

ClinicalFMamba: Advancing Clinical Assessment using Mamba-based Multimodal Neuroimaging Fusion

Meng Zhou^{1,*} and Farzad Khalvati^{2,3}

¹TD Bank Group, Toronto, Canada ²Department of Computer Science, University of Toronto ³Department of Medical Imaging, University of Toronto *Work done while at University of Toronto MLMI 2025 Workshop

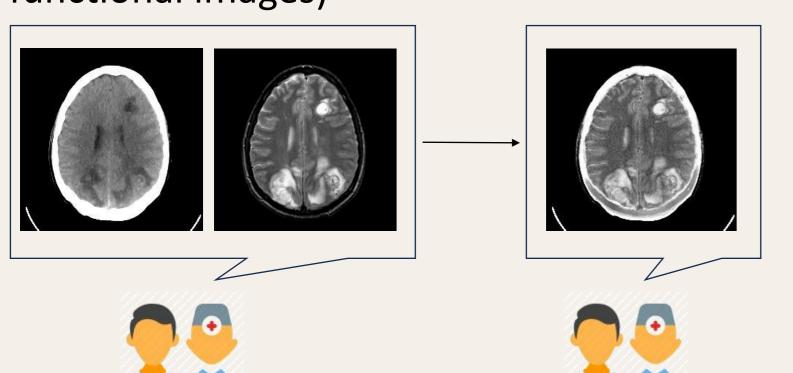




INTODUCTION

Context:

- Multimodal image fusion (MMIF) plays an increasingly prominent role in clinical diagnosis.
- > Aggregate complementary information from different image modalities to produce higherquality fused images (e.g., anatomical and functional images)



Challenges:

- CNN based methods [1] are limited by their inherent local receptive fields, which restrict their ability to capture long-range spatial dependencies.
- > Transformer based methods [2,3] create prohibitive costs for clinical applications with large images due to its self-attention mechanisms, limiting the practical deployment capabilities.
- 3D medical image fusion and the clinical applicability of fused images are still underexplored.

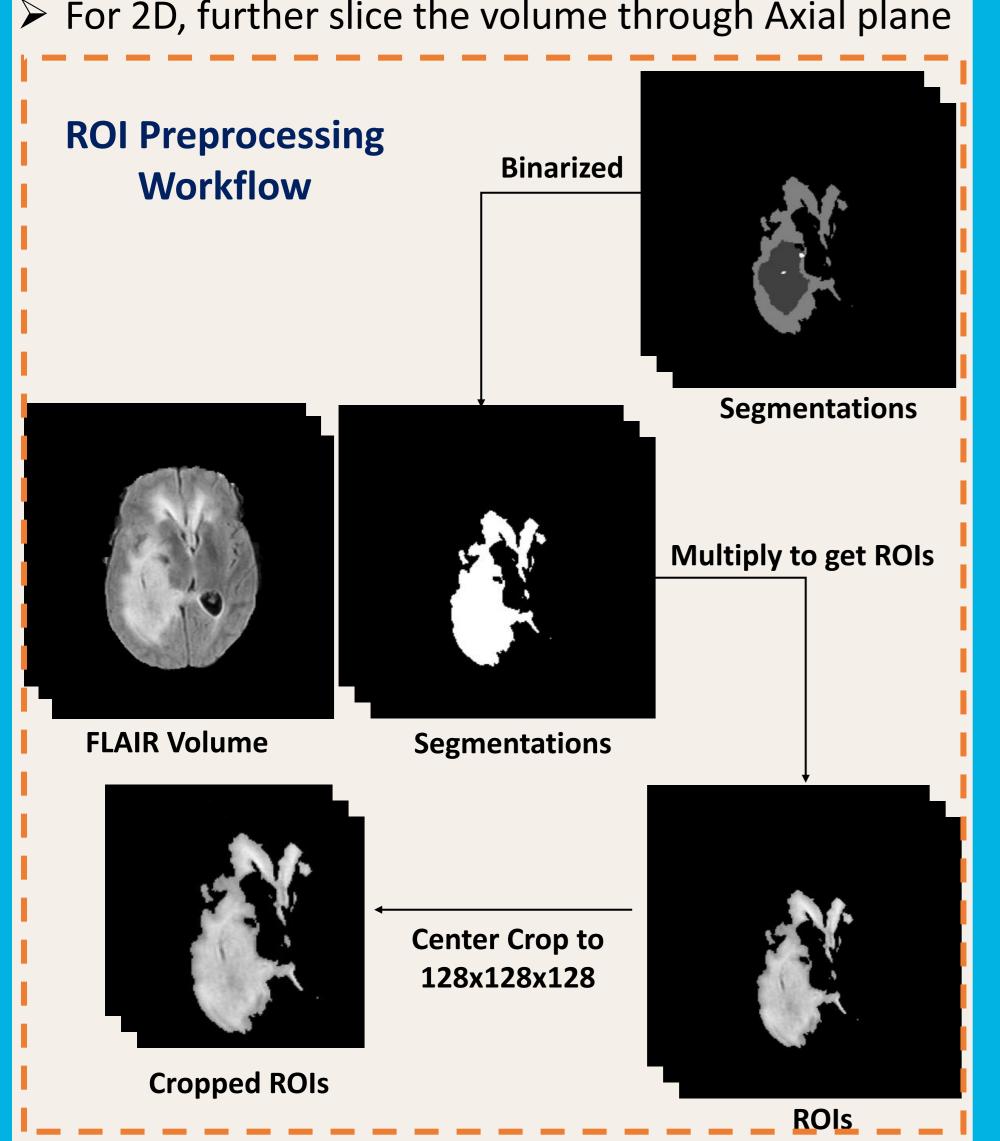
Contributions:

- We present a CNN-Mamba hybrid framework to effectively model local and global features in 2D and 3D medical images.
- We propose dilated gated convnets for multiscale feature learning and cross-modal channel attention for cross-modal information fusion.
- > We are the first to extend the Mamba-based fusion method to 3D medical imaging through a novel tri-plane scanning strategy

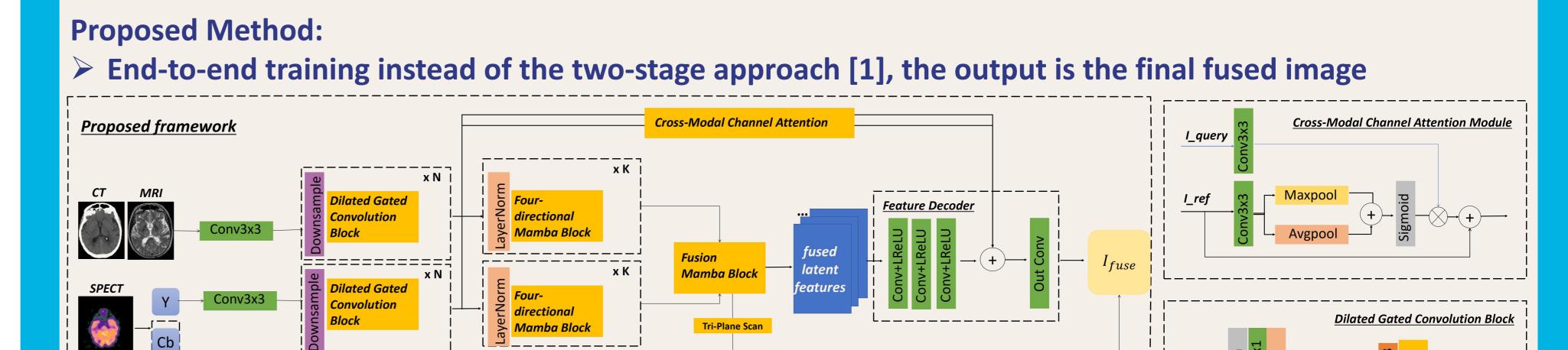
MATERIALS

Dataset and Preprocessing:

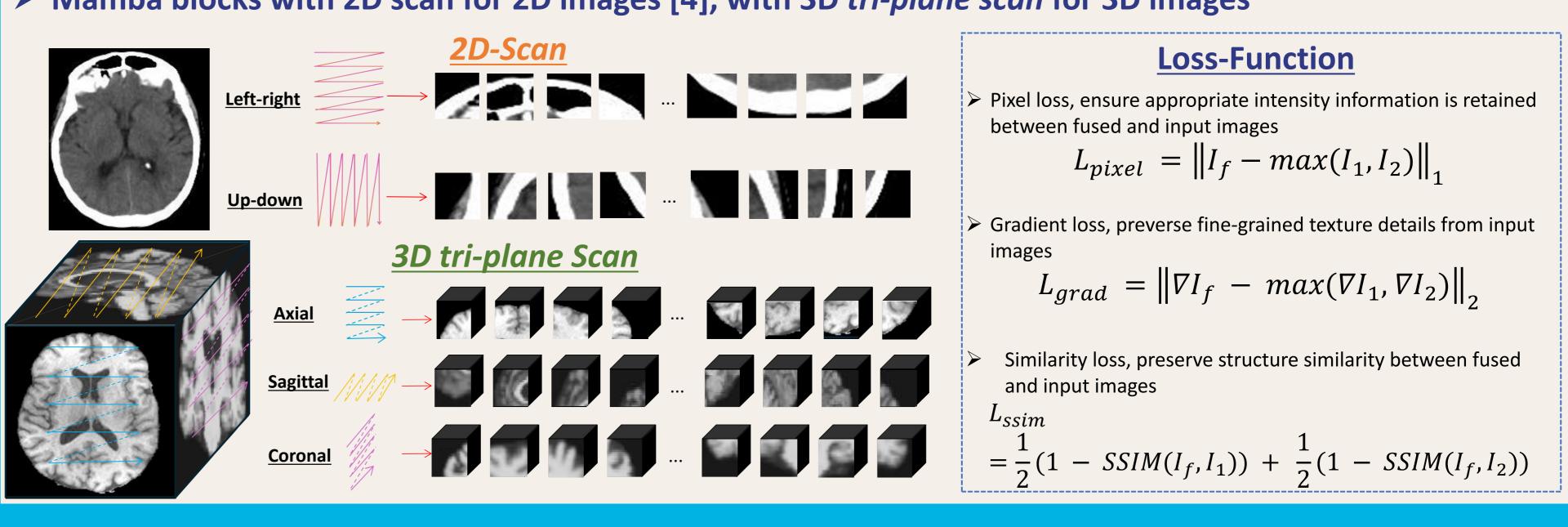
- MRI-CT and MRI-SPECT datasets from The Harvard Whole Brain Atlas dataset, normalized to [0,1]
- BraTS 2019 Dataset with T2 and FLAIR sequences are used for downstream classification task
- Reshape from 240x240x155 to ROI-based volume 128x128x128 [5], normalized to [0,1]
- For 2D, further slice the volume through Axial plane



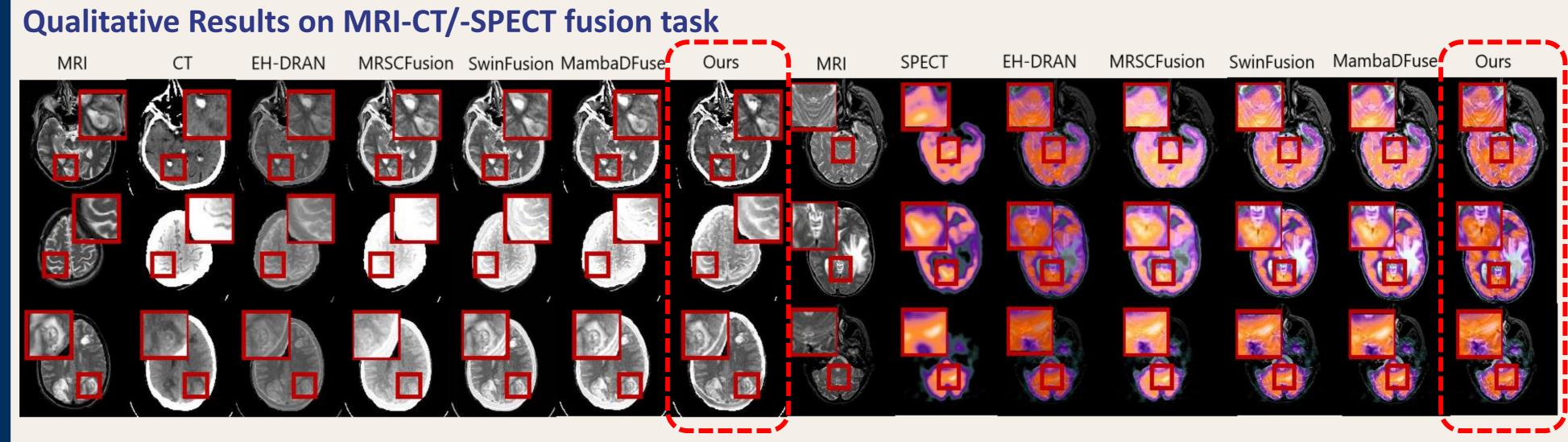
METHODS



- Dilated Gated Conv Block is designed to efficiently learn multi-scale local spatial, descrminative features
- Mamba blocks with 2D scan for 2D images [4], with 3D tri-plane scan for 3D images



