

MedSegMamba: 3D CNN-Mamba Hybrid Architecture for Brain Segmentation

Aaron Cao¹, Zongyu Li², Jordan Jomsky³, Andrew F. Laine², and Jia Guo^{2,4,5*}

¹ University of California, Santa Barbara, Santa Barbara, CA, USA

² Department of Biomedical Engineering, Columbia University, New York, NY, USA

³ Department of Data Science, Columbia University, New York, NY, USA

⁴ Department of Psychiatry, Columbia University, New York, NY, USA

⁵ Mortimer B. Zuckerman Mind Brain Behavior Institute, Columbia University, New York, NY, USA

Abstract. Widely used traditional pipelines for subcortical brain segmentation are often inefficient and slow, particularly when processing large datasets. Furthermore, deep learning models face challenges due to the high resolution of MRI images and the large number of anatomical classes involved. To address these limitations, we developed a 3D patch-based hybrid CNN-Mamba model that leverages Mamba’s selective scan algorithm, thereby enhancing segmentation accuracy and efficiency for 3D inputs. This retrospective study utilized 1784 T1-weighted MRI scans from a diverse, multi-site dataset of healthy individuals. The dataset was divided into training, validation, and testing sets with a 1076/345/363 split. The scans were obtained from 1.5T and 3T MRI machines. Our model’s performance was validated against several benchmarks, including other CNN-Mamba, CNN-Transformer, and pure CNN networks, using FreeSurfer-generated ground truths. We employed the Dice Similarity Coefficient (DSC), Volume Similarity (VS), and Average Symmetric Surface Distance (ASSD) as evaluation metrics. Statistical significance was determined using the Wilcoxon signed-rank test with a threshold of $P < 0.05$. The proposed model achieved the highest overall performance across all metrics (DSC 0.88383; VS 0.97076; ASSD 0.33604), significantly outperforming all non-Mamba-based models ($P < 0.001$). While the model did not show significant improvement in DSC or VS compared to another Mamba-based model (P -values of 0.114 and 0.425), it demonstrated a significant enhancement in ASSD ($P < 0.001$) with approximately 20% fewer parameters. In conclusion, our proposed hybrid CNN-Mamba architecture offers an efficient and accurate approach for 3D subcortical brain segmentation, demonstrating potential advantages over existing methods. Code is available at: <https://github.com/aaroncao06/MedSegMamba>.

Keywords: Biomedical Image Processing · Mamba · Deep Learning · Semantic Segmentation · Neuroimaging

* Correspondence: Jia Guo, jg3400@columbia.edu

1 Introduction

Subcortical brain segmentation is a significant application in medical image processing, as it enables the extraction of quantitative structural information from MRI scans. These localized details can aid in detecting and monitoring morphological deficits in various neuropsychiatric conditions, including Schizophrenia [10], Major Depressive Disorder [15], and Dementia [30]. However, achieving accurate brain segmentation has remained a formidable challenge due to the intricate 3D structures within the brain, the large number of anatomical labels, and the substantial computational resources required to process MRI scans at full resolution.

While manual segmentation stands as the most trusted method, it is a tedious and difficult task, even for experienced clinicians. Automated tools like FreeSurfer [7] have been developed to address these challenges, providing a widely accepted standard for subcortical segmentation. Despite its widespread use, FreeSurfer’s traditional methods can be slow—requiring many hours to process a single scan—and are often sensitive to data quality issues, which limits their reliability when processing large, heterogeneous datasets.

The FastSurfer pipeline [13] [14] is recognized as one of the foremost deep learning-based alternatives to FreeSurfer, capable of executing whole-brain level segmentations. For a 2.5D approach, FastSurfer aggregates three 2D fully convolutional neural networks that utilize the classic encoder-decoder structure originating from the U-Net [26]. Even though FastSurfer has demonstrated superior performance compared to standard alternative models such as 3D U-Net [2], QuickNAT [9], and SDNet [27], its reliance on 2D models inherently limits its ability to capture the full 3D spatial dependencies of the brain’s anatomical structures.

On the other hand, 3D patch-based solutions are better suited to capture such geometries. Although full 3D volume deep learning models for segmenting many classes are currently not feasible due to data and memory constraints, a patch-based approach significantly reduces memory usage and generates more training samples per subject. Additionally, smaller patches allow the model to better capture local 3D information, leading to more accurate segmentation in complex anatomical structures.

Moreover, traditional pure CNN architectures like those used in FastSurfer can suffer from the local receptive fields in each layer, making them susceptible to overlooking the comprehensive global 3D context. The Vision Transformer (ViT) architecture [4] addresses this limitation by employing a self-attention mechanism—originally developed for natural language processing [29]—that achieves state-of-the-art performance in image recognition without relying on convolutions. Each self-attention layer has a global receptive field, enabling the model to extract deeper long-range spatial dependencies. Hybrid CNN-Transformer architectures modeled after the U-Net have become popular for medical image segmentation tasks, showing improved generalization and performance [3] [11]. Specifically for subcortical segmentation, TABSurfer [1], a CNN-Transformer hy-

brid building off of TABS [25] and inspired by TransBTS [32], demonstrated the advantages of combining a hybrid architecture with a 3D patch-based approach.

However, the computational cost of Transformers scales quadratically with sequence length, which leads to substantial hardware memory requirements. This makes their application to dense prediction tasks challenging, particularly for large, high-resolution biomedical images. A shifted window-based self-attention approach [19] can mitigate this issue by reducing the computational burden, enabling the use of pure Transformer models. However, it also restricts the receptive fields in each layer, potentially leading to the loss of global relationships that ViT was designed to capture.

Recent advancements in state space models offer a promising alternative. Mamba [8] introduces the selective scan mechanism (S6) with a hardware-aware algorithm to enhance both training and inference efficiency. This innovation allows the model to scale linearly with respect to sequence length while maintaining state-of-the-art performance across various long-sequence processing tasks.

Since its introduction, Mamba has been adapted to various vision tasks, consistently demonstrating competitive results compared to both Transformer and CNN-based architectures, but with improved memory and parameter efficiency. For 2D image applications, Vision Mamba [36] introduces the Vim block, which incorporates a bidirectional State-Space Model (SSM) combined with positional embeddings. Meanwhile, VMamba [18] introduced the Visual State-Space (VSS) block, centered around the 2D-Selective-Scan (SS2D) module. Because the visual components in images are not sequential, unlike the ordered nature of text and audio, the Vim block unravels the image in two ways (forward and backward) and the SS2D module unravels the image along four traversal paths, to mitigate any unidirectional bias in the final outputs. MambaAD [12] demonstrates that processing additional traversal paths within each Mamba module can enhance performance in 2D anomaly detection tasks. However, in the realm of 3D segmentation, current methods do not fully exploit the richness of the 3D context. For example, U-Mamba [21] unravels the image in only one direction, and Seg-Mamba’s Tri-Oriented Mamba module [34] unfolds the image along just three directions.

