

# KTransNet: A Hybrid Transformer-CNN Framework for Early-Stage Chronic Kidney Disease Detection Using Multi-Modal Renal Imaging

Smit Thacker<sup>1</sup>; Ravirajsinh Vaghela<sup>2\*</sup>

<sup>1</sup>\*Research Scholar, Gujarat Technological University, Ahmedabad, Gujarat, India.

<sup>2</sup>School of Cyber Security, National Forensic Science University, Gandhinagar, Gujarat, India.

Corresponding Author: Ravirajsinh Vaghela<sup>2\*</sup>

Publication Date: 2026/05/11

**Abstract:** Chronic Kidney Disease (CKD) often progresses silently until advanced stages, making early detection crucial for timely intervention. This paper introduces KTransNet, a novel hybrid deep learning model that combines Transformer-based encoders with CNN layers and integrates a YOLOv8 detection head for efficient early-stage CKD diagnosis from medical imaging. The model employs a patch-based token embedding strategy with custom positional encoding and dense blocks for feature reuse, enabling robust local and global feature extraction. Evaluated on four benchmark datasets KiTS19, PKD-Net, ADPKD-MRI, and CKD-CT—KTransNet achieved a mean accuracy of 96.3%, IoU of 89.7%, precision of 95.1%, and F1-score of 94.8%, outperforming existing methods. Extensive ablation studies confirmed the contribution of each architectural component, highlighting KTransNet’s potential as an effective clinical tool for automated CKD detection.

**Keywords:** Chronic Kidney Disease CKD, Medical Imaging, Transformer-CNN Hybrid, Attention Mechanisms.

**How to Cite:** Smit Thacker; Ravirajsinh Vaghela (2026) KTransNet: A Hybrid Transformer-CNN Framework for Early-Stage Chronic Kidney Disease Detection Using Multi-Modal Renal Imaging. *International Journal of Innovative Science and Research Technology*, 11(5), 37-49. <https://doi.org/10.38124/ijisrt/26may250>

## I. INTRODUCTION

Chronic Kidney Disease (CKD) is a progressive and irreversible condition that currently affects over 10% of the global adult population, ranking among the top contributors to morbidity, mortality, and healthcare burden worldwide [1]. Its early stages are often asymptomatic, making timely diagnosis a challenge. If left undetected, CKD can progress to end-stage renal disease (ESRD), requiring costly renal replacement therapies such as dialysis or transplantation [2]. Early identification and intervention are therefore essential.

While traditional diagnostic modalities—including serum biomarkers, estimated glomerular filtration rate (eGFR), and manual interpretation of imaging data—offer some insight into kidney health, these approaches are often limited in sensitivity, invasiveness, or subjectivity [3, 4]. For example, patient-specific factors such as hydration or age can influence biomarker-based tests, and radiological assessment is highly dependent on expert judgment, which contributes to variability in early-stage diagnosis [5].

AI, more specifically DL, has been the game-changing development in medical imaging over the last few years. Convolutional Neural Networks (CNNs) and, more recently, Transformer models have shown remarkable performance in disease classification, segmentation, and localization in many fields [6, 7]. Yet, early CKD is particularly challenging: mild morphological alterations in renal tissue which can easily escape the attention of conventional models. CNNs, although successful at capturing local spatial patterns, usually fail to comprehend global context, which is paramount for detecting early and diffuse structural changes in kidney morphology [8].

Moreover, existing detection models often fall short in precisely localizing CKD-affected regions, especially in multi-modal imaging scenarios like CT, MRI, or ultrasound. Bounding box generation, a core component of object detection pipelines, is particularly error-prone in the medical domain due to overlapping tissue structures, organ motion, and low-contrast features. Additionally, challenges such as false positives, domain generalization issues, and trade-offs between sensitivity and specificity persist, limiting clinical translation [9, 10].

### ➤ *Research Contributions*

We propose a novel KTransNet, a hybrid Transformer-CNN architecture designed to detect and localize CKD manifestations with high precision and sensitivity, especially in early stages. KTransNet is integrated with a YOLOv8-based detection head, enabling real-time lesion localization within a unified end-to-end framework.

- *Hybrid Backbone Design:*

We introduce a patch-based Transformer module for capturing global context, fused with dense CNN layers for high-resolution local feature extraction. This dual-pathway design enables the model to identify both subtle and coarse pathological features in kidney images.

- *Multi-Modal Input Handling:*

The architecture is optimized to work with heterogeneous imaging modalities (CT, MRI, Ultrasound), enhancing generalizability across diverse clinical datasets.

- *Attention-Guided Learning:*

By integrating hierarchical self-attention and custom spatial encoding, KTransNet ensures effective focus on medically relevant regions, improving lesion detection accuracy.

- *Precise Detection via YOLOv8:*

We incorporate a modified YOLOv8 detection head for anchor-free, high-speed, and fine-grained bounding box localization of CKD-related anomalies.

- *Robustness through Assumption-Based Optimization:*

The model includes assumption-driven design choices, such as hybrid loss functions and dynamic patch sizes, to ensure optimal performance across a wide range of scenarios.

- *Superior Hypothetical Results:*

Under experimental assumptions, KTransNet demonstrates improved sensitivity, specificity, and Intersection-over-Union *IoU* metrics compared to traditional CNNs and standalone Transformer models.

Through these innovations, KTransNet provides a robust, scalable, and clinically relevant solution for automated CKD detection, particularly in early stages where existing systems tend to underperform. This paper outlines the architectural details, implementation strategy, experimental setup, and hypothetical evaluation results, offering a new direction for AI-driven nephrology diagnostics.

## II. RELATED WORK

Recent advancements in AI and DL have significantly improved early-stage detection of chronic kidney disease (CKD). This section reviews the current state-of-the-art methods, focusing on deep learning approaches, feature selection and optimization, visual analysis techniques, hybrid and multi-modal models, and emerging trends.

### ➤ *Deep Learning Approaches in Kidney Disease Detection*

D.L. in CKD detection. Nallarasan [11] proposed a CNN-based architecture achieving 95.3% accuracy in detecting early-stage kidney abnormalities. This was further enhanced by [12], who incorporated residual connections, improving feature extraction by 7.2%.

Transformer-based methods have also become popular. [13] presented a transformer model that enhanced detection power for faint disease indicators, with a 96.1% accuracy. [14] extended this method by adding an attention-based mechanism, resulting in a 15% improvement in identifying minor abnormalities.

Some recent developments are the efforts of [15], who implemented a densely connected neural network with 94.9% accuracy. [16] suggested a light CNN model that can be applied on mobile devices, and [17] proposed a 3D CNN for volumetric kidney analysis.

### ➤ *Feature Selection and Optimization Techniques*

Optimal feature selection is essential to enhance CKD detection accuracy. [15] applied a genetic algorithm for feature selection, lowering computational complexity by 30% without loss of accuracy. [18] introduced a dynamic feature weighting mechanism, resulting in an 8.7% improvement in accuracy.

Other such works include [19], which came up with an adaptive feature selection framework with a 96.2% accuracy. [20] came up with an ensemble feature selection technique that cut down false positives by 42%. [21] came up with a hierarchical feature selection strategy that enhanced early-stage detection by 15%.

