# Toward Next-generation Medical Vision Backbones: Modeling Finer-grained Long-range Visual Dependency

Mingyuan Meng

School of Computer Science, The University of Sydney, Sydney, Australia
`mmen2292@uni.sydney.edu.au`

**Abstract.** Medical Image Computing (MIC) is a broad research topic covering both pixel-wise (e.g., segmentation, registration) and image-wise (e.g., classification, regression) vision tasks. Effective analysis demands models that capture both global long-range context and local subtle visual characteristics, necessitating fine-grained long-range visual dependency modeling. Compared to Convolutional Neural Networks (CNNs) that are limited by intrinsic locality, transformers excel at long-range modeling; however, due to the high computational loads of self-attention, transformers typically cannot process high-resolution features (e.g., full-scale image features before downsampling or patch embedding) and thus face difficulties in modeling fine-grained dependency among subtle medical image details. Concurrently, Multi-layer Perceptron (MLP)-based visual models are recognized as computation/memory-efficient alternatives in modeling long-range visual dependency but have yet to be widely investigated in the MIC community. This doctoral research advances deep learning-based MIC by investigating effective long-range visual dependency modeling. It first presents innovative use of transformers for both pixel- and image-wise medical vision tasks. The focus then shifts to MLPs, pioneeringly developing MLP-based visual models to capture fine-grained long-range visual dependency in medical images. Extensive experiments confirm the critical role of long-range dependency modeling in MIC and reveal a key finding: MLPs provide feasibility in modeling finer-grained long-range dependency among higher-resolution medical features containing enriched anatomical/pathological details. This finding establishes MLPs as a superior paradigm over transformers/CNNs, consistently enhancing performance across various medical vision tasks and paving the way for next-generation medical vision backbones.

## 1 Research Problem and Motivation

Modern medical image computing increasingly prioritizes network backbones capable of modeling long-range visual dependency, i.e., the ability to capture relationships between distant anatomical regions or features within or across images. This capability is essential as clinically significant visual patterns often span large areas or involve spatially dispersed structures [1]. Effective modeling
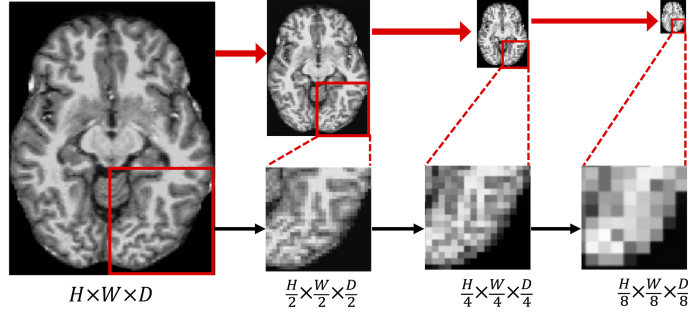
**Fig. 1.** Illustration of a medical image after downsampling. A 2D slice of a 3D brain MRI scan shows white matter and grey matter (cortex) over the brain surface; the white and grey matter are blurred after downsampling and cannot be differentiated.

enhances global context awareness, improving identification of subtle abnormalities within complex anatomy and boosting robustness to anatomical variability.

Convolutional Neural Networks (CNNs) initially dominated medical image computing due to their hierarchical feature extraction via translation-invariant convolutions. Their inherent strengths, such as local connectivity, weight sharing, and spatial hierarchy learning, are suitable for imaging data. However, their intrinsic locality fundamentally limits long-range dependency capture due to constrained receptive fields and absent global connectivity [2].

