

# RETHINKING INTERMEDIATE LAYERS DESIGN IN KNOWLEDGE DISTILLATION FOR KIDNEY AND LIVER TUMOR SEGMENTATION

Vandan Gorade<sup>1</sup>, Sparsh Mittal<sup>2</sup>, Debesh Jha<sup>3</sup>, Ulas Bagci<sup>3</sup>

<sup>1</sup>Artificial Intelligence & Data Science, Jio Institute, Navi Mumbai, India

<sup>2</sup>Mehta Family School of DS&AI, Indian Institute of Technology, Roorkee, India

<sup>3</sup>Machine & Hybrid Intelligence Lab, Department of Radiology, Northwestern University, USA

Knowledge distillation (KD) has demonstrated remarkable success across various domains, but its application to medical imaging tasks, such as kidney and liver tumor segmentation, has encountered challenges. Many existing KD methods are not specifically tailored for these tasks. Moreover, prevalent KD methods often lack a careful consideration of ‘what’ and ‘from where’ to distill knowledge from the teacher to the student. This oversight may lead to issues like the accumulation of training bias within shallower student layers, potentially compromising the effectiveness of KD. To address these challenges, we propose Hierarchical Layer-selective Feedback Distillation (HLFD). HLFD strategically distills knowledge from a combination of middle layers to earlier layers and transfers final layer knowledge to intermediate layers at both the feature and pixel levels. This design allows the model to learn higher-quality representations from earlier layers, resulting in a robust and compact student model. Extensive quantitative evaluations reveal that HLFD outperforms existing methods by a significant margin. For example, in the kidney segmentation task, HLFD surpasses the student model (without KD) by over 10%, significantly improving its focus on tumor-specific features. From a qualitative standpoint, the student model trained using HLFD excels at suppressing irrelevant information and can focus sharply on tumor-specific details, which opens a new pathway for more efficient and accurate diagnostic tools.

## 1. INTRODUCTION

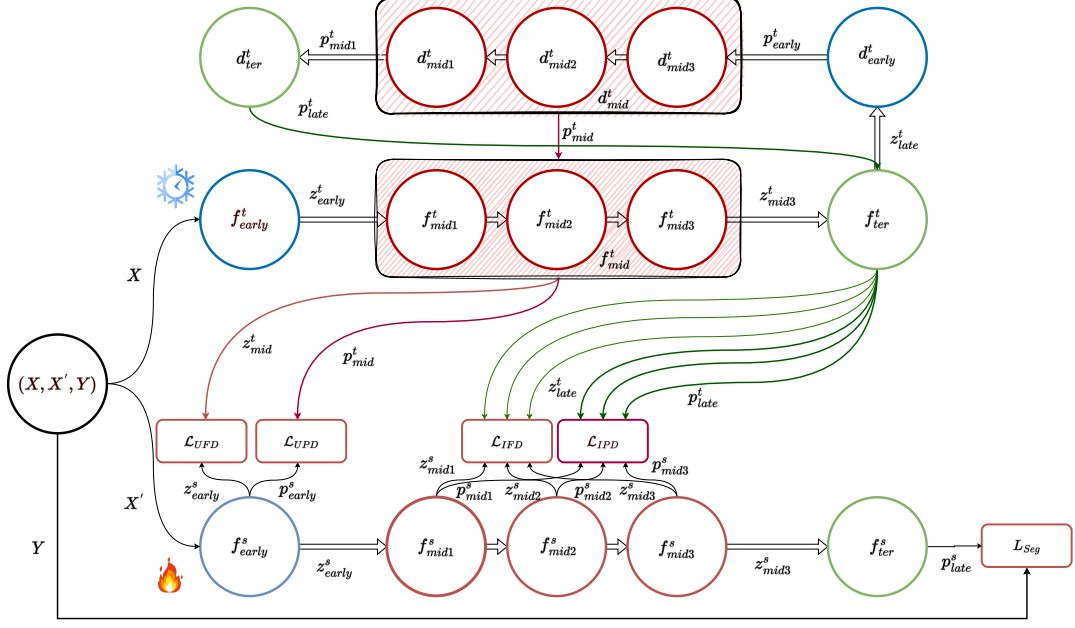
Tumor segmentation in medical imaging enables clinicians to accurately identify, assess, and manage malignancies. Leveraging neural networks, we achieve automated, high-fidelity delineation of tumor boundaries in various imaging modalities, including CT scans and MRIs [1, 2]. This technological breakthrough elevates diagnostic accuracy and efficiency and streamlines treatment planning [3], ultimately leading to enhanced patient care and outcomes.