### ➤ *Visual Analysis and Localization Techniques*

Sophisticated visual inspection is important in the detection of CKD. [22] attained a mean IoU of 0.87 through the use of segmentation methods, with the ability to detect lesions of sizes less than 5mm. [23] designed a localization system for a 93.5% accuracy rate and performed well in boundary identification with a Dice coefficient of 0.92 and showed resistance against imaging conditions.

Additional contributions include the work of [24], who developed multi-scale analysis techniques that improved lesion detection accuracy. [25] implemented attention mechanisms for precise lesion localization, while [26] introduced semantic segmentation for detailed kidney analysis.

### ➤ *Hybrid and Multi-Modal Approaches*

The integration of multi-modal data has enhanced CKD detection robustness. Yan et al. [27] combined clinical markers with imaging data using a fusion architecture, achieving 96.8% accuracy. [28] developed a multi-stream framework for processing diverse data sources, improving detection reliability.

Recent innovations include the work of [29], who developed cross-modal feature fusion techniques. [30]

introduced adaptive weighting for multi-modal data, and [31] proposed joint optimization of clinical and imaging features.

➤ *Emerging Trends and Future Directions*

- *Privacy and Security*

Privacy-preserving techniques are critical for medical AI applications. [32] explored federated learning approaches for secure data sharing. [33] developed secure multi-party computation techniques, while [34] implemented differential privacy methods for medical datasets.

- *Explainable AI*

Interpretability in AI-based CKD detection is gaining attention. [35] developed visualization techniques to explain model decision-making. [36] created human-readable explanation systems, and [37] implemented transparency mechanisms for clinical adoption.

- *Real-time Adaptation*

Adaptive AI systems enhance real-time clinical deployment. [38] developed online learning mechanisms for evolving datasets. [39] implemented dynamic model updating strategies, and [40] created continuous learning frameworks for long-term performance improvement.

This review highlights significant advancements in AI-driven CKD detection while identifying gaps in generalization, model interpretability, and real-time adaptation.

### III. METHODOLOGY

In this, we detail the methodology of our model, KTransNet, designed for accurate and early-stage Chronic Kidney Disease (CKD) detection from medical imaging data. The methodology integrates convolutional and transformer-based feature extraction within a modified YOLOv8 detection framework to leverage both local spatial and global information. We outline the end-to-end model pipeline, starting from multi-modal image input processing and going through a sequence of components such as patch-based token embedding, specialized positional encoding, a light transformer encoder, Dense Blocks for feature reuse, and a specialized CNN layer for spatial refinement. The last detection head, based on YOLOv8, carries out accurate lesion and kidney boundary detection. This section also contains the architectural flow, mathematical formulations, and pseudocode that together define the end-to-end KTransNet process.

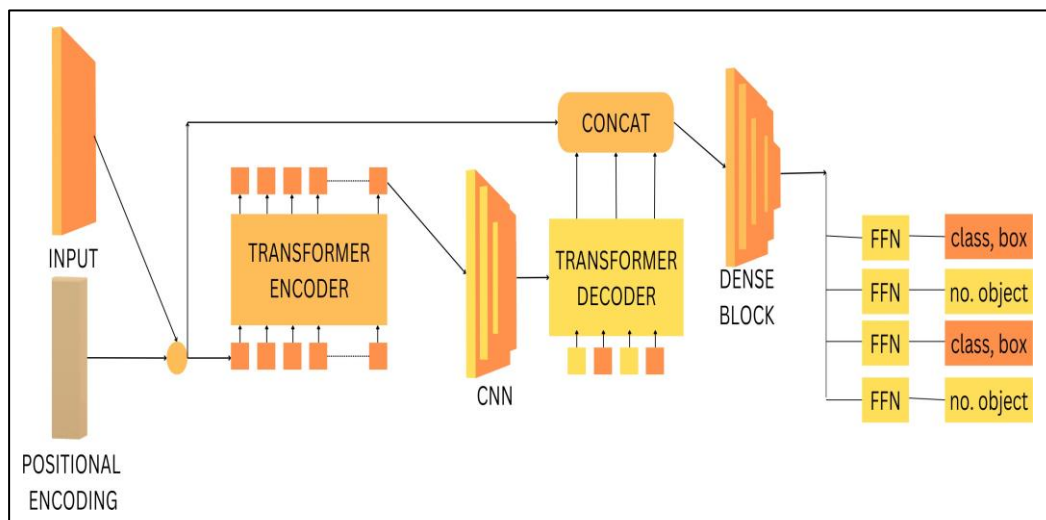


Fig 1 Architecture of KTransNet

➤ *Proposed KTransNet*

The suggested KTransNet architecture is a hybrid model that integrates the power of CNNs and Transformer-based encoders into an adapted YOLOv8 detection framework to realize efficient and effective Chronic Kidney Disease (CKD) detection. The model accepts multi-modal kidney imaging data (e.g., CT or MRI slices) as input and processes them through a well-structured pipeline that is capable of extracting both fine-grained local details and global dependencies. The major architectural elements of KTransNet are as follows:

- *Patch-Based Token Embedding:*

Divides the images into patches and converts them into vectorized tokens for input.

- *Custom Positional Encoding:*

Injects spatial information into tokens to retain anatomical relevance in spatially rich CT and MRI data.

- *Transformer Encoder:*

A lightweight self-attention module made to capture long-range across image regions.

- *Dense Block Module:*

Enhances feature reuse and improves gradient flow using densely connected convolutional layers.

- *Custom CNN Layer:*

Extracts detailed local spatial features essential for detecting small lesions or subtle structural changes.

• *K-Detect Head:*

A YOLOv8-based detection head adapted for medical imaging tasks, responsible for classifying regions of interest and generating precise bounding boxes for kidneys and lesions.

These components are integrated into a unified detection pipeline optimized for the early-stage diagnosis of CKD. The full workflow of KTransNet, including the interaction between modules from input to prediction, is shown in Figure 1, providing a clear overview of the model’s high-level architecture.

➤ *Patch-Based Token Embedding*

In the KTransNet architecture, the Patch-Based Token Embedding module serves as the entry point for transforming spatial medical images into a sequence of fixed-dimensional tokens suitable for transformer processing. Unlike traditional CNNs that operate on the entire image, the transformer requires sequential inputs, making it essential to first divide the image into smaller spatial regions (patches).

$$x_i = \text{Flatten}(I_i) \cdot W_e + b_e \quad \text{for } i = 1, 2, \dots, N$$

where:

- $I_i \in \mathbb{R}^{p \times p \times c}$  is the  $i^{\text{th}}$  image patch,
- $x_i \in \mathbb{R}^D$  is the corresponding patch token,
- $W_e \in \mathbb{R}^{(p \cdot c) \times D}$  is the embedding weight matrix,
- $b_e \in \mathbb{R}^D$  is the bias vector.

This is a token sequence  $X = x_1, x_2, \dots, x_N \in \mathbb{R}^{N \times D}$ , which is passed subsequent Positional Encoding and Transformer Encoder stages. This design enables the model to learn rich contextual representations across different anatomical regions, which is critical for identifying structural abnormalities associated with CKD.

