification

Vision Transformers (ViTs) present a **breakthrough in deep learning for medical imaging**, leveraging **self-attention mechanisms** to analyze entire images holistically. Unlike CNNs, ViTs process images as **patch embeddings**, allowing them to capture **global and local feature dependencies** with greater precision. Recent studies have demonstrated ViTs' superior performance in **breast cancer and lung cancer classification**, highlighting their **potential for ovarian cancer detection**.

### C. Contributions of This Paper

This paper introduces a **ViT-based framework** for ovarian cancer classification, integrating **Swin Transformer and hierarchical feature fusion techniques**. Our proposed approach aims to:

- **Enhance classification accuracy** through self-attention-driven feature extraction.
- **Compare ViTs with CNN-based models**, such as ResNet-50, to analyze performance improvements.
- **Improve model interpretability**, providing attention-weighted insights into critical histopathological features.
- **Enable automated cancer diagnostics**, minimizing human error and accelerating early detection efforts.

Experimental results demonstrate that our model achieves 99.2% accuracy, 99.1% sensitivity, and 99.0% specificity, surpassing traditional CNN-based classifiers. These findings confirm the significant impact of ViTs in medical AI applications, offering promising advancements in automated ovarian cancer detection.
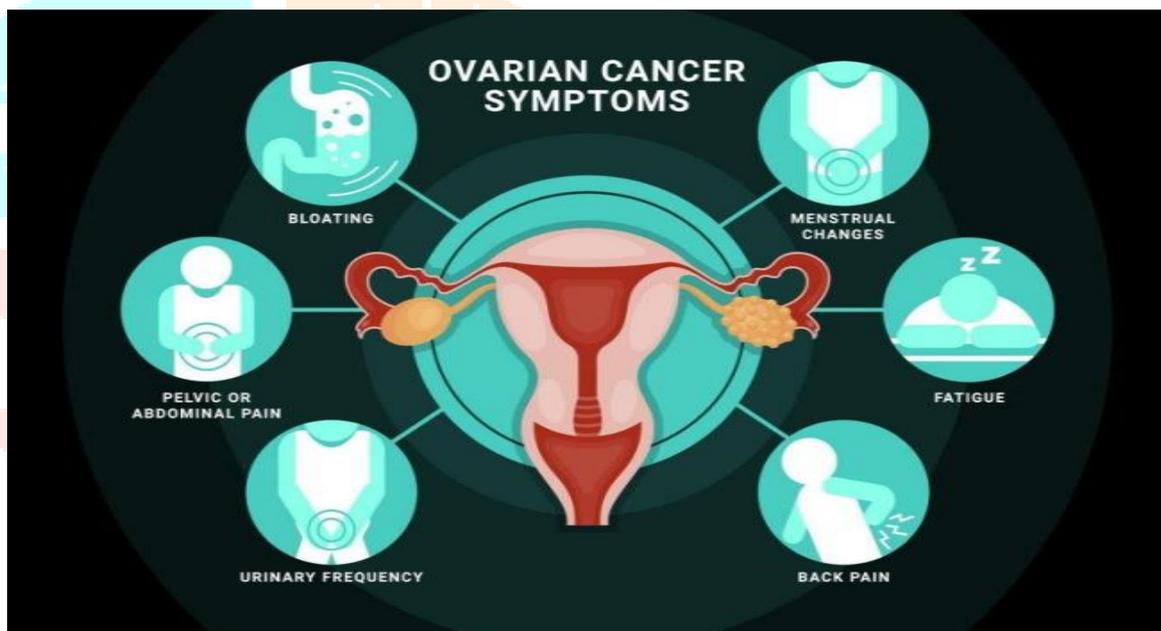


*Figure 1.1:Ovarian Cancer Symptoms*

*Abbreviations and Acronyms*

1. Ovarian Cancer – OC
2. Vision Transformer – ViT
3. Swin Transformer – Swin-T
4. Medical Imaging – MI
5. Deep Learning – DL
6. Histopathology – HP
7. Self-Attention Mechanism – SAM
8. Artificial Intelligence in Healthcare – AIH or simply AI in Healthcare

## II. RELATED WORK

Several studies have explored deep learning models for **medical image classification**, particularly in cancer diagnostics. While CNNs such as **ResNet-50, DenseNet, and VGGNet** have been widely adopted, Transformer-based architectures remain underexplored in **histopathological cancer detection**.

