

Pre-Trained LLM is a Semantic-Aware and Generalizable Segmentation Booster

Fenghe Tang^{1,2†}, Wenxin Ma^{1,2†}, Zhiyang He³,
Xiaodong Tao³, Zihang Jiang^{1,2} ✉, and S. Kevin Zhou^{1,2} ✉

¹ School of Biomedical Engineering, Division of Life Sciences and Medicine, University of Science and Technology of China, Hefei, Anhui, 230026, P.R. China

² Center for Medical Imaging, Robotics, Analytic Computing & Learning (MIRACLE), Suzhou Institute for Advance Research, USTC, 215123, P.R. China

³ Anhui IFlytek CO., Ltd.

fhtan9@mail.ustc.edu.cn, wxma@mail.ustc.edu.cn

<https://github.com/FengheTan9/LLM4Seg>

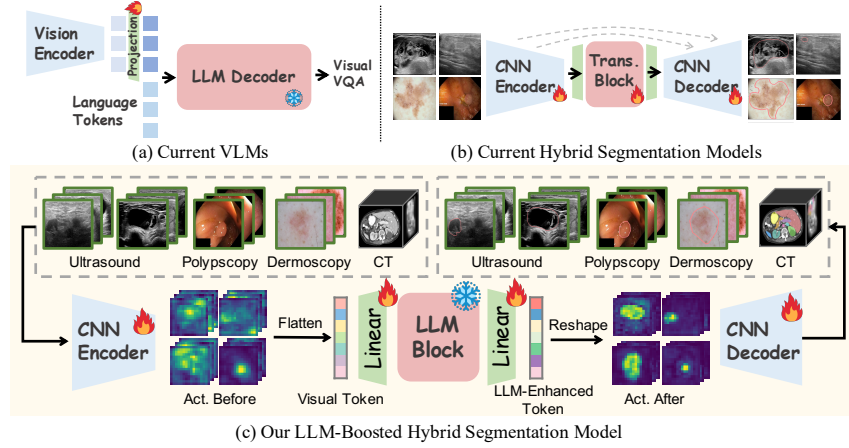


Fig. 1. Comparison of current VLMs, hybrid segmentation models, and our novel hybrid segmentation model. Our LLM4Seg use a frozen LLM layer within CNN encoder-decoder framework to boost global visual understanding.

Abstract. With the advancement of Large Language Model (LLM) for natural language processing, this paper presents an intriguing finding: a frozen pre-trained LLM layer can process visual tokens for medical image segmentation tasks. Specifically, we propose a simple hybrid structure that integrates a pre-trained, frozen LLM layer within the CNN encoder-decoder segmentation framework (LLM4Seg). Surprisingly, this design improves segmentation performance with a minimal increase in trainable parameters across various modalities, including ultrasound, dermoscopy, polyps, and CT scans. Our in-depth analysis reveals the potential of transferring LLM’s semantic awareness to enhance segmentation tasks, offering both improved global understanding and better local modeling capabilities. The improvement proves robust across different LLMs, validated using LLaMA and DeepSeek.

Keywords: Segmentation · Large Language Model · Hybrid Structure.

[†]Equal contribution. [✉]Corresponding author.

1 Introduction

Medical image segmentation is crucial for diagnosis, surgical planning, and disease monitoring. While Convolutional Neural Networks (CNNs) effectively capture local features, Vision Transformers (ViTs) have gained attention for their strong global semantic learning capabilities in segmentation tasks [7,20,2,34,13], driving the development of hybrid structures that leverage their complementary strengths to achieve better performance than pure CNN-/transformer-based methods [4,30,2,32,31]. However, it requires large-scale training data, otherwise it can fall in local shortcuts and struggle to learn robust semantic representations, posing a significant challenge in label-scarcity medical tasks [10].