In this paper, we introduce MedSegMamba, a hybrid CNN-Mamba architecture that fully leverages Mamba’s selective scan algorithm for 3D inputs, aiming to enhance subcortical brain segmentation.

2 Materials and Methods

This study did not require ethical approval as it utilized publicly available anonymous MRI scans previously acquired for studies approved by local institutional review boards, research ethics committees, or human investigation committees.

2.1 Pipeline

The subcortical segmentation pipeline employs a 3D patch-based approach with a hybrid CNN-Mamba model. Initially, the input scans are centered and con-

formed to LIA orientation, followed by intensity rescaling from 0 to 1, consistent with the steps in the FastSurfer pipeline. These input volumes with dimensions 256 x 256 x 256 are cropped and padded before patch extraction. Patches are then extracted with dimensions 96 x 96 x 96, with a step size of 16 voxels between consecutive patches. Each patch is sequentially fed into the model, and the output class probabilities are reconstructed to match the original input image dimensions. The predicted probabilities for each patch are combined through a voting mechanism to determine the class for each voxel, with the final values mapped to the corresponding FreeSurfer labels. As illustrated in Figure 1a, this pipeline ensures that the model can segment an entire scan into 32 classes in under 90 seconds (all 31 subcortical structures covered by FastSurferVNN excluding cortical white matter).

The hippocampal subfield segmentation process is visualized in Figure 1b. Using the left and right hippocampus classes in the subcortical segmentation, a bounding box with dimensions 96 x 96 x 96 is drawn around the hippocampus in the input scan, and this patch is fed to a separate hybrid CNN-Mamba model, which segments the hippocampal subfields. The output is then padded to the original input image dimensions. This segments the hippocampus into 25 classes according to FreeSurfer’s “FS60” level of hierarchy except with left and right labels separated.

2.2 Model Architecture

3D Selective Scan (SS3D) The selected 3D scanning pattern, depicted in Figure 2d, can unravel the volume along 48 unique traversal paths. Our core SS3D module, shown in Figure 2a, is designed to fully exploit each possible unfolded sequence, with each module processing one of six possible groups of eight sequences.

First, the input volume’s axes are transposed in one of the six possible ways (labeled o0 to o5), determined by the assigned orientation index of the SS3D module. The eight sequences are then extracted from the volume by rotating the volume in three different ways, unfolding it along the scanning pattern, and reversing the sequences. These eight sequences are processed independently and in parallel by the S6 blocks, with their outputs reordered and merged to reconstruct the final output volume. Each of the 48 methods of unraveling the volume is intended to train independent S6 blocks to capture the unique geometries best represented by that specific sequence order.

For the S6 blocks, we increase the SSM state dimension from the standard 16 to 64. Although this adjustment slightly reduces processing speed, it enables the model to extract more detailed information about complex structures with minimal impact on parameters and memory usage. The SSM feature expansion factor is set to 1 to keep the dimensionality constant, as the CNN encoder before the bottleneck has already expanded the feature dimension to 1024.

VSS3D block Our 3D Visual State Space (VSS3D) block resembles a traditional Transformer block, with self-attention replaced by our SS3D module,

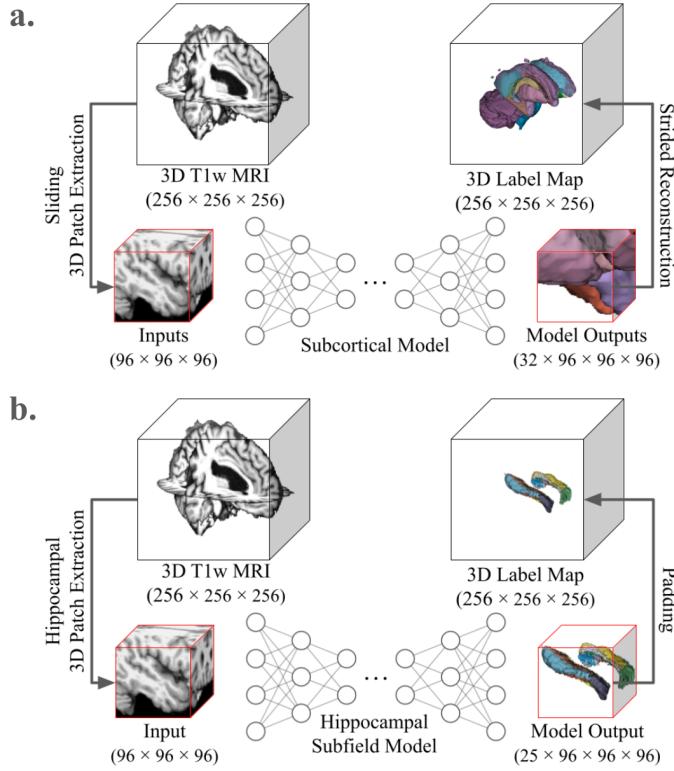


Fig. 1: (a) The subcortical segmentation pipeline extracts 3D patches from the input scan, feeds them into the model, and reconstructs the output predicted label maps to generate the subcortical segmentation of the input scan. (b) The hippocampal subfield segmentation pipeline uses the subcortical segmentation label map to extract one patch centered on the hippocampus region. This patch is fed into the model and the output is padded to the original shape.

inspired by VMamba. It consists of two residual modules. The first module includes a sequence of layer normalization, linear projection, depth-wise convolution, SiLU activation, and SS3D, followed by another layer normalization and linear projection. The second residual module consists of a layer normalization followed by a multilayer perceptron (MLP). This is shown in Figure 2b.

Overall Architecture The overall model architecture features a 3D CNN encoder and decoder with skip connections and a series of VSS3D blocks as the bottleneck. Passing through the encoder, four layers of residual blocks and 3 max pooling operations downsample the input patch for an encoded feature tensor. A subsequent convolution brings the channel size to 1024, before the tensor is fed into the bottleneck. The bottleneck comprises a sequence of nine VSS3D blocks,