Transformers have recently emerged as a powerful alternative that leverages self-attention mechanisms from natural language processing [3]. By enabling interactions between all input tokens (e.g., image patches), they effectively model global relationships, driving rapid adoption in medical imaging. However, the prohibitive computational costs of self-attention operations, especially for 3D medical imaging, remain key limitations [4], incurring difficulties in processing high-resolution image features to capture fine-grained visual dependency [5]. While efficient variants (e.g., Swin Transformers [6]) mitigate computational loads, they still fail to leverage the tissue-level textural information that is only available at high resolutions. Unfortunately, tissue-level textural details are indispensable for accurate medical image comprehension. Unlike natural images, medical images depict standardized anatomical regions across patients, resulting in high structural similarity with subtle anatomical/pathological characteristics discernible only in high resolutions containing tissue-level image textural details. As exemplified in Fig. 1, downsampled brain MRI lacks critical anatomical distinctions, e.g., gray and white matter boundaries.

Under this context, Multi-layer Perceptrons (MLPs) serve as a promising frontier. MLP-based visual models can efficiently capture long-range visual dependency without costly self-attention operations[7,8]. Their efficiency enables modeling of fine-grained long-range dependency among high-resolution features with critical subtle anatomical/pathological details. MLP-based models have been explored and achieved promising performance in natural image tasks, while

their potential for medical image computing has not been fully recognized, as existing models lack the consideration of inductive bias (also known as learning bias) crucial for medical image computing.

This doctoral research [9] aims to uncover the fundamental mechanisms by which long-range visual dependency enhances deep learning-based medical image computing, thereby advancing the field through more effective long-range dependency modeling. This research pioneers novel transformer- and MLP-based network architectures to effectively model long-range visual dependency and rigorously evaluates them across diverse medical image computing tasks.

## 2    Background

Medical image computing has become a pivotal component of modern healthcare, fueled by the critical role of medical imaging in disease diagnosis and prognosis. Medical imaging technologies, such as Magnetic Resonance Imaging (MRI), Computed Tomography (CT), Positron Emission Tomography (PET), and X-rays, generate detailed visualizations of internal anatomy and physiological processes, which are indispensable for detecting abnormalities, evaluating disease severity, and designing personalized treatment strategies. The primary objective of medical image computing is to derive diagnostic and prognostic insights from medical images, empowering clinicians to make precise decisions and develop personalized treatment plans [10,11]. This field encompasses diverse vision tasks, which can be broadly categorized as pixel-wise and image-wise vision tasks: Pixel-wise (also known as dense prediction) tasks require pixel-level predictions for fine-grained analysis, including anatomical structure segmentation [12] or image registration [13], while image-wise tasks derive holistic interpretations from medical images for disease classification [14] or outcome prediction [15]. The clinical importance and widespread nature of medical imaging have motivated intense research and clinical efforts on computational medical image analysis.

In recent years, deep learning has established itself as a transformative force in medical image computing, demonstrating remarkable success across both pixel-wise and image-wise vision tasks [16]. Unlike traditional machine learning approaches that demand handcrafted feature engineering and domain-specific expertise, deep learning automatically extracts high-level pattern representations directly from medical images through deep neural networks. This capability eliminates human bias inherent in handcrafted feature engineering while uncovering clinically relevant semantic features potentially missed by manually-defined feature extraction [17]. The advancement of deep learning has been propelled by increased computational power, expanded datasets, and novel network architectures. Despite these developments, significant architectural challenges persist for medical applications. Medical images exhibit substantial variability across patients and pathologies, necessitating architectures with global perception capabilities to contextualize anatomical structures and suppress irrelevant local variations. Simultaneously, the critical diagnostic importance of subtle textural details demands precise localized perception to capture fine-grained anatomical

and pathological features. Therefore, the pursuit of deep learning-based medical image computing has driven significant architectural evolution, especially in visual backbones capable of long-range visual dependency modeling.

Transformers, with global self-attention mechanisms, have demonstrated substantial success in various medical vision tasks. For medical image segmentation, deep learning models such as TransUNet [18] integrated vision transformers into UNet bottlenecks to enhance contextual awareness. For deformable registration, TransMorph [19] leveraged Swin transformers [6] to capture spatial relationships across images. Despite these advances, transformers face persistent challenges: their computational complexity remains prohibitive for high-resolution 3D images, and even optimized variants (e.g., Swin transformer) struggle to preserve fine-grained anatomical details in full-resolution feature maps. In image-wise tasks such as survival prediction, transformer-based models, e.g., MCTA [20], capture prognostic patterns but often inadequately address multi-modal fusion, relying on simplistic early fusion rather than cross-modal interaction.