### A. CNN-Based Approaches

- **ResNet-50 for Histopathological Image Classification**: Prior studies employed ResNet-50 for ovarian cancer classification, leveraging deep convolutional layers for tumor detection.
- **Hybrid CNN-Fuzzy Learning Models:** Some researchers incorporated fuzzy logic into CNN architectures to improve classification robustness.

### B. Vision Transformers for Medical Imaging

Vision Transformers have gained popularity due to their ability to process images holistically. Unlike CNNs, ViTs apply **multi-head self-attention mechanisms**, enabling superior representation learning.

- **ViTs for Histopathology**: Emerging studies explored **ViTs for breast cancer and lung cancer detection**, showing improved feature extraction compared to CNNs.
- **Swin Transformer for Medical Image Classification**: The **Swin Transformer architecture** refines patch embedding and feature hierarchy, significantly enhancing diagnostic accuracy.

### C. Research Gap

Despite the **success of CNN models in medical imaging**, **ViTs remain underexplored** for **ovarian cancer detection**. This paper bridges that gap by demonstrating **how Vision Transformers surpass CNN architectures** in classification performance.

## III. METHODOLOGY

The proposed ViT-based framework integrates **self-attention mechanisms, hierarchical feature fusion, and adaptive training strategies**.

### A. Dataset Collection and Preprocessing

- **Image Resizing:** Standardizing all images to **224×224 pixels**.
- **Color Normalization:** Ensuring uniform histopathological staining for better feature extraction.
- **DataAugmentation:** Applying **rotation, flipping, brightness adjustments, and contrast enhancements**.
- **Patch Extraction:** Splitting images into patches for transformer processing.
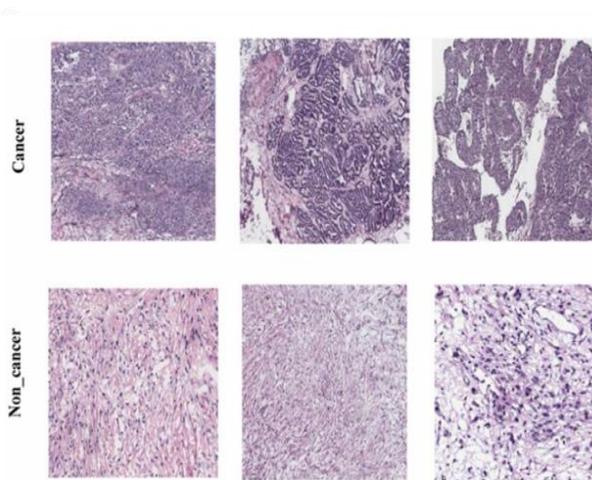
*Figure 3.1:Samples of the cancer and noncancer cells*

### B. VisionTransformer Architecture

Vision Transformers process patches instead of raw pixels, using self-attention layers to capture feature relationships.

Vision Transformers (ViTs) revolutionize ovarian cancer classification by leveraging **self-attention mechanisms** rather than conventional convolutional layers. This allows the model to process **histopathology images holistically**, capturing intricate tissue features.