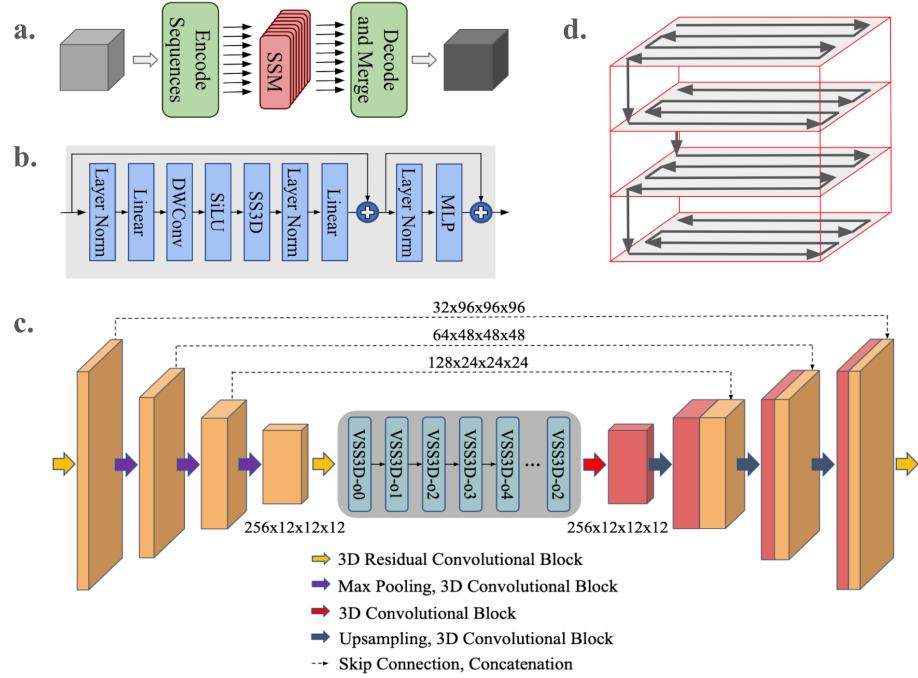


Fig. 2: (a) The SS3D module encodes 8 sequences from the input volume, processes each with an independent S6 block, and then merges the outputs. (b) VSS3D block layout, consisting of SS3D and MLP residual modules. (c) The MedSegMamba model architecture has an encoder-decoder structure with a hybrid of CNN and Mamba-based blocks. Only 6 VSS3D layers are shown here, but the model contains 9 in total. The last Conv1x1 and Softmax layers are also not shown. (d) The SS3D modules unravel the sequences along this continuous 3D scanning pattern.

each assigned an orientation index for its inner SS3D module. The output of the bottleneck is normalized and then passed to the decoder, which reconstructs the image to its original input dimensions. Finally, a convolution operation and a Softmax activation function are applied to generate either a 25 or 32 channel output, where each channel represents the probability of an individual class. The only difference between the subcortical and hippocampal subfield segmentation models is in the channel sizes of these last two layers. Each residual block within the encoder and decoder layers consists of a residual connection and two sequences of 3D Convolution, Group Normalization, and Rectified Linear Unit (ReLU). This is visualized in Figure 2c.

2.3 Data

For the subcortical segmentation dataset, a total of 1784 T1-weighted (T1w) MRI scans were selected from a large-scale heterogeneous dataset representing a uniformly healthy population, compiled from multiple publicly available sources [6]. All scans were acquired at a resolution of 1mm x 1mm x 1mm. Selected subjects came from the Australian Imaging Biomarkers and Lifestyle Study of Ageing (AIBL) [5], Frontotemporal Lobar Degeneration Neuroimaging Initiative (NIFD) database (<http://4rtniftldni.ini.usc.edu>), Information eXtraction from Images (IXI) (<http://brain-development.org/ixi-dataset>), Open Access Series of Imaging Studies-1 (OASIS-1) [23], Open Access Series of Imaging Studies-2 (OASIS-2) [22], Southwest University Adult life-span Dataset (SALD) [33], Southwest University Longitudinal Imaging Multimodal Brain Data Repository (SLIM) [17], Parkinson’s Progression Markers Initiative (PPMI) [24], SchizConnect (SchizConnect) [31], and Consortium for Reliability and Reproducibility (CoRR) [37].

For the hippocampal subfield segmentation dataset, 1221 1mm x 1mm x 1mm T1w scans were selected from the same large-scale multi-site dataset [6]. These subjects came from the Alzheimer’s Disease Neuroimaging Initiative (ADNI) database (adni.loni.usc.edu) and Brain Genomics Superstruct Project (BGSP) [16], in addition to AIBL, IXI, NIFD, and OAS1.

Both datasets were partitioned into training, validation, and test sets with a roughly 3:1:1 ratio. Specifically, the subcortical dataset had a 1076/345/363 split, and the hippocampus dataset had a 719/254/248 split. A roughly balanced age and gender distribution was achieved across the datasets, as shown in Figure 3.

Ground truth segmentations were generated using FreeSurfer, and the input T1w scans were preprocessed with skull-stripping and intensity normalization.

2.4 Model Training

The described model was trained and evaluated on a 24 GB NVIDIA Quadro 6000 GPU. We employed the AdamW optimizer in conjunction with a Cosine Annealing Warm Restarts learning rate scheduler [20]. The loss function used was a combination of Dice Loss and Weighted Cross Entropy during the first epoch, followed by Dice Loss alone in subsequent epochs. Subcortical segmentation training was conducted for 15 epochs on 27 patches per scan (three steps per dimension), utilizing gradient accumulation to simulate each image as a single batch, accounting for the variability in class presence across patches. Hippocampal subfield segmentation training followed the same procedure for 63 epochs except with a batch size of 1. A dropout rate of 0.1 and a drop path rate of 0.3 were applied within the bottleneck.

2.5 Model Evaluation and Statistical Test

To evaluate the efficacy of our novel SS3D module over other 3D mamba-based vision modules for latent space learning, we also trained a model identical to our

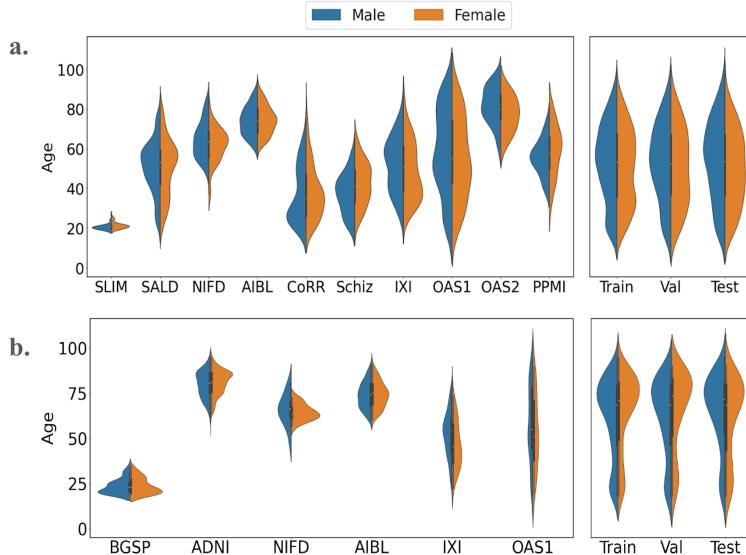


Fig. 3: (a) Age and Gender Distributions for the Subcortical Segmentation Dataset (b) Age and Gender Distributions for the Hippocampal Subfield Segmentation Dataset.

MedSegMamba architecture, except substituting the first residual module in each VSS3D block with SegMamba’s Tri-oriented Mamba module. The resulting block is identical to SegMamba’s Tri-oriented Spatial Mamba block except without the Gated Spatial Convolution module, in order to focus specifically on evaluating the core Mamba modules. This model, referred to as SegMambaBot, was trained following the same procedures as MedSegMamba.