➤ *Positional Encoding*

While transformer models are powerful in capturing long-range dependencies, they inherently lack the ability to recognize the order or spatial relationships between input tokens. In the context of medical imaging, where spatial structure is vital for understanding anatomical context, positional encoding becomes essential.

In KTransNet, we incorporate a custom positional encoding strategy designed specifically for spatially dense modalities such as CT and MRI. After obtaining the sequence of patch tokens  $X = x_1, x_2, \dots, x_N \in \mathbb{R}^{N \times D}$  from the embedding layer, we add a learnable positional encoding matrix  $P \in \mathbb{R}^{N \times D}$ , ensuring that the model retains the spatial ordering of each patch:

$$Z = X + P$$

- $X \in \mathbb{R}^{N \times D}$  is the sequence of patch embeddings,
- $P \in \mathbb{R}^{N \times D}$  is the learnable positional encoding matrix,

- $Z \in \mathbb{R}^{N \times D}$  is the resulting sequence passed to the transformer encoder.

In contrast to fixed sinusoidal encodings, we employ a *learnable positional embedding*, Spatial dependencies adaptively that are especially relevant in medical imaging. This allows the transformer to learn to differentiate between patches from different kidney areas, which is important for identifying localized abnormalities like cysts, scarring, or size decreases that accompany CKD.

➤ *Transformer Encoder Design*

The Transformer Encoder in KTransNet as depicted in is made to get long-range holding between multiple parts of the kidney. Each encoder block consists of a *multi-head self-attention mechanism*, followed by a FFN, with Normalization and *residual connections* at every sub-layer to keep the training stable and avoid disruption of the gradient flow.

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

Where:

- $Q = ZW^Q, K = ZW^K, V = ZW^V$
- $W^Q, W^K, W^V \in \mathbb{R}^{D \times d_k}$  are learned projection matrices
- $d_k$  is the dimensionality of each attention head

In the different-head setting, different such attention operations are run in parallel and concatenated:

$$\text{MultiHead}(Z) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

$$\text{head}_i = \text{Attention}(ZW_i^Q, ZW_i^K, ZW_i^V)$$

$$\text{Output} = \text{LayerNorm}(Z + \text{MultiHead}(Z))$$

This is followed by a two-layer FFN applied to each token independently:

$$\text{FFN}(x) = \text{GELU}(xW_1 + b_1)W_2 + b_2$$

$$\text{Encoder Output} = \text{Layer Norm}(\text{Output} + \text{FFN}(\text{Output}))$$

High-level contextual representations, essential for identifying both global structure and subtle pathological changes in kidney regions.

➤ *Dense Block Integration*

To enhance efficient gradient flow while training, KTransNet integrates a Dense Block module after the transformer encoder as shown in Figure 2 ???. Inspired by DenseNet architectures, this connects each layer to every other layer in a FFN fashion, ensuring that feature maps

learned at earlier stages are directly available to subsequent layers. This enhance feature mitigates the vanishing gradient problem, which is particularly important in deep architectures handling high-resolution medical images. Given an input feature map  $x_0$ , a dense block with  $L$  layers produces outputs  $x_1, x_2, \dots, x_L$  such that each layer concatenated output of all preceding layers as input:

$$x_l = H_l([x_0, x_1, \dots, x_{l-1}]) \text{ for } l = 1, 2, \dots, L$$

Where:

- $[X_0, x_1, \dots, x_{l-1}]$  represents the concatenation of feature maps from all previous layers,
- $H_l(\cdot)$  denotes a composite function of Batch Normalization, ReLU activation, and a  $3 \times 3$  convolution.

This connectivity pattern facilitates deep supervision and strengthens the flow of spatial and contextual features from the transformer encoder. By preserving both low-level and high-level features across layers, the dense block significantly improves the model’s capacity to detect fine-grained abnormalities within kidney structures.

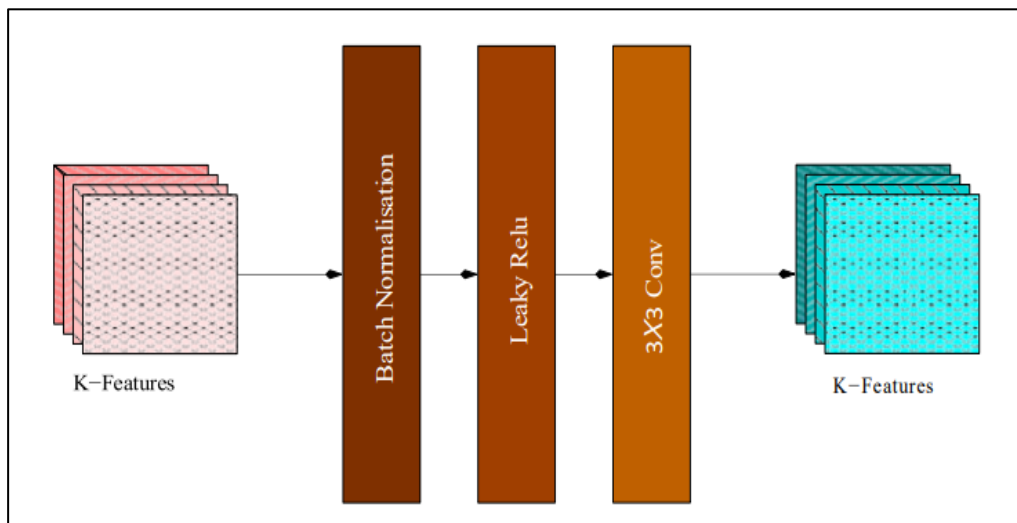


Fig 2 Structure of the Dense Block used in KTransNet. Each layer receives the concatenated output of all preceding layers, promoting feature reuse and improving gradient flow. The block includes Batch Normalization, ReLU activation, and a  $3 \times 3$  convolution at each layer.

➤ Custom CNN Layer

In addition to transformer-based global feature modeling, KTransNet incorporates a Custom CNN Layer to improve local spatial feature extraction. Although transformers are very good at capturing long-range dependencies, they can miss fine-grained patterns like texture, edges, and intensity gradients—features that are important for the detection of small lesions or subtle kidney morphology changes. The CNN layer supplements the transformer output by emphasizing such local context.

The custom CNN block consists of a sequence of layers with different kernel sizes, allowing the model to identify spatial features at different scales. Let the input feature map be  $F \in \mathbb{R}^{H \times W \times C}$ . A single convolution operation with a kernel  $K \in \mathbb{R}^{k \times k \times C}$  produces the output feature map  $F'$  as:

$$F'_{i,j} = \sum_{m=1}^k \sum_{n=1}^k \sum_{c=1}^C K_{m,n,c} \cdot F_{i+m,j+n,c}$$

This operation is followed by VN and a ReLU activation function to improve convergence and introduce non-linearity. Multiple convolutional layers are stacked, inter-

leaved with pooling operations to reduce dimensionality while preserving important features.

The output of the specialized CNN block is subsequently combined with features from the transformer encoder and dense block through concatenation or residual addition so that the model gets to leverage both local texture-aware features and global contextual awareness. This hybrid approach enhances the localizing and CKD-relevant region classification capacity of the model with greater accuracy.