Significant challenges persist despite the remarkable successes of deep-learning models in tumor segmentation. These models demand extensive datasets and substantial computational resources, making deployment on resource-limited

devices a hurdle. Furthermore, the diversity in tumor appearances, irregular sizes, unpredictable locations, and variations amplifies segmentation complexity. To address these challenges, researchers are exploring innovative strategies. For instance, lightweight networks [4, 5] have been explored for real-time semantic segmentation, and recent works have delved into real-time medical image segmentation. However, model simplification may hurt predictive performance. Knowledge distillation (KD) [6] has emerged as a valuable approach, facilitating knowledge transfer from the larger ‘teacher’ models to the leaner ‘student’ models.

Existing works [7, 8, 9, 10, 11, 12] aim to enhance the final representations of the student model by minimizing the difference in softmax representations between the teacher and student models. However, this supervisory signal originates solely from the final student layer. Hence, it tends to attenuate with each layer during backpropagation, accumulating training bias within the shallower student layers. This impairs the efficacy of knowledge transfer. Other works [13, 14, 15, 16] focus on improving the alignment of latent feature maps by mimicking intermediate representations. These intermediate representations serve as solid indicators that facilitate learning the final representation. However, when we replicate intermediate representations, we are limited to capturing the knowledge acquired by that specific layer, potentially missing out on global information. Recognizing this limitation, capturing features from terminal representation at an earlier stage emerges as a valuable strategy [17, 18]. However, these methods give suboptimal results where boundary, shape, texture information, and a combination of low-level features are essential, not only high-level class information.

We have introduced hierarchical layer-selective feedback distillation (HLFD) to address these challenges. HLFD comprises Feature-level LFD (FLFD) and Pixel-level LFD (PLFD). FLFD, in turn, includes Unified Feature-level Distillation (UFD) for unified representations and Individual Feature-level Distillation (IFD) for middle-to-early and later-to-middle layer distillation. PLFD, on the other hand, involves Unified Pixel-level Distillation (UPD) and Individual Pixel-level Distillation (IPD), which transfers pixel-level



**Fig. 1.** The input  $X$  and augmentation  $X'$ , undergo encoding by both a pre-trained teacher encoder and a randomly initialized student encoder, resulting in representations  $z^{t, early}, z^{t, midj}, z^{t, ter}$ , and  $z^{s, early}, z^{s, midj}, z^{s, ter}$ , respectively. These representations contribute to feature-level loss functions,  $\mathcal{L}_{UFD}$  and  $\mathcal{L}_{IFD}$ . Additionally, the teacher decoder decodes  $z^{t, ter}$ , producing representations  $p^{t, early}, p^{t, midj}, p^{t, ter}$ , which are utilized in pixel-level loss functions,  $\mathcal{L}_{UPD}$  and  $\mathcal{L}_{IPD}$ . The training process is further enhanced with the inclusion of a supervised focal dice loss ( $\mathcal{L}_{seg}$ ).

knowledge from the teacher decoder to the student through interpolated features. HLF D integrates both FLFD and PLFD in a multi-task fashion, promoting simultaneous learning of feature-level and pixel-level representations. Our contributions are as follows:

- We rethink the design of layers in the context of distillation and introduce the Hierarchical Layer-selective Feedback Distillation (HLFD) framework.
- We demonstrate HLF D’s capability to capture tumor-specific details from early layers while effectively suppressing irrelevant information flow.
- Extensive experiments conducted on kidney and liver tumor segmentation tasks establish that our proposed method attains state-of-the-art (SOTA) results

## 2. PROPOSED METHOD

### 2.1. Feature-level Layer-selective Feedback Distillation (FLFD)

Given an input  $X$ , transformations occur through both the pre-trained teacher encoder  $f_i^t$  and the random student encoder  $f_i^s$ , denoted by  $i$  for the number of blocks. This yields early representations  $z^{t, early}$  and  $z^{s, early}$ , intermediate representations  $z^{t, midj}$  and  $z^{s, midj}$  (where  $j$  is the number of middle layers), and terminal representations  $z^{t, late}$  and  $z^{s, late}$ . These

representations form the foundation of our framework. We propose an FLFD loss, defined as,  $\mathcal{L}_F = \mathcal{L}_{UFD} + \mathcal{L}_{IFD}$ . These components are defined below.

