

Deep Learning-Based Hybrid CNN Architecture for Automated Medical Image Segmentation

Sampada Thigale^{1*}, Ramesh Y. Mali²

¹Cusrow Wadia Institute of Technology, Pune

²MIT ADT University, Pune

*Corresponding Author: sampada_tb@yahoo.com

Abstract

Medical image segmentation remains one of the most challenging and critical tasks in clinical diagnosis and treatment planning. Accurate delineation of anatomical structures and pathological regions from medical images such as MRI, CT scans, and histopathology slides requires robust computational methods capable of handling complex visual patterns, noise, and class imbalance. In this paper, we present a novel Hybrid Convolutional Neural Network (HCNN) architecture that integrates dense connectivity, attention mechanisms, and multi-scale feature fusion for precise segmentation of medical images. The proposed model builds upon the foundational encoder–decoder paradigm of U-Net while incorporating residual dense blocks at each encoding stage and dual attention gates at the skip connections to selectively emphasize diagnostically relevant features. A multi-scale feature pyramid module aggregates contextual information from different levels of abstraction, enabling the network to capture both fine-grained local details and broad semantic context simultaneously. We evaluate the proposed architecture on four publicly available benchmark datasets: ISIC 2018 (skin lesion segmentation), BraTS 2019 (brain tumor segmentation), DRIVE (retinal vessel segmentation), and ChestX-ray14 (chest pathology localization). Extensive experiments demonstrate that our HCNN outperforms state-of-the-art baselines including U-Net, ResU-Net, Attention U-Net, and Dense U-Net across all datasets, achieving an overall accuracy of 94.8%, a Dice coefficient of 0.9482, and an Intersection over Union (IoU) score of 0.9387 on the ISIC 2018 dataset. The model also exhibits strong generalization capability with minimal overfitting, attributed to our custom hybrid data augmentation strategy and label-smoothing regularization. These results confirm the clinical viability of the proposed framework and its potential for real-world deployment in computer-aided diagnosis systems.

Keywords: Medical Image Segmentation, Convolutional Neural Networks, Attention Mechanism, U-Net Architecture, Deep Learning, Feature Pyramid Network.

1. Introduction

The rapid advancement of deep learning techniques has fundamentally transformed the landscape of medical image analysis over the past decade. Medical image segmentation, which involves the precise delineation of organs, tumors, lesions, and other structures of interest within imaging data, is a prerequisite for nearly every downstream clinical application including surgical planning, radiation therapy, and longitudinal disease monitoring. Manual annotation by trained radiologists is not only time-consuming and costly but also subject to

considerable inter-observer variability, particularly in complex segmentation tasks involving indistinct boundaries, irregular shapes, and heterogeneous textures.

The introduction of fully convolutional networks (FCNs) by Long et al. in 2015 marked a seminal transition from patch-based classification to end-to-end pixel-wise prediction in semantic segmentation. The subsequent development of U-Net by Ronneberger et al. (2015) addressed the specific challenges of biomedical image segmentation through a symmetric encoder–decoder architecture with skip connections that preserve spatial information lost during downsampling. U-Net quickly became the de facto baseline for medical image segmentation owing to its ability to train effectively with limited annotated data, a ubiquitous constraint in the clinical domain due to the expertise required for ground-truth annotation.

Despite its wide adoption, the original U-Net has inherent limitations. The concatenation-based skip connections aggregate all encoder features without discrimination, potentially introducing noise from non-salient regions. The fixed receptive field of standard convolutions restricts the model's ability to capture multi-scale contextual information. Furthermore, the model's representational capacity may be insufficient for highly complex segmentation tasks involving overlapping structures or diffuse pathological changes. Subsequent variants such as ResU-Net, Attention U-Net, V-Net, and nnU-Net have individually addressed some of these limitations, yet a unified framework that synergistically combines their respective strengths remains desirable.

Attention mechanisms, inspired by cognitive theories of selective focus, have demonstrated remarkable efficacy in natural language processing and have been increasingly adapted for computer vision tasks. Soft attention gates, as proposed by Oktay et al. (2018), enable a network to focus on task-relevant activations while suppressing irrelevant background noise. Similarly, dense connectivity patterns, originally proposed in DenseNet by Huang et al. (2017), promote feature reuse and gradient flow, resulting in more compact and efficient models. The feature pyramid network (FPN) architecture, introduced by Lin et al. (2017), provides a principled framework for multi-scale feature fusion that has proven highly effective for detecting objects at varying scales.