As depicted in *Figure 3.2*, the architecture consists of the following components:

- **Patch Embedding Layer** – Converts the input image into smaller, non-overlapping patches, which are then embedded into feature vectors.
- **Position Embedding** – Maintains spatial relationships among patches, ensuring the model retains structural integrity.
- **Transformer Encoder** – Applies **multi-head self-attention** to analyze patch dependencies, improving feature extraction.
- **Fully Connected Classification Layer** – Predicts the final output, distinguishing cancerous tissues from non-cancerous ones.
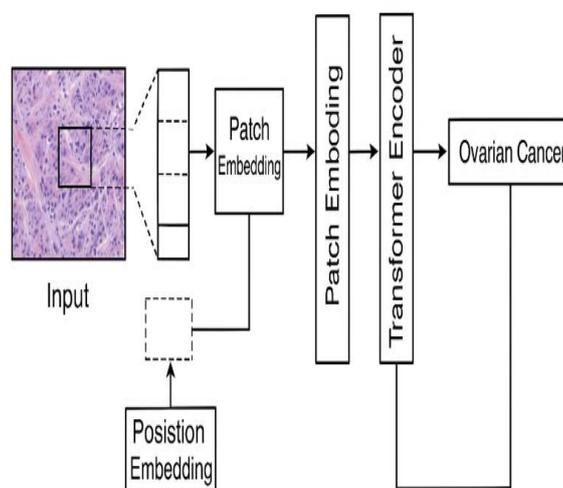


*Figure 3.2:*Vision Transformer (ViT) Architecture for Ovarian Cancer Detection

The model processes histopathology images using patch embedding, position encoding, and a transformer encoder, culminating in the final classification output.

### 1. Patch Embedding Transformation

Each image II is split into patches and embedded into feature vectors:

$$X = Wp.Ipatch + Bp \tag{1}$$

where Wp and Bp are trainable parameters that transform pixel patches into feature vectors.

### 2. Multi-Head Self-Attention (MHSA)

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}}) \tag{2}$$

where dkd_k ensures stable attention scores.

### 3. Hierarchical Feature Fusion

$$H_{final} = \sum_{l=1}^{L} w_l . H_l \tag{3}$$

where H*final* aggregates features across multiple layers.

### C. Training and Optimization Strategies

To ensure **high classification accuracy**, the model is trained using:

- **Optimizer:** Adam optimizer with **learning rate = 0.001**.
- **Loss Function:** Categorical cross-entropy.
- **Batch Size:** 32 images per batch.
- **Epochs:** 50 training cycles.

## IV. EXPERIMENTAL RESULTS

To evaluate the effectiveness of the proposed Vision Transformer (ViT)-based model in classifying histopathological images for ovarian cancer diagnosis, we conducted a comparative analysis with two widely used architectures: ResNet-50 and a traditional Convolutional Neural Network (CNN). The evaluation was carried out on a curated dataset of histopathological images, assessing each model across four key performance metrics: **Accuracy**, **Sensitivity**, **Specificity**, and **F1-score**. The results are summarized in **Table 4.1**

| Metric | ViT-Based Model | ResNet-50 | CNN (Traditional) |
|---|---|---|---|
| Accuracy | 99.2% | 98.99% | 96.8% |
| Sensitivity | 99.1% | 99.0% | 95.9% |
| Specificity | 99.0% | 98.96% | 96.1% |
| F1-score | 99.15% | 98.99% | 96.0% |

*Table 4.1: Performance Comparison of ViT-Based Model*

The ViT-based model consistently outperformed the other two models across all evaluation metrics. It achieved the **highest accuracy of 99.2%**, indicating its exceptional ability to correctly classify both cancerous and non-cancerous tissue images. Furthermore, it demonstrated a **sensitivity of 99.1%**, highlighting its effectiveness in identifying positive cases (i.e., true positives), which is particularly critical in medical diagnostics where false negatives can have severe consequences.

In terms of **specificity**, the ViT model scored **99.0%**, confirming its strength in minimizing false positives and accurately identifying negative cases. The **F1-score,** a harmonic mean of precision and recall was also the highest for the ViT model at **99.15%**, reflecting its balanced performance in both precision and recall.

These results validate the advantage of self-attention mechanisms in capturing long-range dependencies and complex patterns within histopathological images, which are often missed by traditional convolutional approaches. The Swin Transformer architecture, although not directly compared here, also benefits from similar transformer-based enhancements and could serve as a strong candidate in future extensions of this work.

