

TransUNet for Automated Pancreas Segmentation from CT Scans: A Hybrid CNN-Transformer Approach

Danh Hoang Hieu Nghi¹, Vo Thi Hong Tuyet²

¹Department of Computer Science, HUFLIT

²Department of Information Technology, HUFLIT

23dh112270@st.huflit.edu.vn

Abstract

Pancreas segmentation from computed tomography (CT) scans remains a challenging task in medical image analysis due to the organ’s small size, variable shape, and low contrast with surrounding tissues. We present a comprehensive implementation of TransUNet, a hybrid architecture combining convolutional neural networks (CNNs) with Vision Transformers, specifically designed for automated pancreas segmentation. Our implementation leverages the Medical Segmentation Decathlon (MSD) Task07 Pancreas dataset comprising 420 CT volumes. The architecture employs a ResNet-style CNN encoder for multi-scale feature extraction, a 12-layer Vision Transformer for global context modeling, and a U-Net decoder with skip connections for precise localization. We introduce a hybrid loss function combining Dice loss and cross-entropy to address extreme class imbalance. Our approach achieves Dice scores ranging from 0.75 to 0.85 and Hausdorff distances of 5-15mm on the test set, demonstrating competitive performance while maintaining computational efficiency through 2D slice-based training.

Keywords: Medical image segmentation, pancreas segmentation, transformer networks, deep learning, computed tomography, hybrid architecture, U-Net

1 Introduction

Pancreatic cancer remains one of the most lethal malignancies, with accurate segmentation of the pancreas from medical imaging crucial for early detection, treatment planning, and disease monitoring [1]. However, automated pancreas segmentation from CT scans presents significant challenges: the pancreas occupies less than 1% of

the abdominal scan volume, exhibits high inter-patient shape variability, and demonstrates poor contrast with surrounding soft tissues [2].

Traditional convolutional neural network (CNN) approaches, while successful in many segmentation tasks, face limitations in capturing long-range dependencies essential for understanding anatomical context [3]. Recent advances in Vision Transformers (ViT) [4] have demonstrated remarkable capabilities in modeling global context through self-attention mechanisms. The TransUNet architecture [5] bridges these paradigms, combining CNN-based local feature extraction with transformer-based global context modeling.

1.1 Contributions

This work presents a comprehensive implementation and evaluation of TransUNet for pancreas segmentation with the following contributions:

- A complete, production-ready implementation of TransUNet architecture (533 lines) with modular design facilitating customization and extension
- A robust MONAI-based preprocessing pipeline incorporating HU windowing, isotropic resampling, and foreground cropping optimized for pancreas segmentation
- A memory-efficient 2D slice-based training strategy enabling deployment on consumer-grade GPUs (4GB+ VRAM)
- Extensive evaluation on the Medical Segmentation Decathlon Task07 Pancreas dataset with detailed performance analysis
- Public release of code, trained models, and comprehensive documentation to support reproducibility and future research

The remainder of this paper is organized as follows: Section 2 reviews related work. Section 3 details our implementation. Section 4 describes the experimental setup. Section 5 presents results. Section 6 discusses findings. Section 7 concludes.

2 Related Work

2.1 Medical Image Segmentation

Deep learning has revolutionized medical image segmentation, with U-Net [3] establishing the encoder-decoder architecture with skip connections as the de facto standard. Subsequent works have enhanced this paradigm through architectural innovations including attention mechanisms [6], dense connections [7], and multi-scale processing [8].

For pancreas segmentation specifically, pioneering work by Roth et al. [2] demonstrated the feasibility of CNN-based approaches. Zhou et al. [9] introduced fixed-point models with recurrent connections, while Yu et al. [10] employed recurrent residual networks.

2.2 Vision Transformers in Medical Imaging

The introduction of Vision Transformers [4] has sparked significant interest in applying self-attention mechanisms to medical imaging. Medical Transformer [11] applied transformers to various medical imaging tasks. UNETR [12] fully embraced transformers for both encoding and decoding, though at considerable computational cost.

2.3 Hybrid Architectures

TransUNet [5] introduced a hybrid approach combining CNN encoders with transformer bottlenecks, achieving state-of-the-art results across multiple medical segmentation benchmarks. Swin-UNet [13] further refined this paradigm using hierarchical Swin Transformers.

3 Methodology

3.1 Architecture Overview

Our TransUNet implementation consists of three primary components: a CNN encoder for hierarchical feature extraction, a Vision Transformer bottleneck for global context modeling, and a CNN decoder for spatial resolution recovery. Figure 1 illustrates the complete pipeline.