Notably, Large Language Models (LLMs), built upon transformer blocks and trained on massive linguistic datasets, have revolutionized natural language processing (NLP) with their advanced reasoning, contextual understanding, and generalization capabilities [11,12]. While originally designed for NLP, recent multi-modal researches have explored LLM’s interaction with visual representations, developing various Visual Language Models (VLMs) [18,19,22]. Specifically, these methods adapt visual representations to the language space, as shown in Fig. 1(a), which typically use an LLM as a versatile decoder for both linguistic and adapted visual tokens. Surprisingly, as validated in [22], LLMs exhibit remarkable semantic understanding of visual tokens with strong generalizability. This raises an intriguing question as we investigate their potential in medical segmentation: *Is it possible to directly use pre-trained LLMs in hybrid structures to process visual tokens for medical image segmentation?*

To validate the feasibility, unlike conventional multi-modal methods that use LLMs as a shared decoder, we innovatively propose a novel hybrid model with simple yet effective hybrid structure: as illustrated in Fig. 1(c), our design integrates a pre-trained LLM layer for global modeling in a hybrid structure. It takes the CNN-encoder-extracted feature as input, projected by a linear layer, and outputs processed tokens to the CNN decoder. Given that the LLM layer is initialized with pre-trained weights and kept frozen, this design substantially reduces the reliance on large-scale data needed for training ViTs from scratch, while introducing only a minimal increase in the number of trainable parameters. Quantitative results show that this design exhibits strong generalizability and outperforms baseline architectures across various modalities and domains, including ultrasound, dermoscopy, polypscopy, and CT, establishing new state-of-the-art (SOTA) results.

Despite the promising improvement, considering that LLMs are exclusively trained on text data, *what roles does the LLM play during medical image segmentation?* Our analysis demonstrates a possibility that LLM can effectively transfer the semantic awareness learned during pretraining to enhance the understanding of visual semantics. Specifically, in segmentation, a clearer separation between the foreground and background can be found in feature activation maps with the help of an LLM layer, leading to a significant reduction in noise. This semantic refinement, likely stemming from the strong semantic processing ability of LLMs, provides more accurate guidance and enhances local-modeling ability

of CNN. Moreover, the improvement exhibits strong robustness. Regardless of input modality or structural configuration, this enhancement remains effective, highlighting the adaptability of LLMs in diverse medical segmentation tasks.

In summary, our contributions can be summarized as:

- We prove the feasibility of leveraging LLMs to process visual tokens in segmentation and propose a novel hybrid framework that incorporates a pre-trained LLM layer for global understanding in segmentation.
- Our design enhances baseline performances and achieves SOTA results across multiple medical imaging domains, including 2D modalities such as polyps, dermoscopy, and ultrasound, as well as 3D imaging such as CT scans.
- We thoroughly analyze the role of LLM during visual processing and show their potential of generalizable semantic understanding which benefits both global and local modeling in medical image segmentation.

2 Method

In this section, we introduce a novel LLM4Seg framework designed to explore the capabilities of pre-trained LLMs for processing visual tokens, with a focus on segmentation tasks. Unlike conventional approaches, our framework replaces the transformer block in a hybrid pipeline with a frozen LLM layer.

Hybrid segmentation model recap. Hybrid approaches leverage the strengths of both CNN and Transformer architectures, combining local feature extraction with global contextual modeling for improved segmentation performance. A hybrid model for segmentation tasks generally comprises two main stages: a CNN stage and a Transformer stage. Given an input image x , a CNN Encoder(\cdot) extracts position-aware activations, denoted as $t \in R^{C \times H \times W}$, which capture localized spatial features. These activations are then flattened to $t' \in R^{C \times HW}$ and fed into Transformer(\cdot) to enhance global contextual understanding. Finally, the processed tokens are reshaped back and decoded to produce the segmentation prediction \tilde{x} . The whole process is represented as:

$$t = \text{Encoder}(x; \theta_1), \quad t' = \text{Flatten}(t), \quad (1)$$

$$\hat{t} = \text{Transformer}(t'; \gamma), \quad \hat{t}' = \text{Reshape}(\hat{t}), \quad (2)$$

$$\tilde{x} = \text{Decoder}(\hat{t}'; \theta_2). \quad (3)$$

In recent state-of-the-art models [30,33,4,38], the parameters $\theta_1, \gamma, \theta_2$ are typically learned from scratch. However, identifying the optimal parameters for transformer layers often requires access to large-scale training datasets, as their capacity for modeling complex relationships is highly dependent on the volume and diversity of the data.