Concurrently, MLP-based visual models have emerged as efficient alternatives for long-range dependency modeling. Early models (e.g., MLP-Mixer [8]) excelled in classification but lacked multi-scale support for pixel-wise dense prediction. Later hierarchical MLP-based models (e.g., Hire-MLP [21]) enabled pixel-wise tasks, while MAXIM [22] achieved state-of-the-art natural image processing performance by handling high-resolution features. For medical images, the adoption of MLPs remains limited. UNeXt [23] accelerated medical image segmentation with token-shifted MLPs, while a few studies explored other tasks such as registration [24] and reconstruction [25]. Existing approaches adopted generic MLP models without domain-specific adaptations, lacking the consideration of inductive bias (i.e., initial assumptions about the data to be analyzed/generalized) that is crucial for target tasks. Further, existing approaches initiate MLP processing after feature downsampling (e.g., 4×4 patch embedding), discarding tissue-level textures essential for precise medical image comprehension. These limitations highlight the unmet need for medical-optimized MLP frameworks capturing fine-grained long-range dependency at high resolutions.

## 3   Scientific Approach

As illustrated in Fig. 2, this doctoral research advances long-range visual dependency modeling for medical image computing through four studies:

**Study I: Transformers for medical image registration** - This study aims at (i) developing a new transformer-based method for medical image registration and (ii) exploring new insights into how the modeling of long-range visual dependency works in the scenario of image registration. To attain these aims, a transformer-based medical image registration framework (named NICE-Trans) was proposed to address the unmet need for coarse-to-fine dependency modeling. Distinct from previous frameworks, NICE-Trans embeds transformers into a non-iterative coarse-to-fine registration framework to model long-range relevance between images. This framework enables joint modeling of affine and
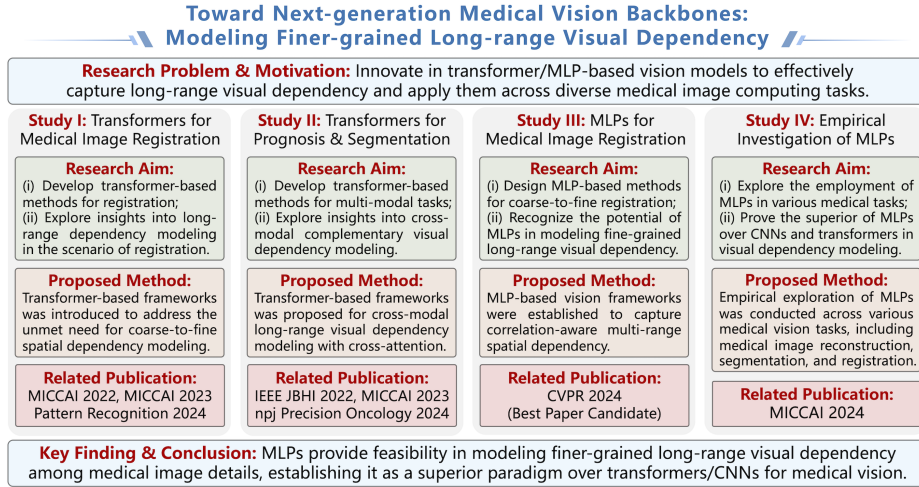
**Toward Next-generation Medical Vision Backbones:**
**Modeling Finer-grained Long-range Visual Dependency**

**Research Problem & Motivation:** Innovate in transformer/MLP-based vision models to effectively capture long-range visual dependency and apply them across diverse medical image computing tasks.