For subcortical segmentation, MedSegMamba was also evaluated against a pretrained FastSurferVINN (obtained from the FastSurfer GitHub repository) and a retrained TABSurfer. Both SegMambaBot and TABSurfer were trained using 3D patch-based approaches with input patch sizes identical to those used in MedSegMamba.

For hippocampal subfield segmentation, MedSegMamba was evaluated against SegMambaBot and TABSurfer, with each model being identical to their subcortical segmentation counterparts except for the final layers generating a 25 instead of 32 channel output. Each model was trained following the same procedure. FastSurfer was not evaluated for this task because only one patch containing the hippocampus needs to be processed, which is small enough for a 3D model to handle in one pass. Aggregating multiple 2D slice-based networks to process this smaller region would be slower and unnecessary.

The evaluation metrics included the Dice Similarity Coefficient (DSC), Volume Similarity (VS) [28], and Average Symmetric Surface Distance (ASSD) [35], assessing both the overall similarity of the segmentations and the contour quality

relative to the ground truth. The Wilcoxon signed-rank test was employed to evaluate the significance of improvements in MedSegMamba’s performance over the other models for each metric. Statistical significance was set to $P < 0.05$.

3 Results

3.1 Subcortical Segmentation

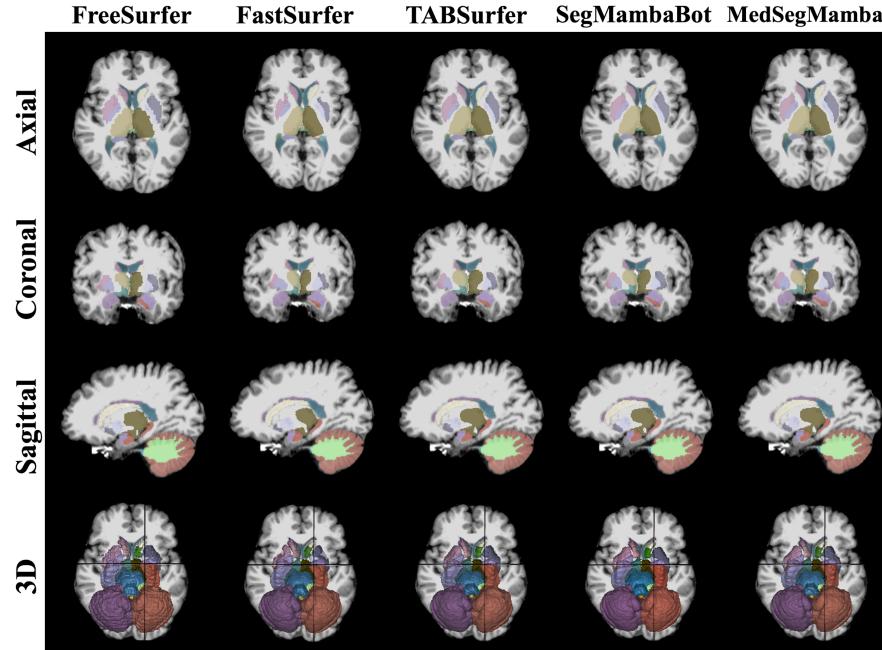


Fig. 4: 2D slices of a sample’s subcortical segmentation by each method.

The average metrics used for evaluating MedSegMamba, SegMambaBot, TAB-Surfer, and FastSurferVNN against the FreeSurfer-generated subcortical segmentation ground truths are displayed in Table 1. MedSegMamba consistently achieved high scores in Dice Similarity Coefficient (DSC), Volume Similarity (VS), and Average Symmetric Surface Distance (ASSD) across all datasets, demonstrating superior overall performance. MedSegMamba’s improvement was statistically significant across all metrics, except for the DSC and VS metrics when compared to SegMamba (P -values of 0.114 and 0.425, respectively).

All three 3D patch-based methods also significantly outperformed the pre-trained 2.5D FastSurfer benchmark. Furthermore, both Mamba-based models significantly outperformed TABSurfer, their transformer-based counterpart.

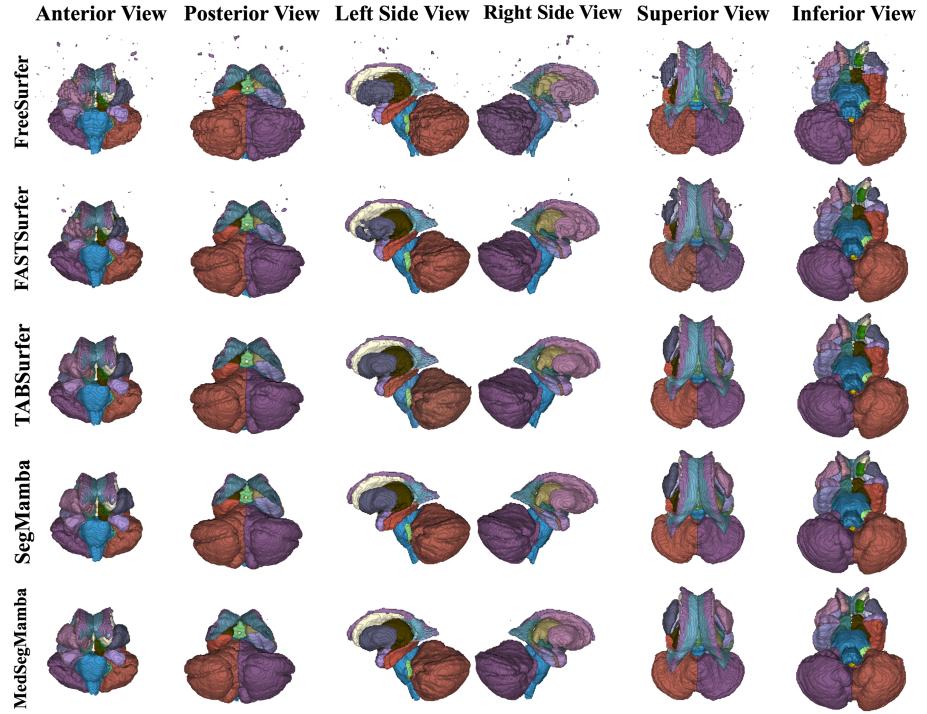


Fig. 5: 3D renderings of a sample's subcortical segmentation by each method.

