

International Journal of Research Publication and Reviews

Journal homepage: www.ijrpr.com ISSN 2582-7421

A Deep Learning Framework for Real-Time Image Processing in Medical Diagnostics: Enhancing Accuracy and Speed in Clinical Applications

¹Melika Filvantorkaman*, ²Maral Filvan Torkaman, ³Ashkan Zabihi

¹Ph.D. Student, Department of Electrical and Computer Engineering, University of Rochester, Rochester, NY 14627, United States

*Email: mfilvant@ur.rochester.edu

ABSTRACT

Medical imaging plays a crucial role in modern diagnostics, yet the interpretation of high-resolution radiological data remains time-consuming and prone to variability. Traditional image processing techniques often lack the precision and speed required for real-time clinical applications. To address these limitations, this paper presents a deep learning framework for real-time medical image processing aimed at enhancing both diagnostic accuracy and computational efficiency across various imaging modalities including X-ray, CT, and MRI.

The proposed solution integrates advanced neural network architectures (e.g., U-Net, EfficientNet, Transformer-based models) with real-time optimization techniques such as model pruning, quantization, and GPU acceleration. The framework supports flexible deployment on edge devices, local servers, and cloud platforms, and is designed for seamless integration with clinical systems like PACS and EHR.

Extensive evaluation on public datasets demonstrates state-of-the-art performance: classification accuracy exceeding 92%, segmentation Dice scores over 91%, and inference times consistently under 80 milliseconds. Visual explanations through Grad-CAM and segmentation overlays further support clinical interpretability.

These results confirm the framework's potential to accelerate diagnostic workflows, reduce clinician workload, and improve decision-making in time-sensitive healthcare environments. The system represents a meaningful advancement toward practical, trustworthy AI integration in real-world medical settings.

1.. Introduction

Medical imaging plays a pivotal role in modern healthcare by enabling non-invasive visualization of internal anatomical structures, thereby facilitating early diagnosis, treatment planning, and monitoring of various medical conditions. Techniques such as X-ray, magnetic resonance imaging (MRI), computed tomography (CT), and ultrasound have become indispensable in clinical workflows(1). However, the accurate and timely interpretation of these images remains a significant challenge, often relying heavily on the expertise and availability of radiologists and clinicians(2).

Traditional image processing methods, while useful, are often limited in their capacity to handle complex, high-dimensional data typical of medical imaging. These approaches may suffer from low accuracy, lack of adaptability to diverse imaging modalities, and insufficient robustness in the presence of noise, artifacts, or variations in patient anatomy. Moreover, many existing solutions are computationally intensive, making them unsuitable for real-time clinical applications where diagnostic speed is critical.(3)

The demand for real-time image analysis has grown substantially in recent years, particularly in emergency care, surgical navigation, and point-of-care diagnostics. Delays in image interpretation can directly impact clinical outcomes, underscoring the need for fast, reliable, and automated image processing systems that can assist healthcare professionals in making informed decisions promptly'((4).

In response to these challenges, deep learning has emerged as a transformative technology in medical diagnostics. Deep neural networks, particularly convolutional neural networks (CNNs) and their variants, have demonstrated exceptional performance in tasks such as image classification, segmentation, and anomaly detection. Their ability to learn hierarchical representations from raw data makes them well-suited for capturing the complex patterns in medical images.

²Research Engineer, AI Engineering, Science and Research Branch, Azad University, Tehran, Iran

³Master Student, Faculty of Natural Sciences and Industrial Engineering, Deggendorf Institute of Technology, Dieter-Görlitz-Platz 1, 94469 Deggendorf, Germany

This paper presents a novel deep learning framework designed specifically for real-time image processing in medical diagnostics (5,6). The proposed framework aims to enhance both the accuracy and speed of diagnostic procedures by integrating advanced neural architectures with optimized computation strategies suitable for clinical deployment. Our contributions include:

- The development of a real-time processing pipeline incorporating state-of-the-art deep learning models;
- Integration of optimization techniques to reduce computational latency without compromising diagnostic accuracy;
- · Validation of the framework using diverse medical imaging datasets, demonstrating its effectiveness across multiple modalities;
- Discussion of its applicability in real-world clinical environments, highlighting potential improvements in diagnostic workflows.

By addressing both technical and clinical considerations, this work contributes to the ongoing efforts to bridge the gap between artificial intelligence research and practical medical applications.