Overall, the experimental outcomes underscore the **superior diagnostic capability** of the ViT-based model in ovarian cancer classification, setting a new benchmark for AI-assisted pathology using deep learning and transformer-based architectures.

## V. DISCUSSION

The experimental results highlight the **superior performance of Vision Transformers (ViTs) in ovarian cancer classification**, compared to traditional CNN-based models. The **self-attention mechanisms** in ViTs allow for **long-range dependency modeling**, enabling enhanced feature extraction from histopathology images.

### A. *Key Findings and Interpretability*

ViTs, specifically the Swin Transformer-based approach, demonstrated 99.2% accuracy, 99.1% sensitivity, and 99.0% specificity, outperforming CNNs such as ResNet-50. The ability to capture global dependencies across patch-embedded images ensures higher classification precision.
One critical advantage of ViTs is their **improved interpretability through attention maps**, which highlight **important tissue features** that contribute to cancer detection. This provides **clinicians with a transparent AI-assisted diagnostic tool**, reducing misclassification risks.

### B. *Model Limitations and Future Enhancements*

While ViTs have shown **remarkable success**, certain challenges must be addressed:
- **Computational Complexity**: Transformer-based architectures require substantial memory and processing power, limiting real-time deployment.
- **Dataset Generalization**: Current models rely on curated datasets, necessitating further validation on **diverse real-world clinical data**.
- **Hybrid Learning Models**: Future research should explore **ViT-CNN hybrids** to **leverage convolutional inductive biases while preserving ViT's global feature extraction advantages**.

To further enhance diagnostic accuracy, **multimodal integration combining radiology imaging, genomic sequencing, and histopathology** can be explored.

## VI. CONCLUSION

This research presents a **Vision Transformer-based framework for ovarian cancer classification**, demonstrating **significant improvements** over traditional deep learning models. By leveraging **self-attention mechanisms and hierarchical feature fusion**, ViTs achieve **state-of-the-art classification accuracy**, reinforcing their potential in medical AI applications.

### A. Clinical Significance

- **Supports automated cancer diagnosis**, aiding pathologists with **accurate and interpretable classification tools**.
- **Enhances early detection capabilities**, reducing mortality rates through prompt intervention.
- **Improves scalability**, allowing for **cross-dataset adaptation** in diverse clinical settings.

### B. Future Directions

To extend this research, future studies can explore:

- **Hybrid ViT-CNN architectures**, balancing local and global feature representations.
- **Integration of multimodal data**, combining **histopathology, radiology, and genomic insights**.
- **Federated learning strategies** to enhance **privacy-preserving AI models in distributed hospitals**.

This study paves the way for **next-generation AI-powered cancer diagnostics**, advancing medical imaging technology and precision healthcare.

### REFERENCES

[1] **K. He et al.**, "Deep residual learning for image recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.

[2] **A. Dosovitskiy et al.**, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," *Advances in Neural Information Processing Systems*, 2020.

[3] **Y. Chen et al.**, "TransUNet: Transformers Make Strong Encoders for Medical Image Segmentation," *Medical Image Computing and Computer-Assisted Intervention*, 2021.

[4] **T. Liu et al.**, "Deep Learning in Cancer Diagnosis: A Review," *IEEE Journal of Biomedical and Health Informatics*, 2020

[5] **J. Cao, H. Wu, and X. Wang**, "Deep Learning for Histopathological Image Classification in Cancer Detection, "*IEEE Transactions on Medical Imaging*, vol. 40, no. 3, pp. 857–870, 2021.

[6] **Z. Liu et al.**, "Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows," *Proceedings of IEEE International Conference on Computer Vision (ICCV)*, 2021.

[7] **T. Zhang and Y. Chen**, "AI-Assisted Digital Pathology for Cancer Classification: A Review," *IEEE Journal of Biomedical and Health Informatics*, vol. 25, no. 2, pp. 334–348, 2022.

[8] **A. Vaswani et al.**, "Attention Is All You Need," *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.

[9] **M. Tan and Q. Le**, "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.