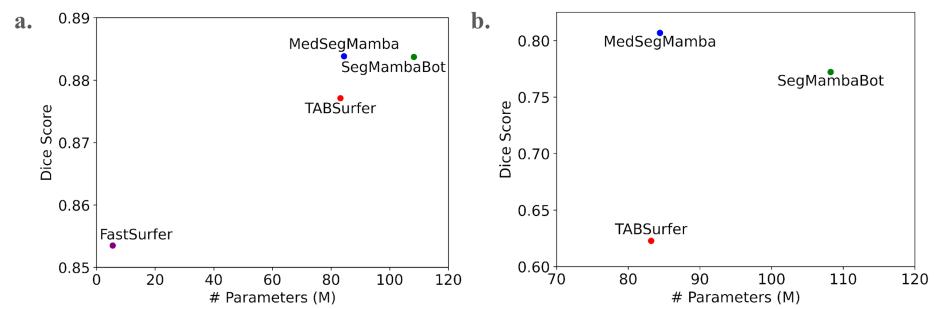


Fig. 6: (a) Plot of millions of parameters versus subcortical segmentation dice score for each model. (b) Plot of millions of parameters versus hippocampal subfield segmentation dice score for each model.

Dataset	Model	DSC \uparrow	VS \uparrow	ASSD \downarrow
AIBL	MedSegMamba	0.89593\pm0.009	0.97406\pm0.006	0.29069 \pm 0.042
	SegMambaBot	0.89580 \pm 0.009	0.97396 \pm 0.006	0.29038\pm0.040
	TABSurfer	0.89094 \pm 0.010	0.97082 \pm 0.007	0.30257 \pm 0.044
	FastSurfer	0.87856 \pm 0.015	0.96485 \pm 0.009	0.33453 \pm 0.059
CoRR	MedSegMamba	0.88635\pm0.020	0.97123\pm0.008	0.31922\pm0.078
	SegMambaBot	0.88627 \pm 0.020	0.97040 \pm 0.008	0.32062 \pm 0.077
	TABSurfer	0.87949 \pm 0.020	0.96765 \pm 0.009	0.33669 \pm 0.074
	FastSurfer	0.86607 \pm 0.027	0.95989 \pm 0.016	0.38024 \pm 0.104
IXI	MedSegMamba	0.86773\pm0.023	0.96485\pm0.010	0.41766\pm0.102
	SegMambaBot	0.86729 \pm 0.023	0.96460 \pm 0.010	0.42425 \pm 0.106
	TABSurfer	0.85877 \pm 0.027	0.96023 \pm 0.011	0.43849 \pm 0.108
	FastSurfer	0.81217 \pm 0.034	0.93340 \pm 0.018	0.61388 \pm 0.140
NIFD	MedSegMamba	0.90060\pm0.007	0.97734 \pm 0.006	0.26797\pm0.030
	SegMambaBot	0.90026 \pm 0.008	0.97745\pm0.005	0.27077 \pm 0.032
	TABSurfer	0.89260 \pm 0.008	0.97361 \pm 0.005	0.29188 \pm 0.035
	FastSurfer	0.88791 \pm 0.008	0.97240 \pm 0.005	0.30494 \pm 0.030
OAS1	MedSegMamba	0.88954 \pm 0.012	0.97351 \pm 0.006	0.30856 \pm 0.050
	SegMambaBot	0.88977\pm0.012	0.97387\pm0.007	0.30839\pm0.050
	TABSurfer	0.88305 \pm 0.011	0.96808 \pm 0.005	0.32129 \pm 0.041
	FastSurfer	0.87504 \pm 0.010	0.96240 \pm 0.005	0.34141 \pm 0.047
OAS2	MedSegMamba	0.88680\pm0.013	0.97017 \pm 0.009	0.31079\pm0.047
	SegMambaBot	0.88670 \pm 0.012	0.97048\pm0.008	0.31382 \pm 0.047
	TABSurfer	0.88247 \pm 0.013	0.96611 \pm 0.009	0.32258 \pm 0.047
	FastSurfer	0.87984 \pm 0.013	0.96454 \pm 0.008	0.32390 \pm 0.049
PPMI	MedSegMamba	0.89373\pm0.007	0.97345 \pm 0.005	0.29935\pm0.037
	SegMambaBot	0.89333 \pm 0.007	0.97422\pm0.005	0.30096 \pm 0.036
	TABSurfer	0.88941 \pm 0.008	0.97158 \pm 0.006	0.30483 \pm 0.032
	FastSurfer	0.87892 \pm 0.008	0.96673 \pm 0.006	0.32847 \pm 0.033
SALD	MedSegMamba	0.87785 \pm 0.028	0.96998 \pm 0.009	0.35457 \pm 0.094
	SegMambaBot	0.87795\pm0.028	0.97015\pm0.009	0.35236\pm0.093
	TABSurfer	0.87046 \pm 0.027	0.96560 \pm 0.009	0.37473 \pm 0.092
	FastSurfer	0.84189 \pm 0.021	0.94426 \pm 0.014	0.48194 \pm 0.089
Schiz	MedSegMamba	0.88204 \pm 0.011	0.97118\pm0.007	0.33188\pm0.044
	SegMambaBot	0.88205\pm0.011	0.97082 \pm 0.008	0.33257 \pm 0.042
	TABSurfer	0.87613 \pm 0.011	0.96747 \pm 0.007	0.34891 \pm 0.053
	FastSurfer	0.83832 \pm 0.022	0.94248 \pm 0.015	0.48469 \pm 0.094
SLIM	MedSegMamba	0.88692 \pm 0.006	0.97101 \pm 0.006	0.30611\pm0.024
	SegMambaBot	0.88740\pm0.007	0.97169\pm0.006	0.30679 \pm 0.025
	TABSurfer	0.88254 \pm 0.006	0.96972 \pm 0.005	0.31617 \pm 0.022
	FastSurfer	0.85483 \pm 0.012	0.94886 \pm 0.008	0.42471 \pm 0.051
Overall	MedSegMamba	0.88383\pm0.021	0.97076\pm0.009	0.33604\pm0.086
	SegMambaBot	0.88372 \pm 0.021	0.97068 \pm 0.009	0.33761 \pm 0.087
	TABSurfer	0.87711 \pm 0.022	0.96684 \pm 0.009	0.35264 \pm 0.088
	FastSurfer	0.85350 \pm 0.035	0.95229 \pm 0.018	0.43554 \pm 0.143

Table 1: Comparing MedSegMamba, SegMambaBot, TABSurfer, and FastSurfer-VINN metrics across datasets. Bold text indicates superior performance. \uparrow indicates that higher numbers correspond to better performance and \downarrow indicates that lower numbers correspond to better performance.

Views of 3D renderings and 2D slices of segmentations produced by each method on a sample subject are shown in Figure 5 and Figure 4. Although FreeSurfer’s atlas-based method yielded the noisiest segmentation, all deep learning methods produced smoother contours. FastSurfer’s 2.5D method retained the most noise from the FreeSurfer output, whereas the 3D patch-based methods achieved similarly smooth segmentations, differing only in minor areas for this sample. However, SegMamba exhibited slight under-segmentation in certain regions, and TABSurfer demonstrated less consistency across different scans.

While both Mamba-based models showed competitive performance, MedSegMamba demonstrated slight improvements over SegMambaBot and has approximately 20% fewer parameters while requiring similar GPU memory during inference (1.988 GiB for MedSegMamba vs. 1.932 GiB for SegMambaBot).