| **Study I:** Transformers for Medical Image Registration | **Study II:** Transformers for Prognosis & Segmentation | **Study III:** MLPs for Medical Image Registration | **Study IV:** Empirical Investigation of MLPs |
|---|---|---|---|
| **Research Aim:** (i) Develop transformer-based methods for registration; (ii) Explore insights into long-range dependency modeling in the scenario of registration. | **Research Aim:** (i) Develop transformer-based methods for multi-modal tasks; (ii) Explore insights into cross-modal complementary visual dependency modeling. | **Research Aim:** (i) Design MLP-based methods for coarse-to-fine registration; (ii) Recognize the potential of MLPs in modeling fine-grained long-range visual dependency. | **Research Aim:** (i) Explore the employment of MLPs in various medical tasks; (ii) Prove the superior of MLPs over CNNs and transformers in visual dependency modeling. |
| **Proposed Method:** Transformer-based frameworks was introduced to address the unmet need for coarse-to-fine spatial dependency modeling. | **Proposed Method:** Transformer-based frameworks was proposed for cross-modal long-range visual dependency modeling with cross-attention. | **Proposed Method:** MLP-based vision frameworks were established to capture correlation-aware multi-range spatial dependency. | **Proposed Method:** Empirical exploration of MLPs was conducted across various medical vision tasks, including medical image reconstruction, segmentation, and registration. |
| **Related Publication:** MICCAI 2022, MICCAI 2023 Pattern Recognition 2024 | **Related Publication:** IEEE JBHI 2022, MICCAI 2023 npj Precision Oncology 2024 | **Related Publication:** CVPR 2024 (Best Paper Candidate) | **Related Publication:** MICCAI 2024 |

**Key Finding & Conclusion:** MLPs provide feasibility in modeling finer-grained long-range visual dependency among medical image details, establishing it as a superior paradigm over transformers/CNNs for medical vision.

**Fig. 2.** Overview of this doctoral research.

deformable transformations while capturing long-range spatial relevance across resolutions, establishing new benchmarks in registration tasks. The research related to this study has been published at *MICCAI 2022* [26], *MICCAI 2023* [27], and *Pattern Recognition* [28].

**Study II: Transformers for survival prediction and segmentation** - This study aims at (i) developing a new transformer-based method for joint survival prediction and tumor segmentation from multi-modality PET-CT images and (ii) exploring new insights into how transformers benefit representation learning by modeling long-range complementary dependency between multi-modality medical images. To attain these aims, a merging-diverging hybrid transformer network (named XSurv) was proposed, where a Hybrid Parallel Cross-Attention (HPCA) block was introduced to effectively integrate complementary PET-CT information via cross-attention transformers. This model enables extracting modality-specific prognostic features while sharing contextual knowledge across modalities. The research related to this study has been published at *IEEE JBHI* [29], *MICCAI 2023* [30], and *npj Precision Oncology* [31].

**Study III: MLPs for medical image registration** - This study aims at (i) developing the first MLP-based method for medical image registration and (ii) revealing the unnoticed potential of MLPs in modeling fine-grained long-range visual dependency for pixel-wise image registration tasks. To attain these aims, a correlation-aware coarse-to-fine MLP-based network (named CorrMLP) was proposed, which introduces the first MLP block that was designed to model correlation-aware multi-range visual dependency for medical image registration. This model unlocks the potential of MLPs for capturing pixel-wise spatial dependency. This study has been published at *CVPR 2024* [32] and was nominated as the *Best Paper Candidate* (Top 24/0.2%).

**Study IV: Empirical investigation of MLPs in medical dense prediction** - This study aims at (i) exploring the advantages of MLPs in modeling fine-grained long-range visual dependency among high-resolution image features and (ii) validating their superiority over transformers in medical image computing. To attain these aims, a comprehensive empirical study was conducted to explore the employment of MLPs on various pixel-wise prediction tasks, including medical image reconstruction, registration, and segmentation. This empirical study also introduces a new feature extraction framework that produces hierarchical feature pyramids with fine-grained long-range dependency modeling and can be generalized to a wide range of medical dense prediction tasks. The research innovating in MLP-based models for medical image reconstruction has been published at *MICCAI 2024* [25].