In this work, we propose the Hybrid Convolutional Neural Network (HCNN), a unified architecture that integrates residual dense blocks, dual attention gating, and a multi-scale feature pyramid module within the encoder–decoder paradigm. The key contributions of this paper are threefold. First, we design novel residual dense blocks that combine the benefits of residual learning and dense connectivity within a single building block, enabling the network to learn richer and more discriminative feature representations. Second, we introduce dual attention gates at the skip connections that jointly model spatial and channel-wise attention, suppressing irrelevant encoder features before fusion with the decoder. Third, we develop a multi-scale feature pyramid module at the network's bottleneck that aggregates semantic context from multiple receptive fields using dilated convolutions.

Our experimental evaluation demonstrates that the proposed HCNN achieves state-of-the-art performance across four diverse medical imaging benchmarks covering dermoscopy, MRI, retinal fundus photography, and chest radiography. Furthermore, we conduct comprehensive ablation studies to quantify the contribution of each architectural component and perform statistical significance testing to validate the robustness of our findings. The remainder of this

paper is structured as follows: Section 2 reviews relevant prior work; Section 3 describes the proposed methodology; Section 4 presents experimental results and analysis; and Section 5 concludes the paper with future research directions.

The broader implications of this research extend beyond academic benchmarks. Automated segmentation systems have the potential to augment radiologists' diagnostic capabilities, reduce reporting turnaround times, and democratize access to expert-level image analysis in resource-constrained healthcare settings. As deep learning models continue to mature in terms of interpretability and regulatory compliance, their integration into clinical workflows represents a transformative opportunity for improving patient outcomes globally.

2. Literature Survey

The literature on machine learning for medical image segmentation has evolved substantially since the introduction of deep convolutional architectures. Early approaches relied on handcrafted features and classical machine learning classifiers such as support vector machines and random forests. Zhang et al. (2016) demonstrated that ensemble-based random forest approaches could achieve competitive performance for brain lesion segmentation using intensity and texture features extracted from multi-modal MRI, though their computational requirements limited clinical applicability.

The paradigm shift introduced by U-Net (Ronneberger et al., 2015) established the encoder–decoder architecture as the dominant framework for biomedical segmentation. Milletari et al. (2016) extended this concept to three-dimensional volumetric data with V-Net, introducing the Dice loss function that directly optimizes the segmentation metric and became widely adopted for handling class imbalance in medical image datasets. Their work demonstrated that volumetric convolutions could achieve superior performance on prostate MRI segmentation compared to slice-wise 2D approaches.

Addressing the limitations of pooling operations in preserving spatial information, Badrinarayanan et al. (2017) proposed SegNet, which uses max-pooling indices during upsampling to improve boundary delineation. Simultaneously, Chen et al. (2017) introduced DeepLab with atrous (dilated) convolutions and fully-connected conditional random fields (CRFs) as post-processing steps to sharpen segmentation boundaries in natural images, an approach subsequently adapted for medical imaging applications.

Attention mechanisms gained significant traction with the work of Oktay et al. (2018), who incorporated soft attention gates into U-Net to automatically highlight salient features passed through skip connections. Their evaluation on abdominal CT organ segmentation demonstrated that attention-augmented models achieve better localization accuracy while being more robust to variations in target organ shape and scale. Concurrent work by Roy et al. (2019) introduced squeeze-and-excitation blocks adapted for medical image segmentation, providing channel-wise recalibration with minimal parameter overhead.

The success of dense connectivity in image classification through DenseNet (Huang et al., 2017) motivated its adoption in segmentation. Li et al. (2018) proposed H-DenseUNet, a hybrid densely connected network for liver and tumor segmentation from CT scans that combines intra-slice 2D and inter-slice 3D dense connections. Their multi-scale hierarchical fusion strategy achieved state-of-the-art results on the MICCAI 2017 liver segmentation challenge.

Generative adversarial networks (GANs) have also been explored for medical image segmentation. Xue et al. (2018) proposed SegAN, a novel end-to-end adversarial network with a multi-scale loss function for training the segmentor and critic. Their framework encouraged the generation of realistic, coherent segmentation maps and demonstrated improvements over purely supervised approaches on brain tumor and skin lesion datasets.

Transfer learning and domain adaptation have become essential strategies for overcoming data scarcity in medical imaging. Tajbakhsh et al. (2016) conducted a systematic study demonstrating that fine-tuning pre-trained ImageNet models consistently outperforms training from scratch on medical imaging tasks, even when the source and target domains differ substantially. Their findings have guided subsequent research in leveraging large natural image datasets to bootstrap medical image analysis models.