**Unified Feature-level Distillation (UFD).** Within this framework, we introduce the concept of distilling the attentive knowledge from the teacher’s unified representation of middle layers,  $z^{t, mid}$ , to the student’s early representation,  $z^{s, early}$ . To achieve this, we propose the following loss function.

$$\mathcal{L}_{UFD} = \frac{\|\mathcal{A}(z^{s, early})\|}{\|\mathcal{A}(z^{s, early})\|_2} - \frac{\|\mathcal{A}(z^{t, mid})\|}{\|\mathcal{A}(z^{t, mid})\|_2} \quad (1)$$

To achieve  $z^{t, mid}$ , we perform interpolation on the middle layers with the larger feature maps to ensure their spatial dimensions match the smallest among them. Next, we concatenate all these interpolated representations along the channel dimension. Finally, operation  $\mathcal{A}(\cdot)$  is employed first to rescale the student’s representation  $z^{s, early}$  to match the spatial dimension of the teacher’s  $z^{t, mid}$ . Additionally, channel normalization is applied to the rescaled student representation, assuming that the absolute value of a neuron activation signifies its importance.

**Individual Feature-level Distillation (IFD).** Within this framework, we introduce the concept of distilling the attentive knowledge from the teacher’s late representation,  $z^{t, late}$ , to each student’s middle layers or intermediate representa-

tion,  $z_{mid_j}^s$ . To achieve this, we propose the following loss function.

$$\mathcal{L}_{IFD} = \sum_{j=1}^N \frac{\|\mathcal{A}(z_{mid_j}^s)\|}{\|\mathcal{A}(z_{mid_j}^s)\|_2} - \frac{\|\mathcal{A}(z_{late}^t)\|}{\|\mathcal{A}(z_{late}^t)\|_2} \quad (2)$$

Here, the operation  $\mathcal{A}(\cdot)$  is same as in Eq.1.

## 2.2. Pixel-level Layer-selective Feedback Distillation (PLFD):

In contrast to feature-level distillation, pixel-level segmentation-map distillation is geared toward conveying pixel-wise predictions. In practice, we distill pixel-level maps generated by the teacher’s decoder to interpolated student maps. First, the teacher encoder output  $z_{late}^t$  is passed through pre-trained teacher decoder  $d_i^t$  resulting in early predictive map  $p_{early}^t$ , intermediate predictive map  $p_{mid_j}^t$  and terminal predictive map  $p_{late}^t$ . For students, we used an interpolated representation map. We propose  $\mathcal{L}_P = \mathcal{L}_{UPD} + \mathcal{L}_{IPD}$ . The components of PLFD are as follows.

**Unified Pixel-level Distillation (UPD).** We propose distilling the precise predictive information from the teacher’s unified pixel-wise predictive maps of middle layers,  $p_{mid_j}^t$ , to the student’s early interpolated representation,  $p_{early}^s$ .

$$\mathcal{L}_{UPD} = \text{KL}(\mathcal{A}(p_{early}^s) || p_{mid_j}^t) \quad (3)$$

**Individual Pixel-level Distillation.** Here, the teacher’s terminal predictive map, denoted as  $p_{late}^t$ , distills precise information to the intermediate predicted maps of the students individually, represented as  $p_{mid_j}^s$ . This allows the student to capture detailed knowledge about the exact pixel locations and their corresponding class assignments within the image from much earlier layers. To achieve this, a KL-divergence loss is employed between these maps:

$$\mathcal{L}_{IPD} = \frac{1}{N} \sum_j \text{KL}(\mathcal{A}(p_{mid_j}^s) || p_{late}^t) \quad (4)$$

Here,  $N$  is the number of middle-layer blocks in the student.

## 2.3. Hierarchical Layer-selective Feedback Distillation (HLFD)

Finally, distilling both feature-level and pixel-level representations allows the student to learn fine-to-coarse hierarchical details at both the feature and pixel levels. The multi-task loss function can be defined as:

$$\mathcal{L}_H = \mathcal{L}_{Seg} + \beta * \mathcal{L}_F + \lambda * \mathcal{L}_P \quad (5)$$

Where  $\mathcal{L}_{seg}$  is the focal dice loss used for training the student network in a supervised fashion. In the inference phase, post-sufficient training, both the teacher network components and distillation modules are discarded.