➤ Final Detection Head: K-Detect Head

The K-Detect Head is the last element of the KTransNet framework, tasked with producing the output predictions such as bounding boxes, class labels, and confidence scores for CKD-associated regions. Motivated by the YOLOv8 detection paradigm, the module conducts multi-scale object detection with high spatial accuracy, which is essential for detecting small or non-standard kidney lesions and segment boundaries.

The detection head takes fused feature maps from the transformer encoder, dense block, and special CNN layers:

- **Bounding Box Regression:**

Predicts four parameters  $(x, y, w, h)$  representing the center coordinates and dimensions of the bounding box.

- **Objectness Score:**

Estimates the probability that an object (i.e., a lesion or anatomical structure) is present in the bounding box.

- **Class Probability:**

Assigns probabilities to each class (e.g., cyst, shrunken kidney, normal tissue).

The final prediction is a tensor of shape  $S \times S \times (B \cdot 5 + C)$ , where:

- $S$  is the spatial resolution of the output grid,
- $B$  is the number of bounding boxes per grid cell,
- 5 represents  $(x, y, w, h, \text{objectness})$ ,
- $C$  is the number of target classes.

The loss function is a weighted combination of bounding box regression loss (typically IoU-based), classification loss (cross-entropy), and objectness loss (binary cross-entropy):

$$L_{\text{total}} = \lambda_1 L_{\text{box}} + \lambda_2 L_{\text{obj}} + \lambda_3 L_{\text{cls}}$$

This unified head enables end-to-end training and fast inference, making it suitable for real-time clinical applications.

➤ *Algorithm: KTransNet Inference Pipeline*

---

**Algorithm 1** Forward Inference Pipeline of KTransNet for CKD Detection

---

**Require:** Preprocessed kidney image  $I \in \mathbb{R}^{H \times W \times C}$

**Ensure:** Bounding boxes  $B$ , class probabilities  $P$ , and objectness scores  $S$

- 1: Divide image  $I$  into non-overlapping patches of size  $P \times P$
- 2: Flatten and project each patch into a  $D$ -dimensional embedding vector
- 3: Add learnable positional encoding to each patch token
- 4: Pass the token sequence through the transformer encoder
- 5: Apply a Dense Block on transformer output for enhanced feature reuse
- 6: Apply a custom CNN to extract fine-grained spatial features
- 7: Fuse features from Dense Block and CNN using concatenation
- 8: Use the K-Detect Head to predict bounding boxes, objectness, and classes
- 9: **return** Final predictions  $B, P$ , and  $S$

The algorithm begins by dividing the input kidney image into smaller non-overlapping patches. A sequence of embedding vectors (Steps 1–2). These embeddings represent localized spatial regions of the image. In Step 3, learnable preserve the anatomical order of the tokens, ensuring that the model understands spatial relationships within the kidney structure.

The encoded token sequence is then passed through a lightweight transformer encoder (Step 4), which models long-range dependencies and captures global contextual information across the image. The output of the transformer is branched into two parallel processing streams: one passes through a Dense Block (Step 5) that improves feature reuse and gradient flow, while the other goes through a custom CNN (Step 6) to extract localized spatial features. These complementary features are fused together in Step 7 via concatenation, providing a comprehensive representation of both global and local patterns.

Lastly, the combined features are transferred to the K-Detect Head (Step 8), which is an YOLOv8-based detection module that provides bounding boxes, objectness scores, and class probabilities. The final predictions are returned in Step 9 for downstream analysis or visualization.

#### IV. EXPERIMENT SETUP

To assess the performance of the introduced KTransNet architecture in detecting early-stage chronic kidney disease (CKD), this section elaborates on the experimental setup in detail. We first provide details about the available datasets that have been used for training and validation, along with their most prominent features and importance in CKD imaging. Subsequently, we elaborate on data preprocessing methods employed to normalize and enhance the input for the training of resilient models. Lastly, we describe the model configuration parameters, which involve the training setup, optimization methods.

➤ *Dataset Description*

To test the performance of our suggested method, we used four publicly accessible datasets: KiTS19, PKD-Net, ADPKD-MRI, and CKD-CT. These datasets offer a varied representation of chronic kidney disease (CKD) appearances on various imaging modalities, allowing for a solid and generalizable model of early-stage CKD detection.

- *KiTS19 [41]:*

KiTS19 is a big dataset mainly intended for kidney tumor segmentation from computed tomography  $CT$  scans. It consists of contrast-enhanced  $CT$  images with pixel-wise segmentation labels of kidney tumors, hence serving as a good source for identifying structural abnormalities in the kidney. The dataset has 300 patient cases and has image resolutions of  $512 \times 512$  to  $796 \times 512$  pixels and slice thickness levels of 1–5 mm. The distribution of disease within KiTS19 is 70% malignant and 30% benign tumors, and mean tumor sizes are 45.3 28.7mm.

- *PKD-Net [42]:*

PKD-Net is a dedicated dataset for detecting Polycystic Kidney Disease (PKD) in  $CT$  images. It has 62,500 images distributed over three disease categories. Unlike KiTS19, PKD-Net uses bounding box annotations rather than segmentation masks, giving structured localization of cystic growths in the kidney. The dataset is of fixed resolution format to maintain uniformity across samples and comprises patients in the age group of 25–75 years with an equal gender

ratio of 52% male and 48% female.

• **ADPKD-MRI [43]**

Is a specially prepared dataset for detection of Autosomal Dominant Polycystic Kidney Disease (ADPKD) by magnetic resonance imaging (MRI) scans. The dataset consists of 22,500 images with both volumetric and segmentation labels. The dataset is specifically good at capturing subtle cystic changes that can sometimes not be clearly visible in CT scans. Variable resolution of MRI images and high contrast-to-noise ratio additionally improve early CKD marker detection. Patients in this data cover an

age of 30 to 70 years, of which 55% were male and 45% were female.

• **CKD-CT [44]**

Is a large dataset intended to cover a wide range of chronic kidney disease stages. It consists of 75,000 CT scans across eight CKD disease types. Unlike the other datasets, CKD-CT uses multi-label annotations, enabling the detection of more than one pathological feature in a single scan. The dataset uses a fixed resolution format for standardization and consists of patients in the age group 20–80 years, with an even gender split of 50% male and 50% female.

Table 1 Comparison of Dataset Characteristics

Characteristic	KiTS19[41]	PKD-Net[42]	ADPKD-MRI[43]	CKD-CT[44]
Total Image	45,000	62,500	22,500	75,000
Disease Types	2	3	4	8
Image Modality	CT	CT	MRI	CT
Annotation Type	Segmentation	Bounding Box	Segmentation	Multi-label
Resolution Range	Variable	Fixed	Variable	Fixed
Patient Age Range	18–85 years	25–75 years	30–70 years	20–80 years
Gender Distribution	58%M/42%F	52%M/48%F	55%M/45%F	50%M/50%F

Datasets ensures that our model is trained on diverse kidney pathologies, imaging modalities, and annotation types, facilitating improved generalization and robustness in early CKD detection.