LLM-boosted hybrid framework. To leverage the effectiveness of LLMs in processing visual semantics, as illustrated in Fig.1, we propose a simple yet effective hybrid architecture. Specifically, we replace the transformer layer with an LLM layer, augmented by linear projection layers before and after the LLM

block. In this design, the original $\text{Transformer}(\cdot)$ in Eq. (2) is substituted with the LLM-based formulation in Eq. (4). This modification leverages the pre-trained capabilities of LLMs for a global semantic understanding while maintaining the compatibility with the existing framework.

$$\begin{aligned}\hat{t} &= \text{Linear}(t', \gamma_1), \\ \hat{t} &= \text{Transformer}(\hat{t}, \lambda_{LLM}), \\ \hat{t} &= \text{Linear}(\hat{t}, \gamma_2).\end{aligned}\tag{4}$$

By leveraging a **pre-trained and fixed** λ_{LLM} , we benefit from a stable and reliable initialization that minimizes the need for extensive retraining. This approach preserves the pre-trained knowledge with minimal modifications, allowing us to effectively assess the role of the LLM in the segmentation task while keeping the increase in computational resources minimal.

3 Experiments

3.1 Setup

Datasets. We use breast ultrasound BUSI [1], thyroid ultrasound TNSCUI [24], dermoscopy ISIC [8], and polypsopy Kvasir [25] for 2D evaluation, and abdomen CT dataset BTCV [16] for 3D evaluation. We use a 7/3 split on BUSI, TNSCUI, and ISIC for training and validation. Following previous works, both the Kvasir and BTCV datasets are partitioned into their official training and validation sets [13], with Kvasir additionally splitting a testing set for evaluation [25].

Evaluation metrics and baseline models. Following pervious work [29, 5, 36, 21], we utilize IoU and F1 score for BUSI, TNSCUI, ISIC, and Kvasir, while adopting the Dice for BTCV as [13, 35, 28]. Trainable params (M) and inference FLOPs (G) are also included for comparison. We select 11 recent SOTA models for comparison, including 2D segmentation networks: U-Net [26], U-Net++ [40], TransUNet [4], UNeXt [36] MissFormer [14], UCTransNet [37], UniRepLKNet [9] and TinyU-Net [5]; 3D segmentation networks: MedNeXt [28] and 3D UX-Net [17].

Implementation details. We use U-Net [26], CMUNeXt [29], and nnUNet [15] as 2D backbones, and MedNeXt [28], 3D UX-Net [17] as 3D backbones. By default, we empirically employ the frozen 15-th layer of LLaMA3.2-1B [11] or 28-th DeepSeek-R1-Distill-Qwen-1.5B [12] as the LLM layer. We also implement a trainable version for it, denoted as “+LLaMA(T) / +DeepSeek(T)”. For fair comparison, we follow the same parameter settings and data augmentation as prior 2D [29, 36] (256×256 as 2D inputs) and 3D [13, 35] (96×96×96 with spacing 1.5×1.5×2.0 mm as voxel inputs) works. To compare with the baseline, we use a randomly initialized transformer with exactly the same structure as the corresponding LLaMA / DeepSeek layer, denoted as “+Transformer / +Transformer”.

3.2 Quantitative Results

Comparison with other models. Surprisingly, Table 1 shows that with a frozen LLaMA layer, the segmentation performance generally improves across

Table 1. Result on BUSI, TNSCUI, ISIC, and Kvasir. Best results are highlighted as **first**, **second** and **third**. Param: trainable parameters.

Method	Computation		BUSI		TNSCUI		ISIC		Kvasir		Avg	
	Param(M)	GFLOPs	IoU	F1	IoU	F1	IoU	F1	IoU	F1	IoU	F1
TinyU-Net	0.48	1.66	66.21	75.01	74.03	82.95	81.95	88.99	86.87	92.51	77.26	84.86
UNeXt	1.47	0.58	65.04	74.16	71.04	80.46	82.10	89.93	88.01	92.05	76.54	84.15
UniRepLK	5.83	9.39	65.26	73.96	67.73	77.20	81.64	88.82	85.30	91.18	74.98	82.78
U-Net++	26.90	37.62	69.49	78.06	76.90	85.13	82.29	89.17	86.60	92.08	78.82	86.11
MissFormer	35.45	7.25	63.29	73.47	68.26	76.71	80.99	88.03	86.37	91.79	74.72	82.50
UCTransNet	66.24	32.98	70.05	78.51	75.51	84.08	82.76	89.50	87.97	92.93	79.07	86.25
TransUnet	105.3	38.52	71.39	79.85	77.63	85.76	82.95	89.68	87.95	92.84	80.03	87.25
UNet	34.52	65.56	68.79	77.00	75.99	84.24	82.28	89.19	86.01	91.49	78.26	85.48
+LLaMA	+4.19	+16.64	72.95	81.46	77.80	85.93	82.88	89.57	88.90	93.53	80.63	87.62
CMUNeXt	3.14	7.39	71.20	79.71	76.95	85.19	82.47	89.36	87.17	92.95	79.44	86.80
+LLaMA	+1.05	+16.03	73.14	81.55	77.51	85.75	82.69	89.45	89.03	93.90	80.59	87.66
nnUNet	7.76	12.12	72.76	80.72	79.53	87.34	83.07	89.51	84.97	91.15	80.08	87.18
+LLaMA	+2.10	+16.10	72.90	80.86	79.93	87.62	83.39	89.79	84.60	90.87	80.21	87.29