RESULTS



Quantitative Results on Downstream Classification task

F1-Score

 0.703 ± 0.018

> 2D and 3D LGG/HGG Pathology type ROI-based classification

Accuracy

 0.604 ± 0.037

FLAIR	0.727 <u>±</u> 0.024	0.701 <u>±</u> 0.008	0.611 <u>+</u> 0.017
T2+FLAIR	0.723±0.028	0.717±0.012	0.640±0.015
EH-DRAN	0.769±0.003	0.723±0.006	0.640 <u>±</u> 0.011
Ours	0.790±0.013	0.778±0.023	0.665±0.004
BraTS-3D	AUC	F1-Score	Accuracy
BraTS-3D T2-3D	AUC 0.647±0.022	F1-Score 0.560±0.029	Accuracy 0.635±0.041
			•
T2-3D	0.647±0.022	0.560 <u>+</u> 0.029	0.635 <u>+</u> 0.041
T2-3D FLAIR-3D	0.647±0.022 0.641±0.110	0.560±0.029 0.529±0.223	0.635±0.041 0.566±0.010

AUC

 0.722 ± 0.021

Ablations

Proposed Cross-modal channel attention module and 3D tri-plane scan strategy

33411341	13614689					
BraTS-2D	PSNR	SSIM	FMI	FSIM	EN	
w/o CMCA	15.967±0.349	0.761±0.007	0.876±0.004	0.813±0.005	14.621±0.075	
with CMCA	16.519±0.352	0.783 <u>+</u> 0.005	0.883±0.003	0.820 <u>+</u> 0.001	15.213±0.069	
	BraTS-3D	PSNR	MS-SSIM	EN		

Bra13-3D	PSINK	IVIS-SSIIVI	EIN
w/o 3D-scan	26.882±0.324	0.833±0.073	19.558±1.374
with 3D-can	33.937±0.361	0.859±0.045	20.468±1.541

Key Takeaways

- Visually, the fused images from our method preserve both modality-specific features and inter-modal contrast.
- Using the fused image, we validate its clinical utility on brain tumor classification task by comparing with other methods such as single-modality, dual-modality and EH-DRAN basline. Our method yields the best performance, demonstrating strong potential for clinical usage.

Conclusion

BraTS-2D

Summary:

- We propose a novel CNN-Mamba hybrid framework for effective multimodal 2D and 3D medical image fusion. Experiments show our method outperforms several baselines.
- The framework is able to generate fused images in real-time, and we valiate the clinical usage of the fused images through brain tumor classification task.

Future Work:

- Extend to other diseases and validation on other tasks such as segmentation.
- Explore pure Mamba-based method to replace CNN.

References

- [1] Zhou M, Zhang Y, Xu X, Wang J, Khalvati F. Edge-Enhanced Dilated Residual Attention Network for Multimodal Medical Image Fusion. In2024 IEEE International Conference on Bioinformatics and Biomedicine (BIBM) 2024 Dec 3 (pp. 4108-4111). IEEE.
- [2] Xie X, Zhang X, Ye S, Xiong D, Ouyang L, Yang B, Zhou H, Wan Y. Mrscfusion: Joint residual swin transformer and multiscale cnn for unsupervised multimodal medical image fusion. IEEE Transactions on Instrumentation and Measurement. 2023 Sep 20;72:1-7.
- [3] Ma J, Tang L, Fan F, Huang J, Mei X, Ma Y. SwinFusion: Cross-domain long-range learning for general
- image fusion via swin transformer. IEEE/CAA Journal of Automatica Sinica. 2022 Jun 30;9(7):1200-17. [4] Peng S, Zhu X, Deng H, Deng LJ, Lei Z. Fusionmamba: Efficient remote sensing image fusion with
- state space model. IEEE Transactions on Geoscience and Remote Sensing. 2024 Nov 11. [5] Zhou M, Khalvati F. Conditional generation of 3d brain tumor regions via vqgan and temporalagnostic masked transformer. In Medical Imaging with Deep Learning 2024 Dec 23.

Acknowledgements

No funding was received for conducting this study. We thank the computational resources provided by Google Colab