## 4   Proposed Solution

Study I introduces NICE-Trans, a non-iterative coarse-to-fine transformer network that unifies affine and deformable medical image registration within a single network. This study makes two technical advances: (i) It extended the coarse-to-fine registration paradigm to jointly model traditionally separated affine and deformable transformations in one network iteration, eliminating multi-stage processing. (ii) It innovatively embedded transformers to progressively capture long-range spatial relevance between images, marking the first integration of transformers into non-iterative coarse-to-fine registration. These dual advances establish NICE-Trans as a pioneering model to unify affine and deformable coarse-to-fine registration while capturing long-range spatial dependency.

Study II introduces XSurv, an X-shaped merging-diverging hybrid transformer network for joint survival prediction and tumor segmentation. The architecture of XSurv features a merging encoder that fuses complementary anatomical (from CT) and metabolic (from PET) information, and a diverging decoder that extracts region-specific prognostic features from primary tumor and metastatic lymph node regions. The technical innovations comprise: (i) a specialized merging-diverging learning framework for joint survival prediction and tumor segmentation, enabling multi-modality exploitation and region-specific feature extraction applicable across survival tasks; (ii) a Hybrid Parallel Cross-Attention (HPCA) block that concurrently learns local intra-modality features through convolutional pathways and global inter-modality dependencies via cross-attention transformers; and (iii) a Region-specific Attention Gate (RAG) that filters lesion-relevant features through anatomical screening.

Study III introduces CorrMLP, the first MLP-based coarse-to-fine deformable registration framework. Its core innovation is the correlation-aware multi-window MLP (CMW-MLP) block, a purpose-built module that computes local feature correlations and captures multi-range spatial dependency through parallel MLP operations. This block was embedded into a novel correlation-aware architecture leveraging both image-level (inter-image feature relationships) and step-level (inter-scale contextual propagation) correlations. This dual-correlation

mechanism provides enriched contextual guidance throughout the coarse-to-fine registration process. CorrMLP thus pioneers three technical contributions: (i) establishing MLPs as promising backbones for deformable registration; (ii) introducing the first registration-optimized MLP block with explicit correlation modeling; and (iii) devising a contextual registration framework where both image-level and step-level correlations actively guide registration.

Study IV uncovers the underexplored capability of MLPs in capturing fine-grained long-range dependency in high-resolution image features, a critical advantage for medical dense prediction. Through a comprehensive empirical investigation, a hierarchical MLP framework was introduced, which extracts multi-scale feature pyramids via hierarchical MLP blocks operating beginning from the full image resolution. Task-specific decoders then leverage these feature pyramids for various medical applications, including medical image reconstruction, deformable registration, and segmentation. Crucially, when evaluating various MLP blocks within this framework, a paradigm-shifting finding was observed: regardless of the specific MLP variants, employing MLPs at full resolution consistently enabled superior performance over CNN- and transformer-based methods across all evaluation tasks, even outperforming task-optimized specialist models.

## 5   Results and Contribution

The methodologies proposed in this research were extensively validated across multi-modal medical imaging datasets (including MRI, CT, and PET) and body regions (including brain, cardiac, and head and neck regions), leveraging public benchmarks [33,34,35,36] to ensure reproducibility and clinical relevance.