## 3. EXPERIMENTAL PLATFORM

**Datasets:** We evaluated our techniques on kidney tumor segmentation (KiTS) [2] and liver tumor segmentation (LiTS) [19] datasets. KiTS comprises 210 abdominal CT scans, where a 168:42 split is used for testing and training. Similarly, the LiTS dataset consists of 201 CT scans and uses the split of 131:70.

**Baselines:** We compare our method with the following SOTA methods: i) Structured Knowledge Distillation (SKD) [11]: Involves pair-wise distillation to capture similarity at feature and pixel level. ii) Intermediate Feature Distiller (IFD) [17]: Distills the teacher’s terminal representation into concatenated branches of the student model. iii) Deep Knowledge Distillation (DKD) [16]: Similar to [16] but without the Relational Knowledge distillation(RKD) module. iv) Hierarchical Individual Feedback Knowledge Distillation (HIFD) [18]: It distills the teacher’s terminal representation to individual layers of the student. We extended this method for segmentation by incorporating pixel-level feedback distillation loss functions. We maintained identical implementation settings across all techniques.

**Implementation Details:** We employed UNet++ [20] architecture (36.1M parameters) as the teacher network and ResNet18 [21] (11.6M parameters) as the student network. Our segmentation networks and distillation processes, inspired by [16], were trained using Adam optimizer with beta1 (0.9) and beta2 (0.999). The learning rate began at 0.001, utilizing CosineAnnealing for rate scheduling, reaching a minimum of 0.000001. Data augmentation techniques such as random rotation and flipping were applied, while experiments revealed that Gaussian noise augmentation is unsuitable for medical images. Most networks processed authentic  $512 \times 512$  CT images, requiring windowing of HU values with radiological standards (e.g., -40 to 160 for the liver and -200 to 300 for the kidney). We use the PyTorch framework. We train all the networks till convergence with up to 120 epochs. We report the result as *mean*  $\pm$  *std* after three runs. For the Dice score (DSC), higher is better. For Relative Volume Difference (RVD), smaller absolute values are desired, indicating a closer match between the predicted and ground truth volumes. When comparing RVD values, a smaller absolute value (closer to zero) is better, regardless of whether the RVD is positive or negative. These metrics provide complementary insights about the performance.

## 4. RESULTS

**Quantitative Results:** As shown in Table 1, our method, HLFD, consistently outshines both the supervised student and the baseline models. Notably, on the KiTS dataset, we observe a substantial enhancement in DSC compared to the student (without KD). Furthermore, both IFD and HIFD exhibit

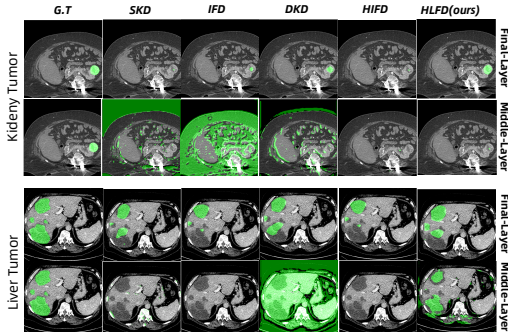
competitive or superior outcomes than other baseline methods. These results underscore the necessity of integrating KD and also emphasize the critical importance of architecting layers that adeptly distill the ‘what’ and ‘where’ dimensions of knowledge from the teacher model. On the RVD metric also, HLF D outperforms baselines, including the student (without KD), by a significant margin. This insight into volume differences holds valuable implications, especially in tasks like tumor segmentation where accuracy in volume is of paramount importance.