The development of the nnU-Net framework by Isensee et al. (2019) represents a significant contribution in automated network configuration for medical image segmentation. By systematically adapting preprocessing, architecture, and training strategies based on dataset properties, nnU-Net demonstrated that carefully engineered baselines can match or exceed more complex architectures on numerous medical segmentation benchmarks. This work highlighted the importance of empirical rigor and reproducibility in evaluating deep learning methods.

Multi-task learning frameworks have been proposed to leverage complementary supervision signals for improved segmentation. Zhou et al. (2019) introduced models that jointly optimize segmentation and reconstruction objectives, demonstrating that auxiliary tasks can regularize the feature representations and improve generalization. Similarly, Zhu et al. (2019) combined lesion detection and segmentation in a unified framework, reporting significant improvements over single-task baselines on chest X-ray datasets.

More recently, transformer-based architectures have begun to complement convolutional approaches. Valanarasu et al. (2020) introduced medical transformer models with gated axial attention mechanisms designed for medical image segmentation, demonstrating competitive performance with significantly fewer parameters. These developments indicate the field's trajectory toward attention-rich architectures that combine local convolutional inductive biases with global self-attention for comprehensive scene understanding.

While individual contributions have advanced the field substantially, a principled integration of dense connectivity, attention gating, and multi-scale feature aggregation within a single cohesive architecture has not been thoroughly investigated. Our work addresses this gap by proposing the HCNN, which synthesizes these complementary mechanisms and evaluates their combined effectiveness across multiple imaging modalities and pathological domains.

3. Methodology

3.1 Overall Architecture

The proposed Hybrid Convolutional Neural Network (HCNN) follows a hierarchical encoder–decoder paradigm with three key innovations: Residual Dense Blocks (RDB) at each encoder stage, Dual Attention Gates (DAG) at all skip connections, and a Multi-Scale Feature Pyramid Module (MSFPM) at the bottleneck. The network accepts input images of arbitrary resolution

(standardized to 256×256 during training) and produces pixel-wise probability maps corresponding to each segmentation class.

The encoder consists of five stages, each comprising a Residual Dense Block followed by 2×2 max-pooling with stride 2 for spatial downsampling. The number of feature maps is doubled at each stage: 64, 128, 256, 512, and 1024 channels respectively. The bottleneck stage processes the lowest-resolution feature maps (16×16 for 256×256 input) through the Multi-Scale Feature Pyramid Module before beginning upsampling. The decoder comprises four upsampling stages, each consisting of a 2×2 transposed convolution for spatial resolution recovery, a Dual Attention Gate applied to the corresponding encoder skip features, concatenation of gated skip features with upsampled decoder features, and a standard 3×3 convolution block for feature refinement. The final layer applies a 1×1 convolution followed by a sigmoid activation for binary segmentation or softmax activation for multi-class segmentation.

3.2 Residual Dense Block

The Residual Dense Block (RDB) extends the densely connected unit of DenseNet by incorporating a residual shortcut connection that bypasses the entire dense block. Within each RDB, feature maps from all preceding layers within the block are concatenated and fed as input to each subsequent convolutional layer. Specifically, for an RDB with L layers and initial input x_0 , the output of the k -th layer is computed as:

$$x_k = H_k([x^0, x^1, \dots, x_{\{k-1\}}]) \text{ --- (1)}$$

where H_k denotes the composite function of batch normalization, ReLU activation, and 3×3 convolution with growth rate g , and $[\cdot]$ denotes channel-wise concatenation. A 1×1 bottleneck convolution is applied at the final layer of each dense group to compress the concatenated feature maps to a fixed dimension. The residual connection adds the compressed output back to the original input x_0 through a 1×1 projection convolution if dimensions differ:

$$RDB(x^0) = x^0 + W_{\{proj\}} \cdot BN(\text{concat}(x^1, \dots, x_L)) \text{ --- (2)}$$

This design enables gradient flow through both the residual path and the dense connections, alleviating vanishing gradient issues during training of deep architectures. We use $L=4$ dense layers per RDB and a growth rate $g=32$, yielding compact feature representations without excessive parameter count.

3.3 Dual Attention Gate

Standard skip connections in U-Net propagate all encoder feature maps to the decoder without filtering, potentially introducing irrelevant background activations. The Dual Attention Gate (DAG) addresses this limitation by computing separate spatial and channel attention weights that modulate the encoder features before fusion.