In Study I, the proposed NICE-Trans achieved state-of-the-art performance in medical image registration by unifying affine and deformable coarse-to-fine registration within a single non-iterative network [27]. In quantitative evaluations, NICE-Trans attained the best Dice Similarity Coefficient (DSC) results across all evaluation datasets, outperforming existing transformer-based methods and CNN-based coarse-to-fine methods. Qualitatively, its alignments show superior anatomical consistency with fixed images. The architecture's joint affine-deformable design proves computationally optimal: affine registration incurs negligible runtime overhead compared to standalone affine methods, while end-to-end training reduces GPU memory burdens versus multi-network pipelines. Ablation studies reveal that embedding transformers in the decoder, not the encoder, drives performance gains by modeling inter-image spatial relevance rather than intra-image representations. This insight counters prior transformer-based registration methods (e.g., TransMorph) and establishes a new principle: transformers in the decoder maximize registration efficacy by explicitly modeling long-range spatial correspondence to handle large deformations between images.

In Study II, the proposed XSurv achieved state-of-the-art performance in head and neck cancer, attaining the highest C-index for survival prediction on the HECKTOR 2022 challenge dataset while simultaneously outperforming multi-task models in both survival prediction and tumor segmentation [30]. This dual

superiority stems from its novel architecture: the merging encoder with Hybrid Parallel Cross-Attention (HPCA) blocks solves critical PET-CT fusion limitations by concurrently extracting intra-modality features and inter-modality relevance, outperforming typical early and late fusion strategies; while the diverging decoder with Region-specific Attention Gates (RAG) achieves precise localization of primary tumors and metastatic lymph nodes, evidenced by attention maps and overall best segmentation DSC among all comparison methods. Crucially, XSurv eliminates reliance on manual segmentation during prognosis, a significant advantage over traditional radiomics-based methods requiring anatomical priors. This integrated approach establishes a new paradigm for joint multi-modal survival modeling, where optimized feature interaction and region-specific decoding drive superior prognostic performance.

In Study III, the proposed CorrMLP established state-of-the-art deformable registration performance, exceeding existing CNN-based, transformer-based, and coarse-to-fine registration methods in both brain and cardiac image datasets [32]. It achieved significantly higher DSC while maintaining competitive transformation smoothness and real-time GPU processing ($<$1s per image pair), with qualitative results demonstrating exceptional anatomical alignment. Three innovations drive these performance gains: (i) The proposed correlation-aware coarse-to-fine framework leverages image-level feature relationships and step-level contextual propagation; (ii) Full-resolution MLP processing outperforms CNNs/transformers by capturing fine-grained dependency in high-resolution features; and (iii) The CMW-MLP block employs optimized multi-window operations ($3\times3\times3$, $5\times5\times5$, $7\times7\times7$) to jointly model subtle and large deformations, outperforming five existing MLP variants.

In Study IV, the conducted empirical study validated MLPs as superior vision backbones for medical dense prediction, demonstrating state-of-the-art performance across reconstruction, registration, and segmentation tasks [5]. For low-dose PET reconstruction, MLP-based MLP-Unet surpassed transformer-based methods by 3.2-8.7% across metrics, with the largest gains in ultra-low dose conditions where sparse textures challenge conventional methods. For deformable registration, MLP-based MLPMorph outperformed transformer-based methods by 2.1-4.3% DSC while maintaining real-time speed, resolving registration's critical reliance on high-resolution textural details. For segmentation, MLP-based MLP-Unet also achieved the highest DSC, excelling particularly at boundary delineation as shown in qualitative analysis. Crucially, ablation studies revealed two paradigm-shifting insights: (i) The advantage of MLPs stems from high-resolution long-range dependency modeling: replacing first-stage MLPs with convolutions degraded all tasks, while transformers failed at full resolution due to GPU constraints; (ii) Performance gains are architecture-agnostic: four distinct MLP blocks outperformed transformers when applied at high resolution, proving the inherent efficacy of MLPs over transformers for fine-grained dependency modeling. This study might redefine the backbone selection principle: under the same computational constraints, MLPs can unlock finer-grained long-range dependency modeling for pixel-wise medical analysis.