Table 2. Comparison with baselines. Improvements over baseline are highlighted as **> 2%**, **> 1%** and **> 0.1%**. **val** (bold): top method. Param: trainable parameters.

Methods	2D	Computation		BUSI		TNSCUI		ISIC		Kvasir		3D	BTCV Dice
		Param(M)	GFLOPs	IoU	F1	IoU	F1	IoU	F1	IoU	F1		
baseline1	UNet	34.52	65.56	68.79	77.00	75.99	84.24	82.28	89.19	86.01	91.49	MedNeXt	82.00
+Transformer		+65.01	+16.64	72.68	81.00	77.81	85.87	82.63	89.29	87.86	92.94		81.85
+LLaMA(T)		+65.01	+16.64	72.63	80.94	77.83	85.84	83.01	89.62	88.97	93.65		82.43
+LLaMA		+4.19	+16.64	72.95	81.46	77.80	85.93	82.88	89.57	88.90	93.53		82.55
baseline2	CMUNeXt	3.14	7.39	71.20	79.71	76.95	85.19	82.47	89.36	87.17	92.95	3D UX-Net	80.78
+Transformer		+61.87	+16.03	71.67	80.26	77.11	85.41	82.62	89.43	87.38	92.61		81.26
+LLaMA(T)		+61.87	+16.03	72.66	81.05	77.54	85.75	82.92	89.70	88.49	92.96		81.71
+LLaMA		+1.05	+16.03	73.14	81.55	77.51	85.75	82.69	89.45	89.03	93.90		81.84
+Transformer		+47.58	+12.38	71.32	79.89	77.28	85.47	82.59	89.41	87.80	93.05		81.23
+DeepSeek(T)		+47.58	+12.38	72.00	80.53	77.18	85.39	82.48	89.30	88.55	93.54		80.86
+DeepSeek		+0.78	+12.38	71.91	80.45	77.39	85.61	82.63	89.45	88.69	93.49		81.48

various benchmarks compared to previous methods, achieving new SOTA average results to 80.63% in average IoU (from the model “UNet + LLaMA”) and 87.66% in average F1 (from the model “CMUNeXt + LLaMA”). Notably, compared to other modality datasets, ISIC features larger foreground regions with distinct edges and lower dependence on global context. Despite this, our design still achieves notable improvements, and it achieves superior average segmentation performance while requiring significantly fewer trainable parameters.

Comparison with baselines. As presented in Table 2, the line labeled “+Transformer” represents the performance obtained by training directly on the visual dataset from scratch, using the same structure and the same number of parameters. In contrast, loading pre-trained weights from LLaMA / DeepSeek consistently improves the performance. These findings indicate that the performance gain is not attributable to an increase in parameters, but rather to the crucial role of pre-trained LLM in long-range modeling, which significantly boosts global semantic representation. Moreover, the improvement is robust across 2D and 3D datasets, showing the strong generalization ability of our findings.

Computational resources. We analyze the number of trainable parameters and GFLOPs, as shown in the “Computation” columns in Tab. 1 and Tab. 2. With the LLM layer frozen, our approach achieves the best performance while incurring only a minimal increase in computational cost. The additional trainable parameters stem solely from the projection layers before and after the transformer layer. Specifically, incorporating a LLaMA / DeepSeek layer adds only 4.19M parameters to UNet and 1.05M / 0.78M to CMUNeXt. Inference runtime increases only slightly compared to without inserting LLM layer (from 5.0ms to 5.9ms), demonstrating the method’s practicality under real-world constraints.