**Table 1.** Quantitative Results ( $\beta = 0.9$  and  $\lambda = 0.1$ )

Method	KiTS		LiTS	
	DSC	RVD	DSC	RVD
Teacher	64.50 ± 1.45	-0.203	57.84 ± 1.55	1.434
Student(w/o KD)	41.30 ± 2.30	-0.421	41.19 ± 1.65	0.701
SKD [11]	38.71 ± 2.53	-0.411	42.09 ± 2.01	<b>0.018</b>
IFD [17]	46.79 ± 1.25	-0.275	45.36 ± 1.20	0.122
DKD [16]	40.21 ± 2.35	-0.496	43.72 ± 0.87	0.234
HIFD [18]	42.50 ± 1.25	-0.434	44.20 ± 1.22	0.187
HLFD(ours)	<b>52.18 ± 2.55</b>	<b>-0.176</b>	<b>48.75 ± 2.23</b>	0.173

On the LiTS dataset, HLF D (our method) consistently outperforms the baselines on the DSC metric, whereas the SKD method is the best on the RVD metric.

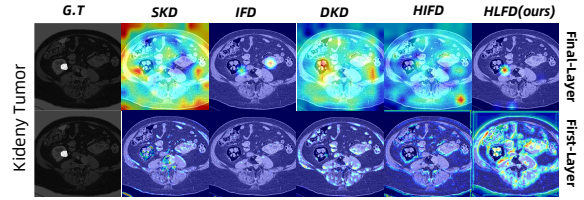
**Qualitative Results:** The visualizations presented in Fig. 2 showcase the superior performance of our proposed HLF D method. HLF D accurately segments the Region of Interest (ROI) while effectively suppressing irrelevant information, even at intermediate layers. The previous techniques fail completely in recognizing the segmentation ROI at intermediate layers, especially for liver segmentation. This underscores the importance of meticulously designing layers to enhance the segmentation task’s representation quality.



**Fig. 2.** The green color highlights the regions of interest (ROI) representing tumors. The segmentation maps are presented for both KiTS (first two rows) and LiTS (last two rows) datasets, with ‘G.T.’ denoting the ground truth.

The GradCAM maps presented in Fig. 3 showcase distinct patterns among methods. SKD, which does not leverage intermediate layers, exhibits a flow of irrelevant information, hindering focus on tumor-specific details. While DKD shows some restriction of information, IFD and HIF D manage to

suppress irrelevant details. However, they face challenges in focusing on tumor-specific information. In contrast, our method distinctly focuses on tumor-specific information without capturing irrelevant details.



**Fig. 3.** Gradient-activated class maps for KiTS19, featuring CAM results for both the final and first layers. Remarkably, HLF D clearly focus on the tumor region, maintaining effectiveness in suppressing irrelevant information as the process advances to the final layer.

**Table 2.** Impact of  $\beta$  and  $\lambda$  on DSC

$\beta$	$\lambda$	KiTS	LiTS
0.9	0.1	52.18 ± 2.55	48.75 ± 2.23
1.8	0.1	51.58 ± 1.25	48.18 ± 1.63
0.9	0.2	51.95 ± 2.35	48.48 ± 2.05

**Sensitivity Analysis:** From Table 2, doubling the  $\beta$  value (from 0.9 to 1.8) while maintaining  $\lambda$  constant led to a slight deterioration in performance. Conversely, increasing the value of  $\lambda$  (from 0.1 to 0.2) while keeping  $\beta$  constant showed a similar trend but with slightly improved performance compared to the previous case. This suggests that  $\mathcal{L}_F$  learns a rich representation of data, while  $\mathcal{L}_P$  learns the essential structure for the segmentation task. Therefore, maintaining  $\beta$  greater than  $\lambda$  is crucial for optimal results. The best performance was achieved with  $\beta = 0.9$  and  $\lambda = 0.1$ .

## 5. CONCLUSION

We introduce a novel Knowledge Distillation (KD) framework for enhancing liver and kidney tumor segmentation, redefining knowledge selection to distill and the distillation source, and transitioning from teacher encoder layers to the student. Quantitatively, HLF D has demonstrated remarkable superiority over existing KD techniques and baseline models. Our method substantially improves DSC, particularly in kidney tumor segmentation, where HLF D surpasses the student model (without KD) by over 10%. Qualitatively, HLF D exhibits exceptional capabilities in suppressing irrelevant information while maintaining a sharp focus on tumor-specific details. The ability of HLF D to accurately segment the region of interest (ROI) at both intermediate and final layers showcases its effectiveness in enhancing the quality of segmentation representations for kidney and liver tumor segmentation.

## 6. COMPLIANCE WITH ETHICAL STANDARDS

This is a numerical simulation study for which no ethical approval was required.