# 6    Open Challenges and Future Work

It is hoped that this doctoral research can inspire new discussions on the employment of MLPs in medical image computing and motivate the community to recognize the potential of MLPs as a superior technical paradigm over CNNs and transformers for capturing finer-grained long-range visual dependency in medical images. However, despite the technical advances in this doctoral research, key barriers remain for clinical adoption: Firstly, future work could prioritize multi-site robustness validation. Large-scale clinical trials using heterogeneous multi-center data are essential to verify performance under real-world variability. Secondly, addressing data bias and fairness is critical. Future efforts could focus on curating inclusive datasets and developing fairness-aware adaptations of MLP networks to ensure equitable performance across populations. Finally, model interpretability remains a critical issue for clinical trustworthiness. Future work could pioneer task-specific interpretability methods, e.g., visualizing fine-grained dependency maps in MLP-based frameworks.

# 7    Long Term Goals

The long-term goals are to (i) motivate the community to rethink the critical roles of long-range vision dependency modeling in medical image computing and (ii) drive a paradigm shift in medical vision backbones by advancing MLPs as a superior alternative to transformers in long-range visual dependency modeling.

# References

1. Li, J., Chen, J., Tang, Y., Wang, C., Landman, B.A., Zhou, S.K.: Transforming medical imaging with transformers? a comparative review of key properties, current progresses, and future perspectives. Medical image analysis **85**, 102762 (2023)
2. Li, Y., Zhang, K., Cao, J., Timofte, R., Van Gool, L.: Localvit: Bringing locality to vision transformers. arXiv preprint arXiv:2104.05707 (2021)
3. Vaswani, A., et al.: Attention is all you need. Advances in neural information processing systems **30** (2017)
4. Shamshad, F., et al.: Transformers in medical imaging: A survey. Medical image analysis **88**, 102802 (2023)
5. Meng, M., Xue, Y., Feng, D., Bi, L., Kim, J.: Full-resolution mlps empower medical dense prediction. arXiv preprint arXiv:2311.16707 (2023)
6. Liu, Z., et al.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 10012–10022 (2021)
7. Liu, R., Li, Y., Tao, L., Liang, D., Zheng, H.T.: Are we ready for a new paradigm shift? a survey on visual deep mlp. Patterns **3**(7) (2022)
8. Tolstikhin, I.O., et al.: Mlp-mixer: An all-mlp architecture for vision. Advances in neural information processing systems **34**, 24261–24272 (2021)
9. Meng, M.: Modeling Fine-grained Long-range Visual Dependency for Deep Learning-based Medical Image Analysis. Ph.D. thesis, The University of Sydney (2025)