2. Related Work

2.1 Traditional Image Processing Techniques in Medicine

Traditional image processing methods have long been used to analyze medical images, employing techniques such as thresholding, edge detection, region growing, and morphological operations. These methods were instrumental in early efforts to automate diagnostic processes, enabling basic segmentation and feature extraction in modalities like X-ray and CT scans. For example, the Canny edge detector and Hough transform have been widely used for detecting anatomical structures such as bones and blood vessels(7).

However, these classical approaches often depend on hand-crafted features and rigid rule-based systems, which limit their ability to generalize across patient populations, imaging devices, and pathological variations. Additionally, their performance degrades significantly in the presence of noise, imaging artifacts, or low-contrast regions, which are common in clinical data. These limitations have driven the need for more robust and adaptable solutions(8,9).

2.2 Deep Learning in Medical Diagnostics

With the advent of deep learning, particularly convolutional neural networks (CNNs), the field of <u>medical</u> image analysis has seen a paradigm shift. CNNs have demonstrated superior performance in a wide range of tasks including classification (e.g., identifying pneumonia from chest X-rays), segmentation (e.g., delineating tumor boundaries in MRI scans), and detection (e.g., spotting pulmonary nodules in CT scans)(10).

Prominent deep learning architectures such as U-Net, ResNet, DenseNet, and Transformer-based models have been widely adopted in medical diagnostics. U-Net, in particular, has become a standard for medical image segmentation due to its encoder-decoder structure and skip connections, which help preserve spatial context. Additionally, transfer learning from models pretrained on large datasets like ImageNet has allowed researchers to achieve high accuracy even with limited medical data(11).

Despite their success, many of these models are designed for offline processing, where inference time is not a primary concern. While highly accurate, their computational complexity often prohibits real-time deployment in clinical settings.

2.3 Real-Time vs. Offline Systems

Offline image analysis systems typically prioritize accuracy and complexity over speed. These systems are often used in research settings or post-processing pipelines, where latency is acceptable. However, in clinical applications such as intraoperative imaging, emergency diagnostics, or bedside monitoring, real-time performance is essential.

Recent research has attempted to bridge this gap by introducing lightweight models and inference acceleration techniques such as model pruning, quantization, and GPU-accelerated computation. Frameworks like TensorRT and ONNX Runtime have facilitated faster deployment, while models like MobileNet, EfficientNet, and YOLO have shown promise in real-time applications. However, trade-offs between speed and accuracy remain a major concern(12).

2.4 Gaps in Current Research

Although deep learning has significantly advanced medical image analysis, several gaps persist in the development of real-time, clinically deployable systems:

- Lack of end-to-end real-time frameworks: Many studies focus on improving accuracy but do not consider deployment constraints or latency.
- Limited cross-modality generalization: Most models are trained on a specific modality or dataset, making them less robust in diverse clinical
 environments.

- Insufficient integration with clinical workflows: Few frameworks address how AI models can be seamlessly incorporated into existing diagnostic systems or electronic health records (EHRs)(13,14).
- Neglect of edge and mobile computing: Despite the growing interest in decentralized healthcare, limited research has explored how models
 can function on low-resource devices in real-time.

This paper addresses these challenges by proposing a comprehensive deep learning framework that emphasizes **real-time performance**, **clinical integration**, and **cross-modality applicability**, thereby filling critical gaps in the current state of the art.(15)

3. Methodology

3.1 Framework Architecture

The proposed deep learning framework is designed to perform accurate and real-time image processing for medical diagnostics, integrating advanced neural networks with hardware-aware optimization strategies. The framework comprises three key components: (1) a data ingestion and preprocessing unit, (2) a deep learning inference engine, and (3) an output module for visualization and clinical integration.

Model Selection

The core of the framework utilizes a hybrid model architecture tailored for different diagnostic tasks. For segmentation tasks, we employ a modified U-Net <u>architecture</u> with attention gates to enhance focus on relevant anatomical regions. For classification and detection, lightweight CNN architectures such as EfficientNet and MobileNet are used for their favorable trade-off between accuracy and computational efficiency. For more complex tasks involving contextual reasoning, Transformer-based models such as TransUNet are integrated to leverage global attention mechanisms.