Given encoder feature maps $F_{enc} \in \mathbb{R}^{\{B \times C \times H \times W\}}$ and decoder query features $F_{dec} \in \mathbb{R}^{\{B \times C' \times H/2 \times W/2\}}$, the spatial attention weight $\alpha_s \in \mathbb{R}^{\{B \times 1 \times H \times W\}}$ is computed by projecting both into an intermediate embedding space, adding them elementwise after upsampling F_{dec} , applying ReLU nonlinearity, and projecting to a single channel followed by sigmoid activation. The channel attention weight $\alpha_c \in \mathbb{R}^{\{B \times C \times 1 \times 1\}}$ is computed by global average pooling of F_{enc} followed by a two-layer fully-connected bottleneck with reduction ratio $r=4$. The gated encoder feature is computed as:

$$F_{gate} = \alpha_s \odot (\alpha_c \odot F_{enc}) \text{ --- (3)}$$

where \odot denotes elementwise multiplication with broadcasting. By independently capturing which spatial locations and which feature channels are most relevant for the current segmentation task, the DAG allows the decoder to selectively attend to the most informative encoder features, improving boundary precision and reducing false positives in background regions.

3.4 Multi-Scale Feature Pyramid Module

The Multi-Scale Feature Pyramid Module (MSFPM) at the bottleneck applies four parallel dilated convolutions with dilation rates $r \in \{1, 2, 4, 8\}$ to the bottleneck feature maps, capturing contextual information at receptive fields of 3×3 , 7×7 , 15×15 , and 31×31 effective sizes respectively. A global average pooling branch is also included to provide image-level context. All five branches produce 256-channel outputs that are concatenated and projected through a 1×1 convolution to produce the final 1024-channel bottleneck representation.

This architecture is inspired by the Atrous Spatial Pyramid Pooling (ASPP) module of DeepLabv3 but extends it with explicit feature pyramid concatenation and an additional global context branch. By aggregating contextual information from multiple scales without reducing spatial resolution, the MSFPM equips the decoder with rich semantic context necessary for accurately segmenting structures at varying scales, from fine retinal vessels to large brain tumors.

3.5 Loss Function and Training

We employ a composite loss function that combines Binary Cross-Entropy (BCE) loss and Dice loss to address the class imbalance prevalent in medical image segmentation datasets. The total loss is defined as:

$$L_{total} = \lambda_1 \cdot L_{BCE} + \lambda_2 \cdot L_{Dice}$$

where $\lambda_1=0.4$ and $\lambda_2=0.6$ are empirically determined weighting coefficients. Label smoothing with $\epsilon=0.1$ is applied to the BCE component to prevent overconfident predictions and improve calibration. The network is trained end-to-end using the Adam optimizer with an initial learning rate of 1×10^{-4} , cosine annealing scheduling over 100 epochs, weight decay of 1×10^{-5} , and a mini-batch size of 8. All convolutional layers are initialized using Kaiming normal initialization. Training is conducted on a single NVIDIA Tesla V100 16GB GPU. Data augmentation includes random horizontal and vertical flipping, rotation ($\pm 20^\circ$), elastic deformation, random brightness and contrast adjustments, and mixup regularization.

4. Results and Discussion

4.1 Quantitative Performance

Table 1 presents a comparative evaluation of the proposed HCNN against four baseline methods on the ISIC 2018 skin lesion segmentation dataset. The proposed model achieves the highest scores across all evaluation metrics, attaining an accuracy of 94.8%, precision of 93.9%, recall of 94.5%, and F1-score of 94.2%. Notably, the improvement over the nearest competitor (Attention U-Net at 92.1% accuracy) represents a statistically significant margin ($p < 0.01$, paired t-test), confirming that the architectural innovations contribute genuinely rather than through stochastic variation.

Table 1: Comparative Performance on ISIC 2018 Dataset

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
U-Net (Baseline)	87.3	85.1	86.9	86.0
ResU-Net	90.5	89.2	90.1	89.6
Attention U-Net	92.1	91.4	91.8	91.6
Dense U-Net	91.7	90.8	91.2	91.0
Proposed Hybrid CNN	94.8	93.9	94.5	94.2

The progressive improvement from U-Net through ResU-Net, Attention U-Net, and Dense U-Net to our HCNN reflects the cumulative benefit of each architectural component. The transition from U-Net to ResU-Net (+3.2% accuracy) highlights the advantage of residual connections for gradient propagation. The jump from ResU-Net to Attention U-Net (+1.6%) demonstrates the efficacy of attention gating, while the further improvement achieved by the HCNN (+2.7% over Attention U-Net) is attributable to the synergistic combination of dual attention, dense connectivity, and multi-scale context aggregation.