As shown in Figure 6a, MedSegMamba demonstrates superior performance while remaining parameter efficient compared to the other 3D patch-based methods. While FastSurfer’s fully convolutional architecture requires very few parameters, each of its three sub-models requires more than double the GPU memory during inference (just under 5 GiB) compared to the 3D models, all of which operate with under 2 GiB of GPU memory.

3.2 Hippocampal Subfield Segmentation

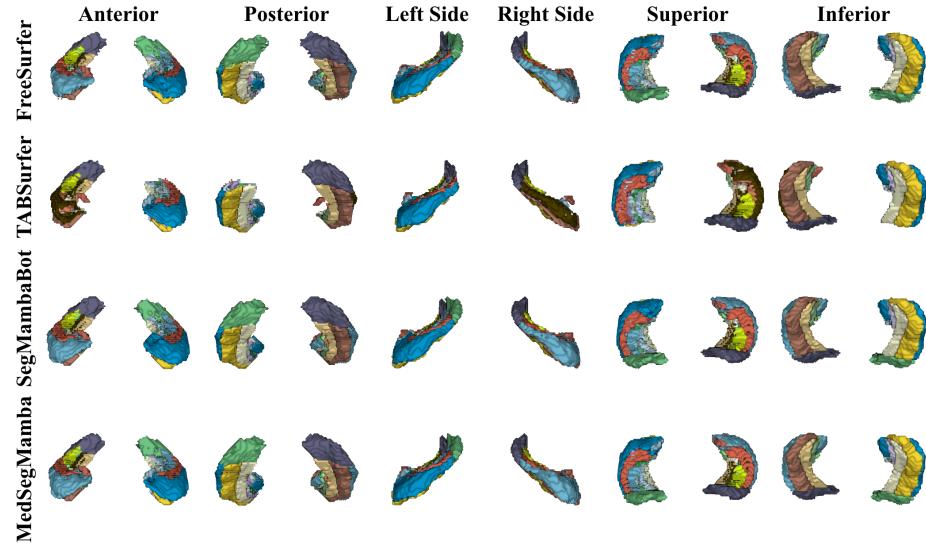


Fig. 7: 3D renderings of a sample’s hippocampal subfield segmentation by each method.

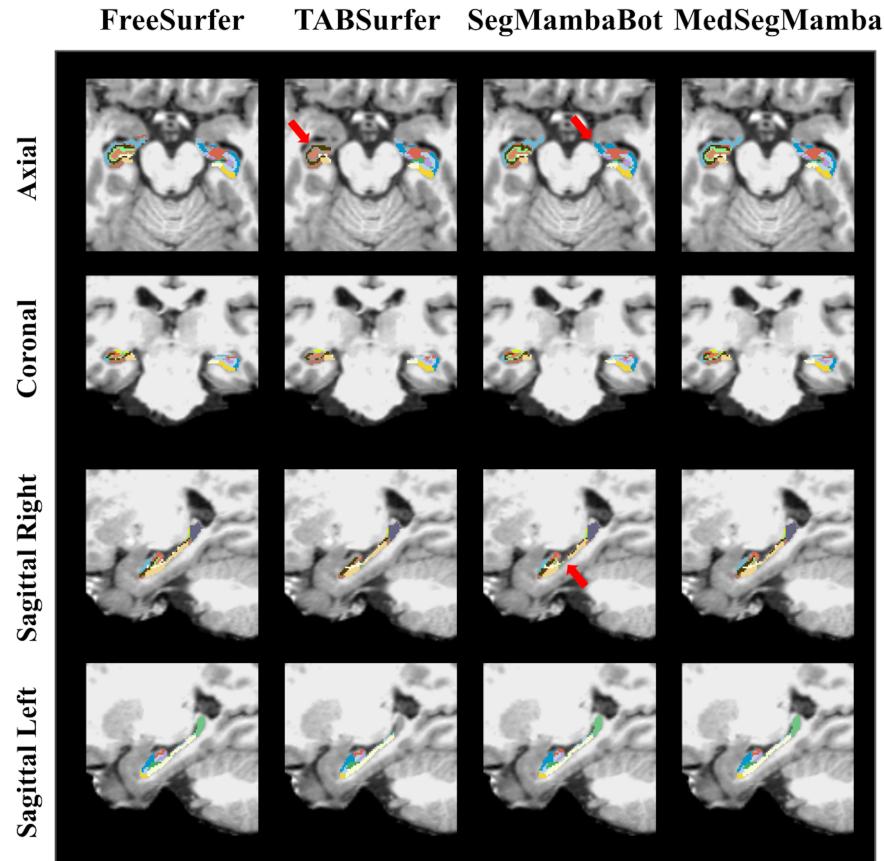


Fig. 8: 2D slices of a sample’s hippocampal subfield segmentation by each method are shown. Large regions in the right hippocampus are missing from TABSurfer’s segmentation, as indicated by the red arrow in the axial view. In SegMambaBot’s segmentation, the left hippocampal amygdala transition area is missing, as shown in the axial view, and the right presubiculum is slightly undersegmented, as shown in the sagittal right view.

Dataset	Model	DSC ↑	VS ↑	ASSD* ↓
ADNI	MedSegMamba	0.80415±0.047	0.96207±0.014	0.22819±0.068
	SegMambaBot	0.76976±0.045	0.92309±0.014	—
	TABSurfer	0.61932±0.037	0.75894±0.011	—
AIBL	MedSegMamba	0.80935±0.037	0.96511±0.012	0.21981±0.055
	SegMambaBot	0.77561±0.032	0.92525±0.010	—
	TABSurfer	0.62564±0.031	0.76095±0.009	—
BGSP	MedSegMamba	0.81948±0.023	0.96375±0.010	0.21061±0.037
	SegMambaBot	0.78562±0.018	0.92611±0.008	—
	TABSurfer	0.63435±0.016	0.76336±0.009	—
IXI	MedSegMamba	0.80858±0.035	0.96461±0.010	0.22449±0.052
	SegMambaBot	0.77445±0.036	0.92566±0.010	—
	TABSurfer	0.62394±0.033	0.75991±0.009	—
NIFD	MedSegMamba	0.81438±0.045	0.96871±0.008	0.22214±0.062
	SegMambaBot	0.77286±0.051	0.92698±0.007	—
	TABSurfer	0.61505±0.043	0.76174±0.006	—
OAS1	MedSegMamba	0.79368±0.077	0.95663±0.027	0.27493±0.288
	SegMambaBot	0.75796±0.072	0.91743±0.027	—
	TABSurfer	0.61299±0.061	0.75517±0.020	—
Overall	MedSegMamba	0.80683±0.046	0.96299±0.015	0.22960±0.115
	SegMambaBot	0.77219±0.043	0.92379±0.014	—
	TABSurfer	0.62276±0.037	0.75979±0.012	—

Table 2: Comparing MedSegMamba, SegMambaBot, and TABSurfer metrics across the hippocampal subfield segmentation datasets. Bold text indicates superior performance. ↑ indicates that higher numbers correspond to better performance and ↓ indicates that lower numbers correspond to better performance. *ASSD is not shown for SegMambaBot and TABSurfer because it cannot be computed for the classes which SegMambaBot and TABSurfer failed to learn.