10. Gu, B., et al.: Multi-task deep learning-based radiomic nomogram for prognostic prediction in locoregionally advanced nasopharyngeal carcinoma. European journal of nuclear medicine and molecular imaging **50**(13), 3996–4009 (2023)
11. Gu, B., et al.: Prediction of 5-year progression-free survival in advanced nasopharyngeal carcinoma with pretreatment pet/ct using multi-modality deep learning-based radiomics. Frontiers in oncology **12**, 899351 (2022)
12. Ye, S., et al.: Enabling text-free inference in language-guided segmentation of chest x-rays via self-guidance. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 242–252. Springer (2024)
13. Meng, M., Bi, L., Fulham, M., Feng, D.D., Kim, J.: Enhancing medical image registration via appearance adjustment networks. NeuroImage **259**, 119444 (2022)
14. Li, M., et al.: Enhancing medical vision-language contrastive learning via inter-matching relation modelling. IEEE Transactions on Medical Imaging **44**(6), 2463–2476 (2025)
15. Meng, M., Peng, Y., Bi, L., Kim, J.: Multi-task deep learning for joint tumor segmentation and outcome prediction in head and neck cancer. In: 3D Head and Neck Tumor Segmentation in PET/CT Challenge, pp. 160–167. Springer (2021)
16. Suganyadevi, S., Seethalakshmi, V., Balasamy, K.: A review on deep learning in medical image analysis. International Journal of Multimedia Information Retrieval **11**(1), 19–38 (2022)
17. Hosny, A., et al.: Handcrafted versus deep learning radiomics for prediction of cancer therapy response. The Lancet Digital Health **1**(3), e106–e107 (2019)
18. 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)
19. Chen, J., et al.: Transmorph: Transformer for unsupervised medical image registration. Medical image analysis **82**, 102615 (2022)
20. Chen, R.J., Lu, M.Y., Weng, W.H., Chen, T.Y., Williamson, D.F., Manz, T., Shady, M., Mahmood, F.: Multimodal co-attention transformer for survival prediction in gigapixel whole slide images. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 4015–4025 (2021)
21. Guo, J., Tang, Y., Han, K., Chen, X., Wu, H., Xu, C., Xu, C., Wang, Y.: Hire-mlp: Vision mlp via hierarchical rearrangement. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 826–836 (2022)
22. Tu, Z., Talebi, H., Zhang, H., Yang, F., Milanfar, P., Bovik, A., Li, Y.: Maxim: Multi-axis mlp for image processing. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 5769–5780 (2022)
23. Valanarasu, J.M.J., Patel, V.M.: Unext: Mlp-based rapid medical image segmentation network. In: International conference on medical image computing and computer-assisted intervention. pp. 23–33. Springer (2022)
24. Wang, Z., et al.: Unsupervised echocardiography registration through patch-based mlps and transformers. In: International Workshop on Statistical Atlases and Computational Models of the Heart. pp. 168–178. Springer (2022)
25. Li, X., Meng, M., et al.: 3dpx: Progressive 2d-to-3d oral image reconstruction with hybrid mlp-cnn networks. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 25–34. Springer (2024)
26. Meng, M., Bi, L., Feng, D., Kim, J.: Non-iterative coarse-to-fine registration based on single-pass deep cumulative learning. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 88–97. Springer (2022)

27. Meng, M., et al.: Non-iterative coarse-to-fine transformer networks for joint affine and deformable image registration. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 750–760. Springer (2023)
28. Meng, M., et al.: Autofuse: Automatic fusion networks for deformable medical image registration. Pattern Recognition **161**, 111338 (2025)
29. Meng, M., Gu, B., Bi, L., Song, S., Feng, D.D., Kim, J.: Deepmts: Deep multi-task learning for survival prediction in patients with advanced nasopharyngeal carcinoma using pretreatment pet/ct. IEEE Journal of Biomedical and Health Informatics **26**(9), 4497–4507 (2022)
30. Meng, M., et al.: Merging-diverging hybrid transformer networks for survival prediction in head and neck cancer. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 400–410. Springer (2023)
31. Meng, M., Gu, B., Fulham, M., Song, S., Feng, D., Bi, L., Kim, J.: Adaptive segmentation-to-survival learning for survival prediction from multi-modality medical images. NPJ Precision Oncology **8**(1), 232 (2024)
32. Meng, M., Feng, D., Bi, L., Kim, J.: Correlation-aware coarse-to-fine mlps for deformable medical image registration. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9645–9654 (2024)
33. Eisenmann, M., et al.: Biomedical image analysis competitions: The state of current participation practice. arXiv preprint arXiv:2212.08568 (2022)
34. Baheti, B., et al.: The brain tumor sequence registration (brats-reg) challenge: Establishing correspondence between pre-operative and follow-up mri scans of diffuse glioma patients. arXiv preprint arXiv:2112.06979 (2021)
35. Meng, M., Bi, L., Feng, D., Kim, J.: Radiomics-enhanced deep multi-task learning for outcome prediction in head and neck cancer. In: 3D Head and Neck Tumor Segmentation in PET/CT Challenge, pp. 135–143. Springer (2022)
36. Meng, M., Bi, L., Feng, D., Kim, J.: Brain tumor sequence registration with non-iterative coarse-to-fine networks and dual deep supervision. In: International MICCAI Brainlesion Workshop. pp. 273–282. Springer (2022)