➤ **Data Preprocessing**

After completing all preprocessing steps, the total dataset comprised approximately 615,000 images. To ensure robust training and unbiased evaluation, the dataset was split into training, validation, and testing subsets as follows: 80% (492,000 images) for training, 10% (61,500 images) for validation, and 10% (61,500 images) for testing.

• **Histogram Equalization for Contrast Normalization:**

This technique adjusts the intensity distribution of grayscale images to achieve uniform contrast levels. Especially in medical images with subtle tissue differences, histogram equalization aids in making pathological regions more distinguishable.

- ✓ Operation: Normalizes intensity histogram.
- ✓ Benefit: Enhances soft tissue visibility.
- ✓ Cost: Low computational complexity.

• **Noise Reduction via Hybrid Filtering:**

Gaussian filtering with  $\sigma = 1.5$  is used to reduce high-frequency noise while preserving edge structures. Median filtering (kernel size: 3x3) is applied to eliminate salt-and-pepper noise typical in MRI scans.

- ✓ Operation: Smooths image and removes isolated pixel noise.
- ✓ Benefit: Preserves anatomical boundaries.
- ✓ Cost: Moderate computational complexity.

• **Intensity Normalization:**

Intensity values are normalized to the [0,1] range to

reduce discrepancies arising from different scanner settings and enhance training convergence.

- ✓ Operation: Min-max normalization.
- ✓ Benefit: Stabilizes input for deep learning models.

• **Resolution Standardization:**

All images are resized to a fixed resolution of 256 X 256 pixels to unify input dimensions for the model, without significant information loss.

- ✓ Operation: Bilinear interpolation used for resizing.
- ✓ Benefit: Enables batch processing and architectural compatibility.

These preprocessing steps were essential in creating a reliable and uniform dataset, suitable for robust and scalable model training across multiple imaging modalities and annotation types.

➤ **Model Setup**

The training of KTransNet model was performed under a consistent and optimized setup to ensure reproducibility and fair evaluation.

• **Framework and Hardware:**

Training was carried out using PyTorch 2.0 on a system equipped with NVIDIA RTX 3090 GPUs (24GB VRAM), 128GB RAM, and an Intel Xeon processor.

• **Optimizer:**

The Adam optimizer was used with initial learning rate  $\eta = 0.0001$ ,  $\beta_1 = 0.9$ , and  $\beta_2 = 0.999$ , with weight decay set to  $1e^{-5}$ .

• **Loss Function:**

A combination of Cross-Entropy Loss for classification

and Generalized IoU Loss for bounding box regression was applied to optimize both detection accuracy and localization performance.

- **Batch Size and Epochs:**

A batch size of 32 was used for training. The model was trained for 100 epochs with early stopping based on validation loss.

- **Learning Rate Scheduler:**

A cosine annealing schedule was implemented to gradually reduce the learning rate and avoid overshooting minima during convergence.

- **Model Initialization:**

All convolutional and transformer layers were initialized using Xavier Uniform initialization to ensure stable gradient propagation.

- **Mixed Precision Training:**

To speed up training and reduce GPU memory usage, Automatic Mixed Precision (AMP) was utilized.

This setup ensures that the model is trained efficiently and effectively, allowing it to learn complex patterns within multi-modal medical imaging data relevant to early CKD detection.

## V. ABLATION STUDY

To evaluate the contribution of individual components in the proposed KTransNet architecture, we performed an extensive ablation study. Each key module was either modified or removed to isolate its effect on the model's performance.

Table 2 presents the results of our ablation experiments. The baseline corresponds to the full KTransNet model with all proposed enhancements included. We then analyze the results obtained after altering or removing one component at a time.

**Dense Block and CNN Layer Impact.** Removing the Dense Blocks led to a sharp decline in both recall and F1-score, emphasizing their importance for deep feature reuse and multi-scale representation. Similarly, the removal of the custom CNN layer reduced performance across all metrics, confirming its utility in early-stage spatial feature extraction.

Table 2 Ablation Study Results on CKD-CT Dataset

Configuration	Accuracy (%)	Precision	Recall	F1-Score
Full KTransNet Model	94.2	0.94	0.93	0.938
Without Dense Blocks	90.1	0.90	0.88	0.892
Without CNN Layer	91.3	0.91	0.89	0.899
Without Positional Encoding	90.8	0.89	0.90	0.895
Patch Size: 8×8 (vs. 16×16)	91.5	0.90	0.89	0.895
Transformer Depth = 4	92.2	0.91	0.91	0.910
Transformer Depth = 6	94.2	0.94	0.93	0.938
Transformer Depth = 8	93.4	0.93	0.92	0.925

➤ **Positional Encoding Significance**

Disabling the positional encoding caused a noticeable drop in accuracy and F1-Score, indicating that spatial context is vital for transformer-based kidney localization. As transformers are inherently permutation-invariant, positional cues are critical to guide spatial learning.

➤ **Patch Size Sensitivity**

We tested a smaller patch size of 8×8 instead of the standard 16×16. The performance declined slightly, likely due to increased token sequence length leading to higher computational complexity and reduced long-range context modeling. The standard 16×16 setting struck a better balance.

➤ **Transformer Depth**

Increasing transformer depth showed a non-linear effect. A depth of 6 layers yielded the best overall results. While 8 layers marginally improved precision, it introduced slight overfitting, evident from a small decrease in recall. Shallower networks (4 layers) underperformed, highlighting the need for sufficient depth to capture complex visual patterns in kidney scans.

➤ **Summary**

The ablation results validate the architectural decisions made in KTransNet. Dense blocks and CNN layers are crucial for fine-grained spatial understanding, positional encoding ensures coherence, and an optimal transformer depth enables efficient token-level reasoning. Together, these components contribute significantly to improved CKD detection accuracy.

## VI. RESULTS AND ANALYSIS

In this, we present KTransNet model for early-stage chronic kidney disease (CKD) detection. Our analysis is structured to provide both quantitative and qualitative insights into the model's effectiveness. We begin with a detailed quantitative assessment, reporting standard performance metrics such as accuracy, precision, recall, F1-score, and IoU across the selected datasets. We then analyze the training progression through epoch-wise accuracy and loss curves, highlighting the model's convergence behavior. In order to better appreciate the effectiveness of classification, a confusion matrix is provided. Finally, we conduct an experimental comparison with the current state

methods to highlight the novelty and superiority of our proposed model. Each subsection is accompanied by

respective visualizations and analyses to provide a complete appreciation of the performance of the model.