4 Understanding the Role of LLM in Segmentation

Despite the intriguing improvement, given that LLMs have never seen visual data during pre-training, our analysis suggests a possibility that the LLM layer can transfer the semantic knowledge acquired during pre-training to enhance the understanding of visual semantics. Simultaneously, the semantic refinement facilitates local-modeling ability of CNN. These indications are supported by: Activation Analysis in Sec. 4.1, Statistic Analysis in Sec. 4.2, and Structural Analysis in Sec. 4.3.

4.1 Activation Analysis

Activation visualization. Following [22], we visualize the feature activation maps before and after the Transformer/LLaMA layer, denoted as “Act. Before” and “Act. After”. As illustrated in Fig. 2(a), a Transformer layer trained from scratch can effectively distinguish the foreground from the background. However, the pre-trained LLaMA layer significantly reduces background noise and produces sharper boundaries, focusing more precisely on lesion regions. This enhanced activation maps underscore the strong transferability of its pre-trained knowledge to novel scenarios, even from text to visual modality, demonstrating the effectiveness of leveraging pre-trained LLaMA layer for visual understanding.

Activation concentration analysis. We further assess the accuracy of activation’s concentration by calculating the IoU between ground truth segmentation masks and the highlighted regions in “Act. Before” and “Act. After”. The results, presented in Fig. 2(b), provide a quantitative measure of how effectively the model focuses on relevant foreground regions.

Considering “Act. Before” (darker-colored columns), the concentration accuracy is improved compared to activations from CMUNeXt. Initializing the LLaMA weights further enhances the average IoU score, underscoring LLaMA’s role in guiding CNN encoding. Additionally, “Act. After” (lighter-colored columns) shows an even greater concentration accuracy compared to those before, reaching its peak with “+ LLaMA” design. This result highlights LLaMA’s high-level semantic refinement, filtering out background noise and generating more precise activations, which in turn provide reliable information for the CNN decoder.

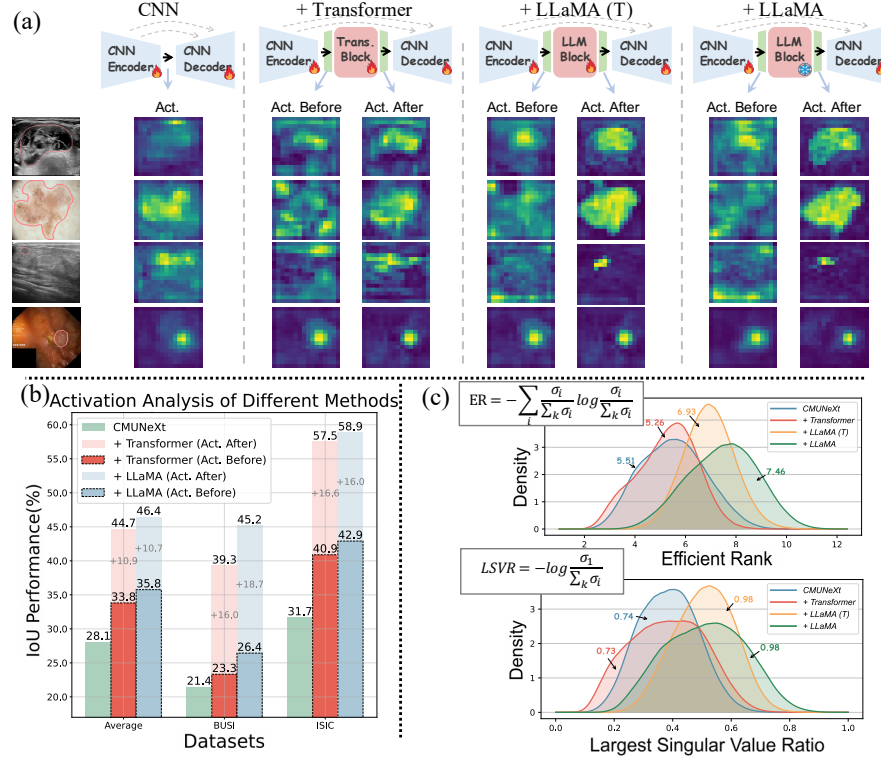


Fig. 2. (a) Activation visualization before and after the Transformer/LLaMA layer. (b) Concentration accuracy of activations. The activations are threshold by 0.4 and compared with ground truth segmentation masks to calculate IoU. (c) The distribution of Effective Rank (ER) and Largest Singular Value Ratio (LSVR) of activations extracted from TNSCUI dataset. Average value of each distribution is denoted by \rightarrow .