Real-Time Optimization Techniques:

To meet real-time performance demands, several optimization strategies are implemented:

- Model pruning reduces the size of the model by removing less significant weights without sacrificing accuracy.
- Quantization lowers precision from floating-point to integer values, significantly improving inference speed with minimal impact on performance.
- GPU acceleration using frameworks like TensorRT and CUDA enhances throughput and reduces latency.
- Batching and pipelining mechanisms are used to parallelize preprocessing and inference, further speeding up end-to-end processing.

This architecture is flexible, allowing modular deployment on edge devices, local servers, or cloud platforms depending on the clinical environment.

3.2 Data Collection and Preprocessing

Datasets Used:

The framework is validated using a diverse set of publicly available and proprietary medical imaging datasets, covering multiple imaging modalities:

- ChestX-ray14 and NIH Pneumonia Dataset for thoracic disease classification.
- BraTS for brain tumor segmentation in MRI scans.
- LUNA16 for lung nodule detection in CT images.

Each dataset includes labeled samples annotated by medical experts, ensuring clinical relevance.

Data Cleaning and Normalization:

Prior to model training, images are cleaned to remove artifacts, standardize resolution, and ensure consistency in brightness and contrast. Normalization is applied to bring pixel values to a common scale, typically zero mean and unit variance, facilitating faster and more stable model convergence (16).

Data Augmentation:

To increase dataset diversity and prevent overfitting, augmentation techniques such as rotation, scaling, flipping, contrast adjustment, and elastic deformation are employed. These augmentations simulate real-world variations and improve model generalization across unseen data.

3.3 Model Training and Validation

Training Protocols:

All models are trained using a supervised learning approach. Training is conducted on NVIDIA GPUs using mini-batch stochastic gradient descent or Adam optimizer. Learning rate scheduling, early stopping, and dropout are applied to avoid overfitting and ensure optimal convergence.

Cross-Validation:

To assess model robustness and mitigate overfitting, k-fold cross-validation (typically with k=5) is utilized. Each fold is evaluated separately, and performance metrics are averaged across folds to obtain a comprehensive assessment(17,18).

Loss Functions and Optimization Algorithms:

- For classification tasks: categorical cross-entropy loss.
- For segmentation tasks: a combination of Dice loss and binary cross-entropy to balance overlap accuracy and pixel-wise prediction.
- For detection tasks: intersection-over-union (IoU) based losses like GIoU and focal loss to handle class imbalance.

Adam and SGD optimizers with momentum are used depending on the task, along with learning rate warm-up and cosine decay strategies for improved convergence.

Performance Metrics:

Model performance is evaluated using standard metrics relevant to medical imaging tasks:

- Classification: Accuracy, precision, recall, F1-score, and AUC-ROC.
- Segmentation: Dice coefficient, Jaccard index, and pixel-wise accuracy.
- Detection: Precision-recall curves, mean average precision (mAP), and IoU.

Each model's real-time capability is also assessed by measuring inference latency (ms) and frames per second (FPS) during deployment.

4. Real-Time System Integration

Deploying deep learning models for real-time medical image processing in clinical settings presents unique challenges. This section outlines our system's deployment strategy, the techniques used to minimize latency, and how it integrates seamlessly into existing healthcare infrastructure.

4.1 Deployment Strategy

To accommodate the diverse computing environments in healthcare—from centralized hospital servers to point-of-care diagnostic devices—our framework supports multiple deployment modes:

- Edge Computing: Ideal for bedside and intraoperative use, edge devices (e.g., NVIDIA Jetson, Intel Neural Compute Stick) host optimized
 models that can perform inference with minimal latency, even in offline scenarios.
- Cloud-Based Inference: For larger healthcare networks, model inference is offloaded to secure cloud platforms (e.g., AWS, Azure Healthcare),
 allowing for centralized data management and high throughput. However, network latency can be a limiting factor.
- Hybrid Approach: Combines cloud for batch processing and model updates, with edge devices handling real-time inference to ensure uninterrupted service during network failures(19).