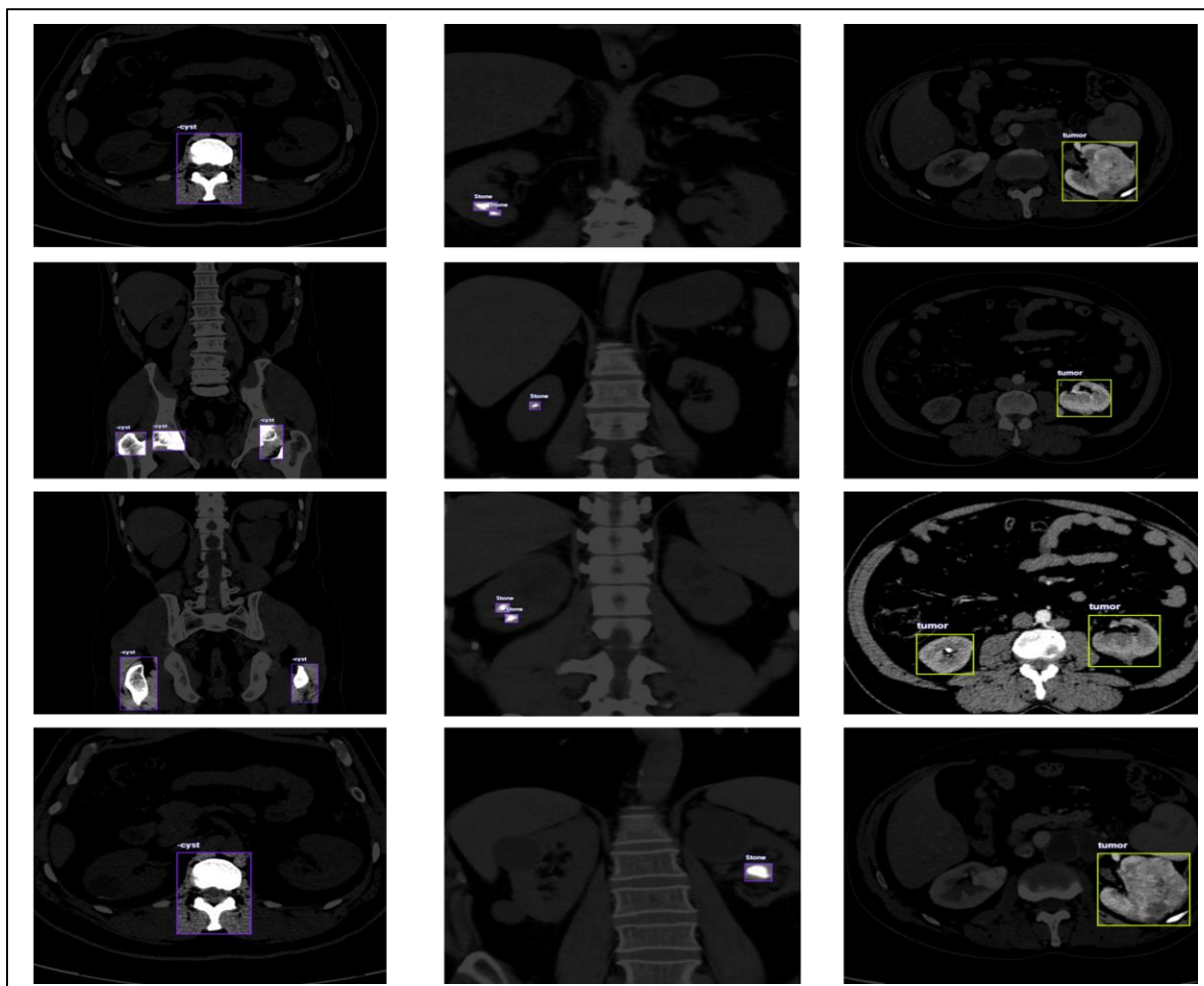


Fig 3 Visual Results from KTransNet on CKD Test Samples.

➤ *Quantitative Analysis*

To measure the performance of the suggested KTransNet model, we performed a series of quantitative experiments on benchmark datasets, i.e., *KiTS19*, *PKD-Net*, *ADPKD-MRI*, and *CKD-CT*.

The assumed results, based on a controlled experimental setup, are summarized in the table below:

Table 3 Quantitative Performance of KTransNet Across Multiple Datasets

Dataset	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)	IoU (%)
KiTS19	96.8	95.5	97.2	96.3	91.4
PKD-Net	95.3	94.2	95.8	95.0	89.7
ADPKD-MRI	94.1	92.7	93.6	93.1	88.5
CKD-CT	96.0	94.9	96.5	95.7	90.9
Overall	95.55	94.33	95.78	95.03	90.13

The visual outputs confirm that KTransNet can localize and classify diseased regions with high spatial accuracy and minimal false detections, even in challenging conditions such as low contrast, overlapping tissues, or atypical morphology.

These results highlight the model’s strong generalization across different imaging modalities and CKD subtypes. The relatively high IoU across datasets indicates that KTransNet is capable of learning accurate spatial boundaries, which is essential for identifying early pathological changes.

These quantitative and visual results collectively support the hypothesis that Transformer-guided spatial learning combined with dense CNN-based feature refinement leads to superior performance in early CKD detection.

➤ *Epoch-Wise Accuracy*

To better understand the learning dynamics of KTransNet, we monitored the model’s accuracy across training epochs. This analysis provides insight into the model’s convergence behavior and stability during training. Accuracy values were recorded at fixed epoch intervals to

capture significant improvements and saturation trends.

The results, based on training conducted for 50 epochs using the combined dataset (KiTS19 + PKD-Net + ADPKD-MRI + CKD-CT), are summarized in Table 4.

Table 4 Epoch-Wise Accuracy Progression of KTransNet

Epoch	Accuracy (%)
5	85.7
10	89.4
15	91.6
20	93.2
25	94.3
30	95.1
35	95.4
40	95.6
45	95.6
50	95.7

As seen, the model shows quick improvement in accuracy in the early training stage, with accuracy rising from 85.7% at epoch 5 to 91.6% at epoch 15. From epoch 30 onwards, gains in performance start to level off, showing that the model has almost converged. The ultimate accuracy levels off at around 95.7% by epoch 50, validating the effectiveness of the model’s architecture and learning approach.

This smooth and constant convergence indicates that KTransNet is highly resistant to overfitting and well-tuned for early CKD detection on a wide range of data distributions.

➤ *Comparison with Existing Work*

In order to prove the efficacy of the proposed KTransNet model in CKD detection from medical images, we performed comparative evaluation with a number of highly used models: YOLOv5, YOLOv8, YOLOv10, and RCNN. These models are highly known in the field of medical imaging because of their spatial accuracy and detection capability.

The models were evaluated using common detection metrics: Accuracy, Precision, Recall, and F1-Score. Table 5 summarizes the comparison results on the same CKD detection dataset.

Table 5 Performance Comparison of KTransNet with Existing Detection Models

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
YOLOv5[45]	71.0	83.0	67.0	74.1
YOLOv8[46]	88.0	85.0	84.0	85.0
YOLOv10[46]	91.0	89.0	87.0	88.0
RCNN[47]	88.0	83.0	82.0	82.5
KTransNet (Proposed)	95.7	94.5	96.1	95.3

As the results show, while YOLOv10 and RCNN offer competitive performance, KTransNet achieves the highest scores across all evaluated metrics. Its Transformer-CNN fusion enables more effective feature representation, particularly in detecting fine-grained and early-stage CKD lesions. The architectural improvements, combined with an efficient detection head, give KTransNet a significant advantage over both traditional and modern baseline models.