4.2 Statistic Analysis

Following previous works [27, 3, 23, 6, 39], we analyze the singular value spectrum of the feature space before decoding. Specifically, we perform channel-wise singular value decomposition (SVD) on the activations $\hat{t}' \in R^{C \times H \times W}$ before decoding, obtaining C singular matrix $\Sigma = [\sigma_1, \sigma_2, \dots, \sigma_k] \in R^{H \times W}$, where each σ_i represents a singular value.

We then calculate two metrics for all singular values: (1) Effective Rank (ER) [27]: the entropy of the singular values, normalized by their sum. Higher ER suggests that the feature space encompasses a greater number of dimensions, thereby more effectively capturing the underlying structure of the data. (2) Largest Singular Value Ratio (LSVR) [6, 3] represents the ratio of the largest singular value to the other singular values, calculated on a negative logarithmic scale. A smaller LSVR indicates greater dominance of the largest singular value, which might overshadow the representation ability from other singular

values and damage the model’s representation ability. As shown in Fig. 2(c), the pre-trained LLaMA layer outperforms the Transformer layer and the baseline CMUNeXt in terms of both average ER and LSVR, with the frozen LLaMA layer achieving the highest ER. This suggests a broader range of singular values of the features after LLaMA, indicating the potentially stronger representation capacity in the feature space. These results validate our suggestion that LLaMA can facilitate visual representation learning in medical image segmentation tasks.

4.3 Structural Analysis

For LLM structure, Fig. 3(Left) reveals that the improvement remains consistent across various layer configurations of LLaMA, provided that the selected layer is sufficiently deep to capture comprehensive semantics for understanding (at least after 3-th layer for LLaMA3.2-1B). Similarly, Fig. 3(Right) demonstrates that inserting a DeepSeek-R1 layer also benefits segmentation with noticeable performance gains when using deeper layers. These results further validate the generalizability of LLM’s knowledge, emphasizing that the effectiveness in medical segmentation is robust to LLM structure and layer selection.

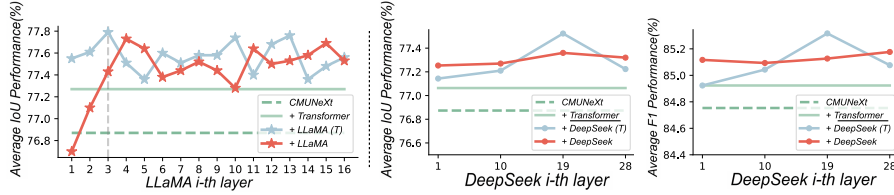


Fig. 3. Impact of different LLM Transformer layers on segmentation performance.

5 Conclusion and Discussion

Our findings reveal a novel and unexpected generalization capability of LLMs: their semantic awareness in medical image segmentation. We propose to integrate a frozen pre-trained LLM layer into a CNN architecture, demonstrating improved segmentation performance across various imaging modalities. Because the encoder is trainable, it learns to project features into the input space of the frozen LLM layer, effectively tapping into the semantic priors learned from large-scale language pretraining. Thorough analysis further validates the effectiveness of this design, highlighting the potential of LLMs in visual processing.

These insights open new avenues for multi-modal synergies between language models and vision tasks and can be potentially extended beyond segmentation to challenges like classification, detection, and anomaly identification in medical imaging. However, as shown in our structural analysis, performance fluctuations

exist, suggesting room for improvement. Future work may explore more robust adaptation strategies or diverse LLM architectures to enhance stability, optimize performance, and broaden applicability.

Acknowledgments. Supported by Natural Science Foundation of China under Grant 62271465, Suzhou Basic Research Program under Grant SYG202338.

Disclosure of Interests. The authors have no competing interests to declare that are relevant to the content of this article.