## 7. REFERENCES

- [1] Abubaker Abdelrahman and Serestina Viriri, “Kidney tumor semantic segmentation using deep learning: A survey of state-of-the-art,” *Journal of Imaging*, vol. 8, no. 3, pp. 55, 2022.
- [2] Nicholas Heller, Niranjan Sathianathen, Arveen Kalapara, Edward Walczak, Keenan Moore, Heather Kaluzniak, Joel Rosenberg, Paul Blake, Zachary Rengel, Makinna Oestreich, et al., “The kits19 challenge data: 300 kidney tumor cases with clinical context, ct semantic segmentations, and surgical outcomes,” *arXiv preprint arXiv:1904.00445*, 2019.
- [3] Philippe Meyer, Vincent Noblet, Christophe Mazzara, and Alex Lallement, “Survey on deep learning for radiotherapy,” *Computers in biology and medicine*, vol. 98, pp. 126–146, 2018.
- [4] Dao-Hui Ge, Hong-Sheng Li, Liang Zhang, R Liu, P Shen, and Qi-Guang Miao, “Survey of lightweight neural network,” *J. Softw*, vol. 31, pp. 2627–2653, 2020.
- [5] Taha Emara, Hossam E Abd El Munim, and Hazem M Abbas, “Liteseg: A novel lightweight convnet for semantic segmentation,” in *2019 Digital Image Computing: Techniques and Applications (DICTA)*, 2019, pp. 1–7.
- [6] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao, “Knowledge distillation: A survey,” *International Journal of Computer Vision*, vol. 129, pp. 1789–1819, 2021.
- [7] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean, “Distilling the knowledge in a neural network,” *arXiv preprint arXiv:1503.02531*, 2015.
- [8] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho, “Relational knowledge distillation,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 3967–3976.
- [9] Jing Yang, Brais Martinez, Adrian Bulat, and Georgios Tzimiropoulos, “Knowledge distillation via softmax regression representation learning,” in *International Conference on Learning Representations*, 2020.
- [10] Yonglong Tian, Dilip Krishnan, and Phillip Isola, “Contrastive representation distillation,” *arXiv preprint arXiv:1910.10699*, 2019.
- [11] Yifan Liu, Ke Chen, Chris Liu, Zengchang Qin, Zhenbo Luo, and Jingdong Wang, “Structured knowledge distillation for semantic segmentation,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 2604–2613.
- [12] Mingi Ji, Byeongho Heo, and Sungrae Park, “Show, attend and distill: Knowledge distillation via attention-based feature matching,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021, vol. 35, pp. 7945–7952.
- [13] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio, “Fitnets: Hints for thin deep nets,” *arXiv preprint arXiv:1412.6550*, 2014.
- [14] Sergey Zagoruyko and Nikos Komodakis, “Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer,” *arXiv preprint arXiv:1612.03928*, 2016.
- [15] Nikolaos Passalis and Anastasios Tefas, “Learning deep representations with probabilistic knowledge transfer,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 268–284.
- [16] Dian Qin, Jia-Jun Bu, Zhe Liu, Xin Shen, Sheng Zhou, Jing-Jun Gu, Zhi-Hua Wang, Lei Wu, and Hui-Fen Dai, “Efficient medical image segmentation based on knowledge distillation,” *IEEE Transactions on Medical Imaging*, vol. 40, no. 12, pp. 3820–3831, 2021.
- [17] Jeongho Kim, Hanbeen Lee, and Simon S. Woo, “Imf: Integrating matched features using attentive logit in knowledge distillation,” in *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence IJCAI*, 2023.
- [18] Shiya Luo, Defang Chen, and Can Wang, “Knowledge distillation with deep supervision,” in *Proceedings of the 2023 International Joint Conference on Neural Networks (IJCNN)*, 2023, pp. 1–8.
- [19] Patrick Bilic, Patrick Christ, Hongwei Bran Li, Eugene Vorontsov, Avi Ben-Cohen, Georgios Kaissis, Adi Szeskin, Colin Jacobs, Gabriel Efrain Humpire Mamani, Gabriel Chartrand, et al., “The liver tumor segmentation benchmark (lits),” *Medical Image Analysis*, vol. 84, pp. 102680, 2023.
- [20] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang, “Unet++: A nested unet architecture for medical image segmentation,” in *Proceedings of the international conference on Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*, 2018, pp. 3–11.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.