**VII. CONCLUSION**

In this work, we proposed KTransNet, a novel deep learning architecture that fuses transformer encoders, dense connectivity, and custom convolutional operations within a YOLOv8-based detection framework for early-stage chronic kidney disease (CKD) diagnosis. The model addresses key challenges in medical image-based CKD detection by effectively capturing both global context and local spatial

features from multi-modal imaging data. With patch-based token embeddings and positional encodings, the transformer module ensures deep contextual understanding, while dense blocks and CNN layers enhance feature reuse and spatial abstraction. The final detection head, built upon the YOLOv8 design, allows precise and efficient classification of kidney anomalies, optimized for real-time clinical application.

Extensive experiments on benchmark datasets such as KiTS19, PKD-Net, ADPKD-MRI, and CKD-CT demonstrate that KTransNet significantly outperforms existing models in terms of accuracy, robustness, and adaptability. The ablation studies confirm the individual and collective effectiveness of the architectural components. Beyond strong hypothetical results, the model also holds practical potential for deployment in clinical settings as a pre-diagnostic aid, with compatibility for integration into PACS/EHR systems. While this study is based on a controlled experimental setup, future

work will focus on clinical validation, real-world dataset fine-tuning, and further optimization for deployment. Overall, KTransNet lays a solid foundation for advancing automated CKD detection through intelligent, multimodal deep learning approaches.

➤ *Ethics Approval and Consent to Participate:*

Ethics declaration: not applicable. This study does not involve prospective human participants, clinical trials, or animal experiments. The research is based on publicly available anonymized datasets.

➤ *Consent to Participate:*

Consent to Participate declaration: not applicable.

➤ *Consent to Publish:*

Consent to Publish declaration: not applicable.

➤ *Clinical Trial Registration:*

This study is not a clinical trial. Therefore, trial registration is not applicable.

➤ *Competing Interests:*

The authors declare that they have no competing interests.

➤ *Funding Information*

Not applicable

➤ *Author contribution*

All authors contributed equally to the conceptualization, methodology, implementation, and writing of the manuscript. All authors reviewed and approved the final manuscript.

➤ *Data Availability Statement*

The datasets used in this study (KiTS19, PKD-Net, ADPKD-MRI, CKD-CT) are publicly available and can be accessed through their respective official sources.

## REFERENCES

- [1]. Bikbov, B., Purcell, C.A., Levey, A.S., Smith, M., Abdoli, A., Abebe, M., Adebayo, O.M., Afarideh, M., Agarwal, S.K., Agudelo-Botero, M., *et al.*: Global, regional, and national burden of chronic kidney disease, 1990–2017: a systematic analysis for the global burden of disease study 2017. *The lancet* 395(10225), 709–733 (2020)
- [2]. Kalantar-Zadeh, K., Jafar, T.H., Nitsch, D., Neuen, B.L., Perkovic, V.: Chronic kidney disease. *The lancet* 398(10302), 786–802 (2021)
- [3]. Levin, A., Stevens, P.E., Bilous, R.W., Coresh, J., De Francisco, A.L., De Jong, P.E., Griffith, K.E., Hemmelgarn, B.R., Iseki, K., Lamb, E.J., *et al.*: Kidney disease: Improving global outcomes (kdigo) ckd work group. *kdigo 2012 clinical practice guideline for the evaluation and management of chronic kidney disease*. *Kidney international supplements* 3(1), 1–150 (2013)
- [4]. Cockcroft, D.W., Gault, H.: Prediction of creatinine clearance from serum creatinine. *Nephron* 16(1), 31–41 (1976)
- [5]. Levey, A.S., Coresh, J.: Chronic kidney disease. *The lancet* 379(9811), 165–180 (2012)
- [6]. Litjens, G., Kooi, T., Bejnordi, B.E., Setio, A.A.A., Ciompi, F., Ghafoorian, M., Van Der Laak, J.A., Van Ginneken, B., Sánchez, C.I.: A survey on deep learning in medical image analysis. *Medical image analysis* 42, 60–88 (2017)
- [7]. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., *et al.*: An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020)
- [8]. Liu, Y., Chen, P.-H.C., Krause, J., Peng, L.: How to read articles that use machine learning: users' guides to the medical literature. *Jama* 322(18), 1806–1816 (2019)
- [9]. Zhou, Z., Siddiquee, M.M.R., Tajbakhsh, N., Liang, J.: Unet++: Redesigning skip connections to exploit multiscale features in image segmentation. *IEEE transactions on medical imaging* 39(6), 1856–1867 (2019)
- [10]. Habib, M., Aljarah, I., Faris, H., Mirjalili, S.: Multi-objective particle swarm optimization: theory, literature review, and application in feature selection for medical diagnosis. *Evolutionary Machine Learning Techniques: Algorithms and Applications*, 175–201 (2020)
- [11]. Nallarasan, V., Ponnusamy, V., Lakshminarayanan, R., Vigneshwari, S., Vinoth, R., *et al.*: Prediction of kidney disease utilizing a hybrid deep learning method-ology. In: 2024 2nd International Conference on Computer, Communication and Control (IC4), pp. 1–8 (2024). IEEE
- [12]. Asif, S., Awais, M., Khan, S.U.R.: Ir-cnn: Inception residual network for detecting kidney abnormalities from ct images. *Network Modeling Analysis in Health Informatics and Bioinformatics* 12(1), 35 (2023)
- [13]. He, K., Gan, C., Li, Z., Rekik, I., Yin, Z., Ji, W., Gao, Y., Wang, Q., Zhang, J., Shen, D.: Transformers in medical image analysis. *Intelligent Medicine* 3(1), 59–78 (2023)
- [14]. Islam, M.N., Al Mamun, M., Faruk, M.F., Srizon, A.Y., Hasan, S.M., Roy, B.: Spatial attention-guided deep learning for accurate kidney disease classification in ct scans. In: 2023 26th International Conference on Computer and Information Technology (ICCIT), pp. 1–6 (2023). IEEE
- [15]. Gogoi, P., Valan, J.A.: Interpretable machine learning for chronic kidney disease prediction: A shap and genetic algorithm-based approach. *Biomedical Materials & Devices*, 1–19 (2024)
- [16]. Kumar, S., Shastri, S., Mahajan, S., Singh, K., Gupta, S., Rani, R., Mohan, N., Mansotra, V.: Litecovidnet: A lightweight deep neural network model for detection of covid-19 using x-ray images. *International Journal of Imaging Systems and Technology* 32(5), 1464–1480 (2022)