References

1. Al-Dhabyani, W., Gomaa, M., Khaled, H., Fahmy, A.: Dataset of breast ultrasound images. *Data in brief* **28**, 104863 (2020)
2. Cao, H., Wang, Y., Chen, J., Jiang, D., Zhang, X., Tian, Q., Wang, M.: Swin-unet: Unet-like pure transformer for medical image segmentation. In: *ECCV*. pp. 205–218. Springer (2022)
3. Chen, H., Wang, J., Shah, A., Tao, R., Wei, H., Xie, X., Sugiyama, M., Raj, B.: Understanding and mitigating the label noise in pre-training on downstream tasks. *arXiv preprint arXiv:2309.17002* (2023)
4. Chen, J., Lu, Y., Yu, Q., Luo, X., Adeli, E., Wang, Y., Lu, L., Yuille, A.L., Zhou, Y.: Transunet: Transformers make strong encoders for medical image segmentation. *arXiv preprint arXiv:2102.04306* (2021)
5. Chen, J., Chen, R., Wang, W., Cheng, J., Zhang, L., Chen, L.: Tinyu-net: Lighter yet better u-net with cascaded multi-receptive fields. In: *MICCAI*. pp. 626–635. Springer (2024)
6. Chen, X., Wang, S., Long, M., Wang, J.: Transferability vs. discriminability: Batch spectral penalization for adversarial domain adaptation. In: *ICML*. pp. 1081–1090. PMLR (2019)
7. Cheng, B., Misra, I., Schwing, A.G., Kirillov, A., Girdhar, R.: Masked-attention mask transformer for universal image segmentation. In: *CVPR*. pp. 1290–1299 (2022)
8. Codella, N., Rotemberg, V., Tschandl, P., Celebi, M.E., Dusza, S., Gutman, D., Helba, B., Kalloo, A., Liopyris, K., Marchetti, M., et al.: Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic). *arXiv preprint arXiv:1902.03368* (2019)
9. Ding, X., Zhang, Y., Ge, Y., Zhao, S., Song, L., Yue, X., Shan, Y.: Unireplknet: A universal perception large-kernel convnet for audio video point cloud time-series and image recognition. In: *CVPR*. pp. 5513–5524 (2024)
10. Dosovitskiy, A.: An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020)
11. Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al.: The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024)
12. Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al.: Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025)