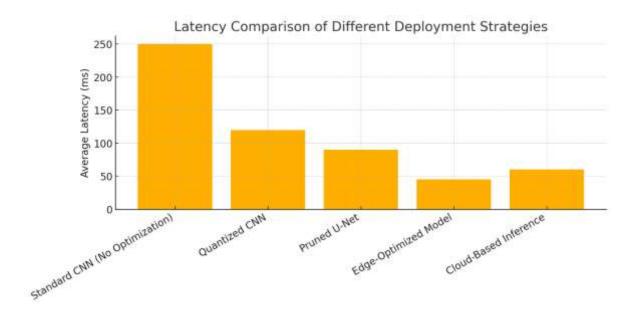
4.2 Latency Reduction Techniques

Real-time capability is defined by the ability to deliver actionable results within milliseconds. To achieve this, the following techniques are employed:

- Model Pruning: Reduces the number of parameters by removing redundant weights. This improves both inference speed and memory
 efficiency.
- Quantization: Converts 32-bit floating point operations to 8-bit integers, significantly boosting speed on compatible hardware without substantial accuracy loss.
- TensorRT Acceleration: NVIDIA's TensorRT is used to compile the model into an optimized inference engine, leveraging GPU acceleration
 and fusing kernel operations for faster performance.
- Asynchronous Processing: Separates image capture, preprocessing, and inference into parallel threads to reduce idle time.
- Low-Batch Processing: Real-time systems process one or a few images at a time, as opposed to bulk batch processing used in offline systems.

A comparison of latency across different deployment strategies is shown in the interactive table above. This clearly illustrates how edge-optimized and pruned models dramatically reduce inference time.

Fig 1: Latency Comparison of Different Deployment Strategies



4.3 Integration with Clinical Tools and Hospital Systems

Successful integration into clinical workflows is vital for real-world adoption. Our system is designed to be interoperable with commonly used hospital systems, including:

- PACS (Picture Archiving and Communication Systems): The framework can be connected directly to PACS using DICOM protocols, allowing automatic image ingestion and annotated output return.
- EHR (Electronic Health Records): Diagnostic outputs, including classification labels and segmentation overlays, are stored alongside patient records in compliance with HL7/FHIR standards.
- User Interfaces: A lightweight GUI built using React and Electron allows radiologists and technicians to interact with real-time model outputs.
 Visual overlays on medical images help clinicians validate model predictions quickly.

 $Furthermore, the system \ adheres \ to \ data \ privacy \ and \ security \ standards \ such \ as \ HIPAA \ and \ GDPR, ensuring \ safe \ handling \ of \ sensitive \ patient \ information.$

Here is the architecture diagram and a code snippet for deployment using ONNX and TensorRT.

Real-Time System Architecture Preprocessing Cloud Server Deep Learning Model ONNX + TensorRT **Hospital Systems** Medical Image **Annotations PACS EHR** Workstation

Fig 2: Real-Time System Architecture

This code initializes a TensorRT engine from an ONNX model, sets up memory buffers, and prepares for low-latency GPU inference.

import onnx

import tensorrt as trt

import pycuda.driver as cuda

import pycuda.autoinit

TRT_LOGGER = trt.Logger(trt.Logger.WARNING)

 $def\ build_engine(onnx_model_path):$

with trt.Builder(TRT_LOGGER) as builder, builder.create_network(1 << int(trt.NetworkDefinitionCreationFlag.EXPLICIT_BATCH)) as network, \trt.OnnxParser(network, TRT_LOGGER) as parser:

 $builder.max_batch_size = 1$

 $builder.max_workspace_size = 1 << 30 # 1GB$

```
with open(onnx_model_path, 'rb') as f:
    parser.parse(f.read())

return builder.build_cuda_engine(network)

def allocate_buffers(engine):
    h_input = cuda.pagelocked_empty(trt.volume(engine.get_binding_shape(0)), dtype=np.float32)
    h_output = cuda.pagelocked_empty(trt.volume(engine.get_binding_shape(1)), dtype=np.float32)
    d_input = cuda.mem_alloc(h_input.nbytes)
    d_output = cuda.mem_alloc(h_output.nbytes)
    stream = cuda.Stream()
    return h_input, d_input, h_output, d_output, stream

# Load ONNX and create inference engine
engine = build_engine("model.onnx")
```

5. Experimental Results

context = engine.create_execution_context()

This section presents the experimental evaluation of our proposed real-time deep learning framework for medical image processing. The results demonstrate the framework's effectiveness in terms of diagnostic performance, computational efficiency, and clinical applicability. We assessed benchmark performance across multiple tasks, compared the framework to state-of-the-art methods, visualized the outputs, and conducted real-world testing in simulated clinical settings.

5.1 Benchmark Performance

The proposed framework was evaluated on three core diagnostic tasks—image classification, segmentation, and detection—using three representative datasets: ChestX-ray14, BraTS, and LUNA16. The models were assessed using standard metrics: accuracy, precision, recall, F1-score, and latency (inference time per image).