Metrics evaluating MedSegMamba, SegMambaBot, and TABSurfer’s hippocampal segmentations are shown in Table 2. MedSegMamba’s improvements over both SegMambaBot and TABSurfer are more notable for this task than for subcortical segmentation, as seen in Figure 6b. Likely due to the smaller structures present in this task, MedSegMamba was the only network that learned to predict all 25 classes, while SegMambabot only learned 24 classes, and TABSurfer converged on predicting just 20 classes. This resulted in significant improvements over both models in DSC and VS. ASSD cannot be computed for missing classes, so statistical tests are not applicable between the models’ overall ASSD. However, even when ignoring the missing classes, SegMambaBot and TABSurfer get an average ASSD of 0.23380 and 0.25512, respectively, which are both still inferior to MedSegMamba’s average of 0.22960. The 3D and 2D visualizations of a sample’s hippocampus segmentation by each method are visualized in Figure 7 and Figure 8. TABSurfer failed to segment numerous classes, leaving large regions missing. While SegMambaBot and MedSegMamba appear similar,

SegMambaBot is missing the left hippocampal amygdala transition area and slightly under-segments some regions, leaving MedSegMamba as the only model to succeed in segmenting each subregion.

4 Discussion

This study introduces the SS3D module and VSS3D block, an extension of Mamba’s selective scan operation tailored for complex 3D image processing tasks. Integrated within the MedSegMamba hybrid CNN-Mamba architecture, the VSS3D block demonstrates robust latent-space learning capabilities, outperforming existing traditional and deep learning approaches across multiple datasets.

The transition from a Transformer bottleneck to a Mamba-based bottleneck significantly enhances memory efficiency, thereby enabling the use of more convolutional layers for the encoder and decoder on identical hardware. This enhancement boosts local feature extraction capabilities while preserving global context extraction, resulting in both SegMambaBot and MedSegMamba outperforming TABSurfer. Comparing SegMambaBot and MedSegMamba in subcortical segmentation, the SS3D module allows MedSegMamba to attain comparable overall performance in terms of DSC and VS metrics but offers superior region boundary delineation as reflected by the ASSD, all while utilizing significantly fewer parameters than the Tri-oriented Mamba module. Increasing the number of selective scan operations to process variously unraveled sequences can achieve performance parity with larger modules that process fewer sequences.

Compared to FastSurfer’s standard 2.5D approach, our utilization of 3D inputs preserves more intricate spatial relationships within the anatomy’s continuity than 2D slices. The 3D strided reconstruction process, which aggregates model outputs from shifted overlapping patches, further mitigates noise artifacts present in the ground truth. This improvement is observed qualitatively when comparing the noisier segmentations of FastSurfer to the smoother outputs of the 3D patch-based models.

The hippocampal subfield segmentation experiment also provides evidence of the SS3D module’s ability to process finer details compared to the Tri-oriented Mamba module. While SegMambaBot and TABSurfer failed this task, MedSegMamba was able to reliably segment all the small and intricate subregions.

4.1 Limitations

The primary limitation of this approach is that MedSegMamba exhibits slower performance than SegMambaBot despite its lower parameter count, most notably during the subcortical segmentation pipeline. However, this issue can be mitigated by increasing the step size during the strided reconstruction process from 16 to 32, resulting in a roughly 4x speedup, reducing the processing time to around 22 seconds per scan with minimal impact on performance. When both use a step size of 32, MedSegMamba’s subcortical segmentation is only a few

seconds slower than SegMambaBot per scan, making them roughly comparable in terms of overall efficiency while retaining the superior performance afforded by the SS3D module.

4.2 Conclusion

These results showcase the advantages of this novel hybrid architecture combined with 3D patch-based processing. Future research should explore the application of our SS3D module across different architectures, including its integration at various latent space levels throughout the encoder, as seen in SegMamba’s original design. Ablation studies focused on determining the optimal number of scanning directions or unraveling methods for 3D segmentation may further enhance our model’s performance.

5 Acknowledgments

No funding was received for conducting this study and there are no relevant financial or non-financial interests to disclose.

References

1. Cao, A., Rao, V.M., Liu, K., Liu, X., Laine, A.F., Guo, J.: Tabsurfer: a hybrid deep learning architecture for subcortical segmentation. arXiv preprint arXiv:2312.08267 (2023)
2. Çiçek, Özgün and Abdulkadir, Ahmed and Lienkamp, Soeren S and Brox, Thomas and Ronneberger, Olaf: 3d u-net: Learning dense volumetric segmentation from sparse annotation. In: Medical Image Computing and Computer-Assisted Intervention – MICCAI 2016. pp. 424–432. Springer (2016)
3. 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 (2021)
4. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. ICLR (2021)
5. Ellis, K.A., Bush, A.I., Darby, D., De Fazio, D., Foster, J., Hudson, P., Lautenschlager, N.T., Lenzo, N., Martins, R.N., Maruff, P., et al.: The australian imaging, biomarkers and lifestyle (aibl) study of aging: methodology and baseline characteristics of 1112 individuals recruited for a longitudinal study of alzheimer’s disease. International psychogeriatrics **21**(4), 672–687 (2009)
6. Feng, X., Lipton, Z.C., Yang, J., Small, S.A., Provenzano, F.A.: Estimating brain age based on a uniform healthy population with deep learning and structural magnetic resonance imaging. Neurobiology of Aging **91**, 15–25 (2020)
7. Fischl, B., Salat, D.H., Busa, E., Albert, M., Dieterich, M., Haselgrove, C., van der Kouwe, A., Killiany, R., Kennedy, D., Klaveness, S., et al.: Whole brain segmentation: Automated labeling of neuroanatomical structures in the human brain. Neuron **33**(3), 341–355 (2002)