13. Hatamizadeh, A., Tang, Y., Nath, V., Yang, D., Myronenko, A., Landman, B., Roth, H.R., Xu, D.: Unetr: Transformers for 3d medical image segmentation. In: WACV. pp. 574–584 (2022)
14. Huang, X., Deng, Z., Li, D., Yuan, X., Fu, Y.: Missformer: An effective transformer for 2d medical image segmentation. TMI **42**(5), 1484–1494 (2022)
15. Isensee, F., Jaeger, P.F., Kohl, S.A., Petersen, J., Maier-Hein, K.H.: nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. Nature methods **18**(2), 203–211 (2021)
16. Landman, B., Xu, Z., Igelsias, J., Styner, M., Langerak, T., Klein, A.: Miccai multi-atlas labeling beyond the cranial vault–workshop and challenge. In: Proc. MICCAI Challenge. vol. 5, p. 12 (2015)
17. Lee, H.H., Bao, S., Huo, Y., Landman, B.A.: 3d ux-net: A large kernel volumetric convnet modernizing hierarchical transformer for medical image segmentation. In: ICLR
18. Li, C., Wong, C., Zhang, S., Usuyama, N., Liu, H., Yang, J., Naumann, T., Poon, H., Gao, J.: Llava-med: Training a large language-and-vision assistant for biomedicine in one day. NeurIPS **36** (2024)
19. Li, J., Li, D., Savarese, S., Hoi, S.: Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In: ICML. pp. 19730–19742. PMLR (2023)
20. Lin, A., Chen, B., Xu, J., Zhang, Z., Lu, G., Zhang, D.: Ds-transunet: Dual swin transformer u-net for medical image segmentation. TIM **71**, 1–15 (2022)
21. Ma, W., Yao, Q., Zhang, X., Huang, Z., Jiang, Z., Zhou, S.K.: Towards accurate unified anomaly segmentation. arXiv preprint arXiv:2501.12295 (2025)
22. Pang, Z., Xie, Z., Man, Y., Wang, Y.X.: Frozen transformers in language models are effective visual encoder layers. arXiv preprint arXiv:2310.12973 (2023)
23. Parkhi, O.M., Vedaldi, A., Zisserman, A., Jawahar, C.: Cats and dogs. In: CVPR. pp. 3498–3505. IEEE (2012)
24. Pedraza, L., Vargas, C., Narváez, F., Durán, O., Muñoz, E., Romero, E.: An open access thyroid ultrasound image database. In: 10th International symposium on medical information processing and analysis. vol. 9287, pp. 188–193. SPIE (2015)
25. Pogorelov, K., Randel, K.R., Griwodz, C., Eskeland, S.L., de Lange, T., Johansen, D., Spampinato, C., Dang-Nguyen, D.T., Lux, M., Schmidt, P.T., et al.: Kvasir: A multi-class image dataset for computer aided gastrointestinal disease detection. In: ACMMM. pp. 164–169 (2017)
26. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: MICCAI. pp. 234–241. Springer (2015)
27. Roy, O., Vetterli, M.: The effective rank: A measure of effective dimensionality. In: 2007 15th European signal processing conference. pp. 606–610. IEEE (2007)
28. Roy, S., Koehler, G., Ulrich, C., Baumgartner, M., Petersen, J., Isensee, F., Jaeger, P.F., Maier-Hein, K.H.: Mednext: transformer-driven scaling of convnets for medical image segmentation. In: MICCAI. pp. 405–415. Springer (2023)
29. Tang, F., Ding, J., Quan, Q., Wang, L., Ning, C., Zhou, S.K.: Cmunext: An efficient medical image segmentation network based on large kernel and skip fusion. In: ISBI. pp. 1–5. IEEE (2024)
30. Tang, F., Nian, B., Ding, J., Quan, Q., Yang, J., Liu, W., Zhou, S.K.: Mobileutr: Revisiting the relationship between light-weight cnn and transformer for efficient medical image segmentation. arXiv preprint arXiv:2312.01740 (2023)
31. Tang, F., Nian, B., Li, Y., Jiang, Z., Yang, J., Liu, W., Zhou, S.K.: Mambamim: Pre-training mamba with state space token interpolation and its application to medical image segmentation. Medical Image Analysis p. 103606 (2025)

32. Tang, F., Xu, R., Yao, Q., Fu, X., Quan, Q., Zhu, H., Liu, Z., Zhou, S.K.: Hyspark: Hybrid sparse masking for large scale medical image pre-training. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 330–340. Springer (2024)
33. Tang, F., Xu, R., Yao, Q., Fu, X., Quan, Q., Zhu, H., Liu, Z., Zhou, S.K.: Hyspark: Hybrid sparse masking for large scale medical image pre-training. In: MICCAI. pp. 330–340. Springer (2024)
34. Tang, F., Yao, Q., Ma, W., Wu, C., Jiang, Z., Zhou, S.K.: Hi-end-mae: Hierarchical encoder-driven masked autoencoders are stronger vision learners for medical image segmentation. arXiv preprint arXiv:2502.08347 (2025)
35. Tang, Y., Yang, D., Li, W., Roth, H.R., Landman, B., Xu, D., Nath, V., Hatamizadeh, A.: Self-supervised pre-training of swin transformers for 3d medical image analysis. In: CVPR. pp. 20730–20740 (2022)
36. Valanarasu, J.M.J., Patel, V.M.: Unext: Mlp-based rapid medical image segmentation network. In: MICCAI. pp. 23–33. Springer (2022)
37. Wang, H., Cao, P., Wang, J., Zaiane, O.R.: Uctransnet: rethinking the skip connections in u-net from a channel-wise perspective with transformer. In: AAAI. vol. 36, pp. 2441–2449 (2022)
38. Wenxuan, W., Chen, C., Meng, D., Hong, Y., Sen, Z., Jiangyun, L.: Transbts: Multimodal brain tumor segmentation using transformer. In: MICCAI. pp. 109–119 (2021)
39. Xue, Y., Whitecross, K., Mirzasoleiman, B.: Investigating why contrastive learning benefits robustness against label noise. In: ICML. pp. 24851–24871. PMLR (2022)
40. Zhou, Z., Siddiquee, M.M.R., Tajbakhsh, N., Liang, J.: Unet++: Redesigning skip connections to exploit multiscale features in image segmentation. TMI **39**(6), 1856–1867 (2019)