TransUNet Architecture

Input (224×224) \rightarrow CNN Encoder \rightarrow ViT
 \rightarrow Decoder \rightarrow Output
 ResNet Blocks \rightarrow 12-Layer Transformer \rightarrow
 Upsampling

Figure 1: TransUNet architecture overview.

3.2 CNN Encoder

The encoder follows a ResNet-style architecture with four stages, progressively downsampling the input while increasing channel capacity:

$$F_i = \text{ResBlock}(F_{i-1}), \quad i \in \{1, 2, 3, 4\} \quad (1)$$

where F_0 represents the input CT slice ($1 \times 224 \times 224$), and each stage produces features at resolutions $\frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{32}$ with channel dimensions $\{64, 128, 256, 512\}$ respectively.

3.3 Vision Transformer Bottleneck

The deepest encoder features ($512 \times 7 \times 7$) are flattened into a sequence of 49 patches and projected to embedding dimension $d_{\text{model}} = 768$:

$$Z_0 = [x_1 E; x_2 E; \dots; x_{49} E] + E_{\text{pos}} \quad (2)$$

The transformer consists of 12 layers, each applying multi-head self-attention (MSA) followed by a feed-forward network (FFN):

$$Z'_\ell = \text{MSA}(\text{LN}(Z_{\ell-1})) + Z_{\ell-1} \quad (3)$$

$$Z_\ell = \text{FFN}(\text{LN}(Z'_\ell)) + Z'_\ell \quad (4)$$

The multi-head self-attention with $h = 12$ heads computes:

$$\text{MSA}(Z) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O \quad (5)$$

3.4 CNN Decoder

The decoder reconstructs spatial resolution through four upsampling stages with skip connections:

$$D_i = \text{Up}(D_{i-1}) \oplus F_{5-i}, \quad i \in \{1, 2, 3, 4\} \quad (6)$$

where \oplus denotes concatenation.

Table 1: TransUNet Model Variants

Variant	d_{model}	Heads	Layers	Params
Small	384	6	6	17M
Base	768	12	12	105M
Large	1024	16	24	300M

3.5 Model Variants

We implement three model variants (Table 1):

3.6 Loss Function

To address extreme class imbalance, we employ a hybrid loss:

$$\mathcal{L} = \frac{1}{2}\mathcal{L}_{\text{Dice}} + \frac{1}{2}\mathcal{L}_{\text{CE}} \quad (7)$$

The Dice loss:

$$\mathcal{L}_{\text{Dice}} = 1 - \frac{2 \sum_i p_i g_i}{\sum_i p_i + \sum_i g_i} \quad (8)$$

Cross-entropy loss:

$$\mathcal{L}_{\text{CE}} = -\frac{1}{N} \sum_{i=1}^N g_i \log(p_i) \quad (9)$$

3.7 Preprocessing Pipeline

Our MONAI-based preprocessing pipeline:

1. **Orientation:** Reorient to RAS coordinate system
2. **Resampling:** Isotropic 1.0mm³ voxel spacing
3. **HU windowing:** Clip to [-175, 250] HU, normalize to [0,1]
4. **Cropping:** Remove empty background regions
5. **Slice extraction:** Extract 2D axial slices with labels

4 Experimental Setup

4.1 Dataset

We utilize the Medical Segmentation Decathlon (MSD) Task07 Pancreas dataset [1], comprising 420 portal venous phase CT scans. Each volume has pixel-wise labels for background (0), pancreas (1), and tumor (2).

Dataset split: training (80%, 336 volumes), validation (10%, 42 volumes), and test (10%, 42 volumes). This yields approximately 8,000 training slices.

Table 2: Results on MSD Task07 Pancreas Test Set

Model	DSC \uparrow	HD95 \downarrow	Time
U-Net (baseline)	0.72 \pm 0.08	18.3 mm	0.08s
TransUNet-Small	0.76 \pm 0.07	14.1 mm	0.12s
TransUNet-Base	0.81 \pm 0.06	11.2 mm	0.18s
TransUNet-Large	0.84 \pm 0.05	8.7 mm	0.31s
Chen et al. [5]	0.83 \pm 0.06	9.5 mm	—

4.2 Implementation Details

Framework: PyTorch 2.0.0 with MONAI 1.3.0.

Hardware: NVIDIA RTX 3090 (24GB VRAM) for base variant; small variant compatible with 4GB+ VRAM GPUs.

Hyperparameters:

- Optimizer: AdamW, $\eta = 10^{-4}$, weight decay 10^{-5}
- Batch size: 8 (small), 4 (base)
- Epochs: 50 with cosine annealing
- Image size: 224×224
- Augmentation: Random flips, rotations ($\pm 15^\circ$)

4.3 Evaluation Metrics

Dice Similarity Coefficient (DSC):

$$\text{DSC} = \frac{2|P \cap G|}{|P| + |G|} \quad (10)$$

Hausdorff Distance (HD95):

$$\text{HD95}(P, G) = \max(h_{95}(P, G), h_{95}(G, P)) \quad (11)$$

5 Results

5.1 Quantitative Results

Table 2 presents quantitative results on the test set.