8. Gu, A., Dao, T.: Mamba: Linear-time sequence modeling with selective state spaces. arXiv preprint arXiv:2312.00752 (2023)
9. Guha Roy, A., Conjeti, S., Navab, N., Wachinger, C.: Quicknat: A fully convolutional network for quick and accurate segmentation of neuroanatomy. NeuroImage **186**, 713–727 (2019)
10. Gutman, B.A., Van Erp, T.G., Alpert, K., Ching, C.R.K., Isaev, D., Ragothaman, A., Jahanshad, N., Saremi, A., Zavaliangos-Petropulu, A., Glahn, D.C., et al.: A meta-analysis of deep brain structural shape and asymmetry abnormalities in 2,833 individuals with schizophrenia compared with 3,929 healthy volunteers via the enigma consortium. Human Brain Mapping **43**(1), 352–372 (2022)
11. Hatamizadeh, A., Tang, Y., Nath, V., Yang, D., Myronenko, A., Landman, B., Roth, H., Xu, D.: Unetr: Transformers for 3d medical image segmentation. In: Proceedings of the IEEE/CVF winter conference on applications of computer vision. pp. 574–584 (2022)
12. He, H., Bai, Y., Zhang, J., He, Q., Chen, H., Gan, Z., Wang, C., Li, X., Tian, G., Xie, L.: Mambaad: Exploring state space models for multi-class unsupervised anomaly detection. arXiv preprint arXiv:2404.06564 (2024)
13. Henschel, L., Conjeti, S., Estrada, S., Diers, K., Fischl, B., Reuter, M.: Fastsurfer - a fast and accurate deep learning based neuroimaging pipeline. NeuroImage **219**, 117012 (Oct 2020)
14. Henschel, L., Kügler, D., Reuter, M.: Fastsurfervinn: Building resolution-independence into deep learning segmentation methods—a solution for highres brain mri. NeuroImage **251**, 118933 (2022)
15. Ho, T.C., Gutman, B., Pozzi, E., Grabe, H.J., Hosten, N., Wittfeld, K., Völzke, H., Baune, B., Dannlowski, U., Förster, K., et al.: Subcortical shape alterations in major depressive disorder: Findings from the enigma major depressive disorder working group. Human Brain Mapping **43**(1), 341–351 (2022)
16. Holmes, A.J., Hollinshead, M.O., O'keefe, T.M., Petrov, V.I., Fariello, G.R., Wald, L.L., Fischl, B., Rosen, B.R., Mair, R.W., Roffman, J.L., et al.: Brain genomics superstruct project initial data release with structural, functional, and behavioral measures. Scientific data **2**(1), 1–16 (2015)
17. Liu, W., Wei, D., Chen, Q., Yang, W., Meng, J., Wu, G., Bi, T., Zhang, Q., Zuo, X.N., Qiu, J.: Longitudinal test-retest neuroimaging data from healthy young adults in southwest china. Scientific data **4**(1), 1–9 (2017)
18. Liu, Y., Tian, Y., Zhao, Y., Yu, H., Xie, L., Wang, Y., Ye, Q., Liu, Y.: Vmamba: Visual state space model. arXiv preprint arXiv:2401.10166 (2024)
19. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 10012–10022 (2021)
20. Loshchilov, I., Hutter, F.: Sgdr: Stochastic gradient descent with warm restarts. arXiv preprint arXiv:1608.03983 (2016)
21. Ma, J., Li, F., Wang, B.: U-mamba: Enhancing long-range dependency for biomedical image segmentation. arXiv preprint arXiv:2401.04722 (2024)
22. Marcus, D.S., Fotenos, A.F., Csernansky, J.G., Morris, J.C., Buckner, R.L.: Open access series of imaging studies: longitudinal mri data in nondemented and demented older adults. Journal of cognitive neuroscience **22**(12), 2677–2684 (2010)
23. Marcus, D.S., Wang, T.H., Parker, J., Csernansky, J.G., Morris, J.C., Buckner, R.L.: Open access series of imaging studies (oasis): cross-sectional mri data in young, middle aged, nondemented, and demented older adults. Journal of cognitive neuroscience **19**(9), 1498–1507 (2007)

24. Marek, K., Jennings, D., Lasch, S., Siderowf, A., Tanner, C., Simuni, T., Coffey, C., Kieburtz, K., Flagg, E., Chowdhury, S., et al.: The parkinson progression marker initiative (ppmi). *Progress in neurobiology* **95**(4), 629–635 (2011)
25. Rao, V.M., Wan, Z., Arabshahi, S., Ma, D.J., Lee, P.Y., Tian, Y., Zhang, X., Laine, A.F., Guo, J.: Improving across-dataset brain tissue segmentation for mri imaging using transformer. *Frontiers in Neuroimaging* **1** (2022)
26. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*. pp. 234–241. Springer (2015)
27. Roy, A.G., Conjeti, S., Sheet, D., Katouzian, A., Navab, N., Wachinger, C.: Error corrective boosting for learning fully convolutional networks with limited data. In: Descoteaux, M., Maier-Hein, L., Franz, A., Jannin, P., Collins, D.L., Duchesne, S. (eds.) *Medical Image Computing and Computer Assisted Intervention - MICCAI 2017*. pp. 231–239. Springer International Publishing, Cham (2017)
28. Taha, A.A., Hanbury, A.: Metrics for evaluating 3d medical image segmentation: analysis, selection, and tool. *BMC Medical Imaging* **15**(1), 29 (Aug 2015)
29. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. *Advances in Neural Information Processing Systems* **30** (2017)
30. van der Velpen, I.F., Vlasov, V., Evans, T.E., Ikram, M.K., Gutman, B.A., Roshchupkin, G.V., Adams, H.H., Vernooij, M.W., Ikram, M.A.: Subcortical brain structures and the risk of dementia in the rotterdam study. *Alzheimer's & Dementia* **19**(2), 646–657 (2023)
31. Wang, L., Alpert, K.I., Calhoun, V.D., Cobia, D.J., Keator, D.B., King, M.D., Kogan, A., Landis, D., Tallis, M., Turner, M.D., et al.: Schizconnect: Mediating neuroimaging databases on schizophrenia and related disorders for large-scale integration. *Neuroimage* **124**, 1155–1167 (2016)
32. Wang, W., Chen, C., Ding, M., Yu, H., Zha, S., Li, J.: Transbts: Multimodal brain tumor segmentation using transformer. In: *Medical Image Computing and Computer Assisted Intervention – MICCAI 2021*. pp. 109–119. Springer (2021)
33. Wei, D., Zhuang, K., Ai, L., Chen, Q., Yang, W., Liu, W., Wang, K., Sun, J., Qiu, J.: Structural and functional brain scans from the cross-sectional southwest university adult lifespan dataset. *Scientific data* **5**(1), 1–10 (2018)
34. Xing, Z., Ye, T., Yang, Y., Liu, G., Zhu, L.: Segmamba: Long-range sequential modeling mamba for 3d medical image segmentation. arXiv preprint arXiv:2401.13560 (2024)
35. Yeghiazaryan, V., Voiculescu, I.: Family of boundary overlap metrics for the evaluation of medical image segmentation. *Journal of Medical Imaging* **5**, 1 (02 2018)
36. Zhu, L., Liao, B., Zhang, Q., Wang, X., Liu, W., Wang, X.: Vision mamba: Efficient visual representation learning with bidirectional state space model (2024)
37. Zuo, X.N., Anderson, J.S., Bellec, P., Birn, R.M., Biswal, B.B., Blautzik, J., Breitner, J., Buckner, R.L., Calhoun, V.D., Castellanos, F.X., et al.: An open science resource for establishing reliability and reproducibility in functional connectomics. *Scientific data* **1**(1), 1–13 (2014)