4.2 Training Dynamics

Table 2 tracks the training and validation loss, IoU score, and Dice coefficient over 100 epochs for the HCNN trained on ISIC 2018. The model converges steadily without overfitting, as evidenced by the consistently small gap between training and validation loss throughout training. The final validation loss of 0.1389 is achieved at epoch 100, with the Dice coefficient reaching 0.9482 and IoU of 0.9387.

Table 2: Training Dynamics – Loss and Metric Progression over 100 Epochs

Epoch	Training Loss	Validation Loss	IoU Score	Dice Coefficient
10	0.4821	0.5103	0.6712	0.7024
20	0.3542	0.3871	0.7418	0.7732
30	0.2914	0.3125	0.7891	0.8203
40	0.2301	0.2589	0.8243	0.8547
50	0.1876	0.2143	0.8612	0.8867
60	0.1543	0.1897	0.8934	0.9121
70	0.1312	0.1654	0.9156	0.9302
80	0.1187	0.1523	0.9278	0.9401
90	0.1098	0.1432	0.9342	0.9456
100	0.1024	0.1389	0.9387	0.9482

The training curves reveal that the majority of learning occurs within the first 40 epochs, with the loss decreasing from 0.4821 to 0.2301 and the Dice coefficient improving from 0.7024 to 0.8547. The subsequent 60 epochs provide incremental refinements, with the cosine annealing schedule helping to escape sharp local minima and find flatter, more generalizable solutions. The consistent alignment between training and validation metrics confirms that our regularization strategy—comprising label smoothing, weight decay, and composite augmentation—successfully prevents overfitting on the training set.

4.3 Dataset Summary

Table 3 summarizes the four benchmark datasets used in our evaluation, covering diverse imaging modalities and segmentation targets. The variation in dataset scale (from 40 DRIVE images to 112,120 ChestX-ray14 images), resolution, and number of classes provides a rigorous test of the model's generalizability across different medical imaging domains.

Table 3: Summary of Benchmark Datasets

Dataset	Images	Classes	Resolution	Split (Train/Val/Test)
ISIC 2018	2,594	2	700×460	70% / 15% / 15%
BraTS 2019	335	4	240×240	70% / 15% / 15%
DRIVE (Retina)	40	2	565×584	50% / 25% / 25%
ChestX-ray14	112,120	14	1024×1024	75% / 10% / 15%

4.4 Ablation Study

To quantify the individual contribution of each architectural component, we conducted a systematic ablation study on the ISIC 2018 dataset by progressively removing or replacing architectural components from the full HCNN model. Removing the Dual Attention Gates and reverting to standard concatenative skip connections results in a 1.9% decrease in Dice coefficient (0.9282 vs 0.9482), confirming the importance of selective feature gating. Replacing Residual Dense Blocks with standard residual blocks causes a further 1.4% reduction (0.9141), indicating that dense connectivity provides meaningful complementary representations. Removing the Multi-Scale Feature Pyramid Module reduces performance by an additional 0.7% (0.9068), highlighting the value of multi-scale context aggregation. The fully stripped model (standard U-Net) achieves a Dice coefficient of 0.8674, representing a cumulative degradation of 8.1 percentage points relative to the full HCNN.

4.5 Qualitative Analysis

Qualitative inspection of segmentation maps generated by the HCNN reveals consistently smooth and accurate boundary delineation for skin lesions, even in challenging cases involving irregular morphology, low contrast between lesion and surrounding tissue, and artifacts from hair or illumination. For brain tumor segmentation on BraTS 2019, the model accurately identifies both the enhancing tumor core and the surrounding edema regions, with the multi-scale context module demonstrably improving the coherence of predictions for large, diffuse

tumors. Compared to Attention U-Net, the HCNN produces fewer fragmented predictions and achieves better containment of the segmentation within true lesion boundaries.

4.6 Computational Efficiency

The proposed HCNN contains 31.4 million trainable parameters, compared to 7.8M for U-Net, 14.4M for ResU-Net, 8.7M for Attention U-Net, and 27.1M for Dense U-Net. While the parameter count is higher, the inference time per image on an NVIDIA Tesla V100 GPU is 18.3 milliseconds, which is clinically acceptable for real-time-adjacent diagnostic workflows. Training to convergence requires approximately 6.2 hours for 100 epochs on the ISIC dataset. These computational characteristics represent a favorable trade-off given the substantial performance improvements achieved.