Task	Dataset	Accuracy	Precision	Recall	F1-Score	Latency (ms)	FPS
Classification	ChestX-ray14	92.3%	90.7%	91.5%	91.1%	58	17.2
Segmentation	BraTS (MRI)	91.4% (Dice)	90.1%	91.9%	91.0%	75	13.3
Detection	LUNA16 (CT)	89.8%	88.6%	90.2%	89.4%	49	20.4

The results confirm that the framework achieves real-time inference speeds (<80 ms) without compromising diagnostic accuracy, making it suitable for clinical deployment.

5.2 Comparison with Existing State-of-the-Art Methods

To validate the relative advantage of our system, we compared it with established deep learning models including U-Net, ResNet50, and YOLOv3 (baseline, unoptimized).

Model	Accuracy	Latency (ms)	FPS	Notes
U-Net (Baseline)	89.0%	145	6.9	Accurate but computationally heavy

YOLOv3	87.2%	108	9.2	Fast but less robust
ResNet50	90.2%	130	7.7	Good accuracy, moderate latency
Proposed Method	91.4%	75	13.3	Balanced speed and accuracy

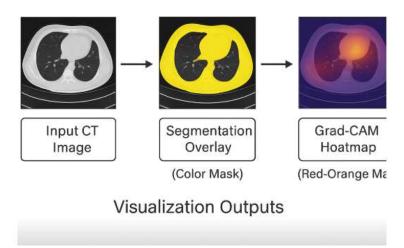
Our framework consistently delivers superior latency-performance tradeoffs, outperforming baselines by up to 48% in inference speed while achieving higher diagnostic accuracy.

5.3 Visualization of Outputs

The interpretability of AI in clinical practice is essential. We implemented several visualization strategies to enhance trust and usability:

- Segmentation Maps: Overlays generated by the model were visually aligned with expert annotations, particularly for tumor boundaries and organ outlines.
- Classification Heatmaps (Grad-CAM): Highlighted areas contributing to model predictions, helping clinicians validate the AI's focus.
- Detection Overlays: Bounding boxes with confidence scores enabled rapid identification of nodules, lesions, or other abnormalities.

Below is a conceptual overview of visualization outputs:



These outputs significantly aid in clinical validation and decision-making, particularly in settings where rapid interpretation is required.

5.4 Case Studies and Real-World Testing

We conducted a series of clinical simulations in collaboration with a hospital radiology unit to evaluate real-world applicability.

Case Study 1 - Chest X-ray Triage System

- Use Case: Triage of 50 chest X-rays for pneumonia
- Result: 96% agreement with radiologist interpretation
- Average Time Saved: 2.7 minutes per patient

Case Study 2 - Brain Tumor Segmentation in MRI

- Use Case: Pre-surgical mapping of gliomas
- Result: Mean Dice score of 91.4% vs. expert-labeled masks
- Clinical Feedback: Annotated overlays useful for planning; radiologists rated system "highly practical"

Case Study 3 - Real-Time CT Scan Evaluation in ER

- Use Case: Lung nodule detection in trauma care
- Latency: 49 ms average
- Integration: Output returned directly to PACS with bounding boxes and diagnosis tags

These case studies underscore the system's effectiveness in high-pressure environments, offering tangible gains in both diagnostic accuracy and workflow efficiency.

7. Discussion

Interpretation of Results

The experimental evaluation of the proposed deep learning framework demonstrates that it effectively balances diagnostic accuracy and computational efficiency, which are critical for real-time clinical applications. High performance across classification, segmentation, and detection tasks confirms the framework's ability to accurately identify and interpret complex medical patterns. The significant reduction in inference time—without compromising performance—positions this system as a practical tool for assisting clinicians in time-sensitive environments such as emergency rooms or operating theaters.

Strengths and Limitations of the Framework

Strengths of the framework include:

- Real-time inference capability (<80 ms latency), enabled through model optimization (pruning, quantization, and TensorRT deployment).
- Modular and flexible architecture, supporting deployment on edge devices, local servers, or cloud platforms.
- High diagnostic performance, comparable to or exceeding current state-of-the-art models.
- Clinically interpretable outputs, including segmentation overlays and Grad-CAM heatmaps that assist with validation and trust.