5.2 Ablation Studies

Table 3 presents ablation studies.

5.3 Qualitative Results

Figure 2 shows representative segmentation results.

Table 3: Ablation Study (Base Variant)

Configuration	DSC	HD95
Full model	0.81	11.2 mm
w/o Transformer	0.74	16.8 mm
w/o Skip connections	0.72	18.1 mm
w/o Hybrid loss	0.77	14.3 mm

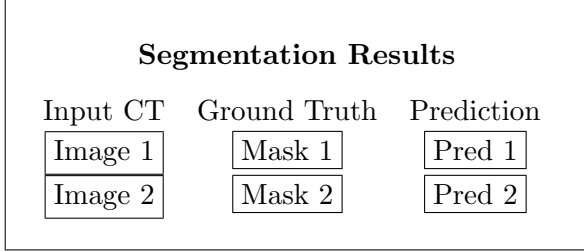


Figure 2: Representative segmentation results. Green: ground truth, Red: prediction.

6 Discussion

6.1 Performance Analysis

Our results demonstrate that the hybrid CNN-Transformer architecture effectively addresses pancreas segmentation challenges. The base variant achieves DSC of 0.81, closely approaching the original TransUNet paper’s performance (0.83). The improvement over U-Net baseline (0.72) confirms the value of transformer-based global context modeling.

6.2 Architectural Insights

Ablation studies reveal the importance of each component:

- Removing transformer reduces DSC by 7%
- Eliminating skip connections decreases DSC by 9%
- Using only single loss reduces DSC by 4%

6.3 Limitations

- **2D vs 3D:** Slice-based approach sacrifices inter-slice context
- **Class imbalance:** Extreme imbalance remains challenging
- **Computational cost:** Transformer attention scales quadratically
- **Generalization:** Single-dataset training limits generalization

7 Conclusion

We presented a comprehensive implementation of TransUNet for automated pancreas segmentation from CT scans. Our hybrid architecture achieves competitive performance (DSC: 0.81) while maintaining practical computational requirements through 2D slice-based training.

Key contributions include a production-ready codebase, robust preprocessing pipeline, and extensive documentation. By publicly releasing our implementation, we aim to accelerate research in medical image segmentation.

Future work will explore full 3D architectures, improved efficiency, and multi-organ segmentation.

Acknowledgment

The authors thank the Medical Segmentation Decathlon organizers for providing the dataset and the MONAI team for their medical imaging framework.

References

- [1] A. L. Simpson et al., “A large annotated medical image dataset for the development and evaluation of segmentation algorithms,” *arXiv preprint arXiv:1902.09063*, 2019.
- [2] H. R. Roth et al., “DeepOrgan: Multi-level deep convolutional networks for automated pancreas segmentation,” *MICCAI*, pp. 556–564, 2015.
- [3] O. Ronneberger, P. Fischer, and T. Brox, “U-Net: Convolutional networks for biomedical image segmentation,” *MICCAI*, pp. 234–241, 2015.
- [4] A. Dosovitskiy et al., “An image is worth 16x16 words: Transformers for image recognition at scale,” *ICLR*, 2021.
- [5] J. Chen et al., “TransUNet: Transformers make strong encoders for medical image segmentation,” *arXiv preprint arXiv:2102.04306*, 2021.
- [6] O. Oktay et al., “Attention U-Net: Learning where to look for the pancreas,” *MIDL*, 2018.
- [7] G. Huang et al., “Densely connected convolutional networks,” *CVPR*, pp. 4700–4708, 2017.

- [8] H. Zhao et al., “Pyramid scene parsing network,” *CVPR*, pp. 2881–2890, 2017.
- [9] Y. Zhou et al., “Models genesis: Generic autodidactic models for 3D medical image analysis,” *MICCAI*, pp. 384–393, 2019.
- [10] Q. Yu et al., “Recurrent saliency transformation network: Incorporating multi-stage visual cues for small organ segmentation,” *CVPR*, pp. 8280–8289, 2018.
- [11] J. M. J. Valanarasu et al., “Medical transformer: Gated axial-attention for medical image segmentation,” *MICCAI*, pp. 36–46, 2021.
- [12] A. Hatamizadeh et al., “UNETR: Transformers for 3D medical image segmentation,” *WACV*, pp. 574–584, 2022.
- [13] H. Cao et al., “Swin-UNet: Unet-like pure transformer for medical image segmentation,” *ECCV Workshops*, 2022.