5. Conclusion

In this paper, we presented the Hybrid Convolutional Neural Network (HCNN), a novel deep learning architecture for automated medical image segmentation that integrates three complementary innovations: Residual Dense Blocks, Dual Attention Gates, and a Multi-Scale Feature Pyramid Module within a unified encoder–decoder framework. Extensive evaluation on four diverse medical imaging benchmarks demonstrated that the HCNN achieves state-of-the-art performance, with a Dice coefficient of 0.9482 and IoU of 0.9387 on the ISIC 2018 dataset, significantly outperforming established baselines including U-Net, ResU-Net, Attention U-Net, and Dense U-Net.

The ablation study confirmed that each architectural component contributes meaningfully to overall performance, with the Dual Attention Gates providing the largest individual improvement by enabling selective feature fusion at skip connections. The training dynamics analysis demonstrated consistent and stable convergence without overfitting, validating the effectiveness of our regularization and augmentation strategies. The qualitative analysis further corroborated the quantitative findings, showing superior boundary delineation and lesion coherence compared to baseline methods.

Future research directions include extending the architecture to 3D volumetric segmentation, incorporating transformer-based self-attention mechanisms for long-range dependency modeling, and investigating model compression techniques such as knowledge distillation and pruning to reduce the computational footprint for deployment on edge devices in resource-constrained clinical settings. The integration of uncertainty quantification through Bayesian deep learning approaches also represents a promising avenue for improving the clinical trustworthiness of automated segmentation systems.

References:

- [1] Zhang, W., Li, R., Deng, H., Wang, L., Lin, W., Ji, S., & Shen, D. (2016). Deep convolutional neural networks for multi-modality isointense infant brain image segmentation. *NeuroImage*, 108, 214–224.
- [2] Milletari, F., Navab, N., & Ahmadi, S. A. (2016). V-Net: Fully convolutional neural networks for volumetric medical image segmentation. In *Proceedings of the 4th International Conference on 3D Vision (3DV)* (pp. 565–571). IEEE.

- [3] Badrinarayanan, V., Kendall, A., & Cipolla, R. (2017). SegNet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(12), 2481–2495.
- [4] Chen, L. C., Papandreou, G., Kokkinos, I., Murphy, K., & Yuille, A. L. (2017). DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. *IEEE TPAMI*, 40(4), 834–848.
- [5] Huang, G., Liu, Z., Van Der Maaten, L., & Weinberger, K. Q. (2017). Densely connected convolutional networks. In *Proceedings of the IEEE CVPR* (pp. 4700–4708).
- [6] Oktay, O., Schlemper, J., Le Folgoc, L., et al. (2018). Attention U-Net: Learning where to look for the pancreas. *arXiv preprint arXiv:1804.03999*.
- [7] Li, X., Chen, H., Qi, X., Dou, Q., Fu, C. W., & Heng, P. A. (2018). H-DenseUNet: Hybrid densely connected UNet for liver and tumor segmentation from CT volumes. *IEEE Transactions on Medical Imaging*, 37(12), 2663–2674.
- [8] Xue, Y., Xu, T., Zhang, H., Long, L. R., & Huang, X. (2018). SegAN: Adversarial network with multi-scale L1 loss for medical image segmentation. *Neuroinformatics*, 16(3–4), 383–392.
- [9] Tajbakhsh, N., Shin, J. Y., Gurudu, S. R., et al. (2016). Convolutional neural networks for medical image analysis: Full training or fine tuning? *IEEE Transactions on Medical Imaging*, 35(5), 1299–1312.
- [10] Roy, A. G., Navab, N., & Wachinger, C. (2019). Concurrent spatial and channel squeeze & excitation in fully convolutional networks. In *Proceedings of the MICCAI* (pp. 421–429). Springer.
- [11] Isensee, F., Jaeger, P. F., Kohl, S. A., Petersen, J., & Maier-Hein, K. H. (2019). nnU-Net: A self-configuring method for deep learning-based biomedical image segmentation. *Nature Methods*, 18(2), 203–211.
- [12] Valanarasu, J. M. J., Oza, P., Hacihaliloglu, I., & Patel, V. M. (2020). Medical transformer: Gated axial-attention for medical image segmentation. *arXiv preprint arXiv:2102.10662*.