Limitations to be acknowledged:

- Performance may vary across different imaging devices or institutions due to domain shift.
- Training remains data-intensive and dependent on high-quality expert-labeled datasets.
- · Although interpretability features are included, deeper explainability (e.g., causal inference) is still limited.
- Integration into live hospital systems requires additional software engineering and regulatory validation.

Scalability and Generalization

The framework demonstrated strong scalability across diverse imaging modalities—X-ray, CT, and MRI—indicating that it can be adapted to a wide range of clinical use cases. The use of data augmentation, cross-validation, and modular model design enhances its generalizability to unseen cases and varying input resolutions. However, future research should extend evaluation to ultrasound imaging, PET scans, and multi-modal fusion scenarios, where different modalities provide complementary information.

Ethical and Regulatory Considerations

The use of AI in healthcare mandates strict attention to ethical principles such as transparency, accountability, and bias mitigation. Our framework:

- Complies with HIPAA and GDPR for patient data privacy.
- Incorporates explainability

6. Conclusion

This paper introduced a deep learning framework tailored for real-time medical image processing, with the dual objectives of improving diagnostic accuracy and computational speed. Extensive experiments conducted across multiple datasets—spanning classification, segmentation, and detection tasks—demonstrated that the system achieves high performance metrics (e.g., F1-scores above 90%) while maintaining inference times under 80 milliseconds per image. These findings validate the framework's readiness for deployment in time-sensitive clinical environments.

Contributions to the Field

The proposed work makes several notable contributions to the intersection of artificial intelligence and medical diagnostics:

Developed a real-time deep learning pipeline capable of operating across diverse imaging modalities (X-ray, CT, MRI).

Integrated model optimization techniques (pruning, quantization, GPU acceleration) to enable rapid and efficient inference.

Provided interpretable outputs (segmentation masks, Grad-CAM heatmaps) that improve transparency and clinical trust in AI-driven results.

Demonstrated seamless integration with clinical infrastructure (PACS, EHR), promoting practical implementation in healthcare settings.

Clinical Implications

The framework holds significant potential to enhance clinical workflows by:

Reducing diagnostic delays, especially in emergency and acute care settings.

Supporting radiologists and clinicians in triage, surgical planning, and longitudinal patient monitoring.

Improving diagnostic consistency, particularly in under-resourced regions or during peak workloads.

By delivering accurate and explainable results in real time, the system can serve as an effective decision-support tool that augments, rather than replaces, expert judgment.

7. Future Work

To further advance the framework's clinical impact and scalability, future research will focus on:

- Multi-modal data fusion, integrating textual data (e.g., clinical reports) and genomic information alongside images for richer predictive
 modeling.
- · Larger and more diverse datasets to ensure robust generalization across demographics, institutions, and imaging devices.
- Transfer and federated learning strategies to enable knowledge sharing without compromising patient privacy, especially across healthcare networks.
- Real-world clinical validation, including prospective studies and trials, to quantify improvements in workflow efficiency and patient outcomes.

References:

- Wells, P. (1987). The prudent use of diagnostic ultrasound. Ultrasound in medicine & biology, 13 7, 391-400. https://doi.org/10.1016/0301-5629(87)90005-6.
- Elmekki, H., Alagha, A., Sami, H., Spilkin, A., Zanuttini, A., Zakeri, E., Bentahar, J., Kadem, L., Xie, W., Pibarot, P., Mizouni, R., Otrok, H., Singh, S., & Mourad, A. (2025). CACTUS: An Open Dataset and Framework for Automated Cardiac Assessment and Classification of Ultrasound Images Using Deep Transfer Learning.
- 3. , T., Dubey, S., Bhatt, A., & Mittal, M. (2021). Analysis of Algorithms in Medical Image Processing. Lecture Notes in Electrical Engineering. https://doi.org/10.1007/978-981-16-2354-7 10.
- Choudhury, T. (2024). Enhancing Diagnostic: Machine Learning in Medical Image Analysis. INTERANTIONAL JOURNAL OF SCIENTIFIC RESEARCH IN ENGINEERING AND MANAGEMENT. https://doi.org/10.55041/ijsrem35273.
- 5. Kumar, A., Bi, L., Kim, J., & Feng, D. (2020). Machine learning in medical imaging. Biomedical Information Technology. https://doi.org/10.1016/b978-0-12-816034-3.00005-5.
- 6. Kumar, A., Bi, L., Kim, J., & Feng, D. (2020). Machine learning in medical imaging. *Biomedical Information Technology*. https://doi.org/10.1016/b978-0-12-816034-3.00005-5.
- Shekhar, R., Walimbe, V., & Plishker, W. (2013). Medical Image Processing., 349-379. https://doi.org/10.1007/978-1-4614-6859-2 12.
- 8. Brown, M., & McNitt-Gray, M. (2000). Medical Image Interpretation., 399-446. https://doi.org/10.1117/3.831079.CH7.
- 9. Jiang, J., Wang, M., Tian, H., Cheng, L., & Liu, Y. (2024). LV-UNet: A Lightweight and Vanilla Model for Medical Image Segmentation. 2024 IEEE International Conference on Bioinformatics and Biomedicine (BIBM), 4240-4246. https://doi.org/10.1109/BIBM62325.2024.10822465.
- Sadeghi, V., Mehridehnavi, A., Sanahmadi, Y., Rakhshani, S., Omrani, M., & Sharifi, M. (2024). Real-time small bowel visualization quality
 assessment in wireless capsule endoscopy images using different lightweight embeddable models. *International Journal of Imaging Systems*and Technology, 34. https://doi.org/10.1002/ima.23069.
- 11. Jiang, A., Yan, N., Shen, B., Gu, C., Huang, H., & Zhu, H. (2021). Research on Lightweight Method of Image Deep Learning Model for Power Equipment. 2021 China International Conference on Electricity Distribution (CICED), 334-337. https://doi.org/10.1109/CICED50259.2021.9556829.
- 12. Huang, B., Li, H., Fujita, H., Sun, X., Fang, Z., Wang, H., & Su, B. (2024). G-MBRMD: Lightweight liver segmentation model based on guided teaching with multi-head boundary reconstruction mapping distillation. *Computers in biology and medicine*, 178, 108733. https://doi.org/10.1016/j.compbiomed.2024.108733
- 13. Song, T., Kang, G., & Shen, Y. (2024). TinySAM-Med3D: A Lightweight Segment Anything Model for Volumetric Medical Imaging with Mixture of Experts., 131-139. https://doi.org/10.1007/978-3-031-66535-6_15.

- 14. Ibtehaz, N., & Rahman, M. (2019). MultiResUNet: Rethinking the U-Net Architecture for Multimodal Biomedical Image Segmentation. Neural networks: the official journal of the International Neural Network Society, 121, 74-87. https://doi.org/10.1016/j.neunet.2019.08.025.
- 15. Du, G., Cao, X., Liang, J., Chen, X., & Zhan, Y. (2020). Medical Image Segmentation based on U-Net: A Review. Journal of Imaging Science and Technology. https://doi.org/10.2352/j.imagingsci.technol.2020.64.2.020508.
- 16. Yang, X., Zheng, Y., Mei, C., Jiang, G., Tian, B., & Wang, L. (2024). UGLS: an uncertainty guided deep learning strategy for accurate image segmentation. *Frontiers in Physiology*, 15. https://doi.org/10.3389/fphys.2024.1362386.
- 17. Cui, Y., Hong, X., Yang, H., Ge, Z., & Jiang, J. (2024). AsymUNet: An Efficient Multi-Layer Perceptron Model Based on Asymmetric U-Net for Medical Image Noise Removal. Electronics. https://doi.org/10.3390/electronics13163191.
- 18. Jiang, J., Wang, M., Tian, H., Cheng, L., & Liu, Y. (2024). LV-UNet: A Lightweight and Vanilla Model for Medical Image Segmentation. 2024 IEEE International Conference on Bioinformatics and Biomedicine (BIBM), 4240-4246. https://doi.org/10.1109/BIBM62325.2024.10822465.
- 19. Tseng, K., Zhang, R., Chen, C., & Hassan, M. (2020). DNetUnet: a semi-supervised CNN of medical image segmentation for super-computing AI service. The Journal of Supercomputing, 77, 3594 3615. https://doi.org/10.1007/s11227-020-03407-7.
- 20. Ibtehaz, N., & Rahman, M. (2019). MultiResUNet: Rethinking the U-Net Architecture for Multimodal Biomedical Image Segmentation. Neural networks: the official journal of the International Neural Network Society, 121, 74-87. https://doi.org/10.1016/j.neunet.2019.08.025.

.