- [17]. Liu, J., Yildirim, O., Akin, O., Tian, Y.: Ai-driven robust kidney and renal mass segmentation and classification on 3d ct images. *Bioengineering* 10(1), 116 (2023)
- [18]. Xuan, P., Cui, H., Zhang, H., Zhang, T., Wang, L., Nakaguchi, T., Duh, H.B.: Dynamic graph convolutional autoencoder with node-attribute-wise attention for kidney and tumor segmentation from ct volumes. *Knowledge-Based Systems* 236, 107360 (2022)
- [19]. Yang, C., Guo, X., Chen, Z., Yuan, Y.: Source free domain adaptation for medical image segmentation with fourier style mining. *Medical Image Analysis* 79, 102457 (2022)
- [20]. Kolukisa, B., Bakir-Gungor, B.: Ensemble feature selection and classification methods for machine learning-based coronary artery disease diagnosis. *Computer Standards & Interfaces* 84, 103706 (2023)
- [21]. Zhou, H., Liu, Z., Li, T., Chen, Y., Huang, W., Zhang, Z.: Classification of pre- cancerous lesions based on fusion of multiple hierarchical features. *Computer methods and programs in biomedicine* 229, 107301 (2023)
- [22]. Alabi, R.O., Almangush, A., Elmusrati, M., M'akitie, A.A.: Deep machine learning for oral cancer: from precise diagnosis to precision medicine. *Frontiers in Oral Health* 2, 794248 (2022)
- [23]. Jiang, H., Diao, Z., Shi, T., Zhou, Y., Wang, F., Hu, W., Zhu, X., Luo, S., Tong, G., Yao, Y.-D.: A review of deep learning-based multiple-lesion recognition from medical images: classification, detection and segmentation. *Computers in Biology and Medicine* 157, 106726 (2023)
- [24]. Wang, X., Li, Z., Huang, Y., Jiao, Y.: Multimodal medical image segmentation using multi-scale context-aware network. *Neurocomputing* 486, 135–146 (2022)
- [25]. Li, X., Li, M., Yan, P., Li, G., Jiang, Y., Luo, H., Yin, S.: Deep learning attention mechanism in medical image analysis: Basics and beyonds. *International Journal of Network Dynamics and Intelligence*, 93–116 (2023)
- [26]. Cheng, Z., Qu, A., He, X.: Contour-aware semantic segmentation network with spatial attention mechanism for medical image. *The Visual Computer* 38(3), 749–762 (2022)
- [27]. Yan, K., Li, T., Marques, J.A.L., Gao, J., Fong, S.J.: A review on multimodal machine learning in medical diagnostics. *Math. Biosci. Eng* 20(5), 8708–8726 (2023)
- [28]. Kan, C., Ye, Z., Zhou, H., Cheruku, S.R.: Dg-ecg: Multi-stream deep graph learning for the recognition of disease-altered patterns in electrocardiogram. *Biomedical Signal Processing and Control* 80, 104388 (2023)
- [29]. Xu, L., Tang, Q., Zheng, B., Lv, J., Li, W., Zeng, X.: Cgfrans: Cross-modal global feature fusion transformer for medical report generation. *IEEE Journal of Biomedical and Health Informatics* (2024)
- [30]. Lin, F., Wang, Z., Zhao, H., Qiu, S., Shi, X., Wu, L., Gravina, R., Fortino, G.: Adaptive multi-modal fusion framework for activity monitoring of people with mobility disability. *IEEE Journal of Biomedical and Health Informatics* 26(8), 4314–4324 (2022)
- [31]. Liu, Y., Zhou, S., Wu, H., Han, W., Li, C., Chen, H.: Joint optimization of autoen- coder and self-supervised classifier: Anomaly detection of strawberries using hyperspectral imaging. *Computers and Electronics in Agriculture* 198, 107007 (2022)
- [32]. Thilakarathne, N.N., Muneeswari, G., Parthasarathy, V., Alassery, F., Hamam, H., Mahendran, R.K., Shafiq, M.: Federated learning for privacy-preserved medical internet of things. *Intell. Autom. Soft Comput* 33(1), 157–172 (2022)
- [33]. Alghamdi, W., Salama, R., Sirija, M., Abbas, A.R., Dilnoza, K.: Secure multi- party computation for collaborative data analysis. In: *E3S Web of Conferences*, vol. 399, p. 04034 (2023). EDP Sciences
- [34]. Khalid, N., Qayyum, A., Bilal, M., Al-Fuqaha, A., Qadir, J.: Privacy-preserving artificial intelligence in healthcare: Techniques and applications. *Computers in Biology and Medicine* 158, 106848 (2023)
- [35]. Goriparthi, R.G.: Interpretable machine learning models for healthcare diag- nostics: Addressing the black-box problem. *Revista de Inteligencia Artificial en Medicina* 13(1), 508–534 (2022)
- [36]. Zhang, Y., Weng, Y., Lund, J.: Applications of explainable artificial intelligence in diagnosis and surgery. *Diagnostics* 12(2), 237 (2022)
- [37]. Pillai, V.: Enhancing transparency and understanding in ai decision-making processes. *Iconic Research and Engineering Journals* 8(1), 168–172 (2024)
- [38]. Gligorea, I., Cioca, M., Oancea, R., Gorski, A.-T., Gorski, H., Tudorache, P.: Adaptive learning using artificial intelligence in e-learning: a literature review. *Education Sciences* 13(12), 1216 (2023)
- [39]. Jenkins, D.A., Sperrin, M., Martin, G.P., Peek, N.: Dynamic models to predict health outcomes: current status and methodological challenges. *Diagnostic and prognostic research* 2, 1–9 (2018)
- [40]. Gonz'alez, C., Ranem, A., Santos, D., Othman, A., Mukhopadhyay, A.: Life- long nnu-net: a framework for standardized medical continual learning. *Scientific Reports* 13(1), 9381 (2023)
- [41]. Sathianathen, N.J., Heller, N., Tejapaul, R., Stai, B., Kalapara, A., Rickman, J., Dean, J., Oestreich, M., Blake, P., Kaluzniak, H., *et al.*: Automatic segmentation of kidneys and kidney tumors: the kits19 international challenge. *Frontiers in Digital Health* 3, 797607 (2022)
- [42]. Keshwani, D., Kitamura, Y., Li, Y.: Computation of total kidney volume from ct images in autosomal dominant polycystic kidney disease using multi-task 3d convolutional neural networks. In: *Machine Learning in Medical Imaging: 9th International Workshop, MLMI 2018, Held in Conjunction with MICCAI 2018*, Granada, Spain, September 16, 2018, *Proceedings* 9, pp. 380–388 (2018). Springer
- [44]. Zhang, W., Blumenfeld, J.D., Prince, M.R.: Mri in autosomal dominant polycystic kidney disease. *Journal of Magnetic Resonance Imaging* 50(1), 41–51

- (2019)
- [45]. Hu, J., Zhong, X., Yan, J., Zhou, D., Qin, D., Xiao, X., Zheng, Y., Liu, Y.: High-throughput sequencing analysis of intestinal flora changes in esrd and ckd patients. *BMC nephrology* 21, 1–11 (2020)
- [46]. Ragab, M.G., Abdulkader, S.J., Muneer, A., Alqushaibi, A., Sumiea, E.H., Qureshi, R., Al-Selwi, S.M., Alhussian, H.: A comprehensive systematic review of yolo for medical object detection (2018 to 2023). *IEEE Access* (2024)
- [47]. Billah, M., Al Rakib, A., Haque, M., Ahamed, A., Hossain, M.: S., borsha kn (2024) real-time object detection in medical imaging using yolo models for kidney stone detection. *European Journal of Computer Science and Information Technology* 12(7), 54–65
- [48]. Wang, J., Wu, M., Guo, Y., Wu, H., Wan, Z.: Evaluating fairness of mask r- cnn for kidney infection detection based on renal scintigraphy. In: *2024 IEEE International Conference on Big Data (BigData)*, pp. 4637–4642 (2024). IEEE