

BrainSegFounder: Towards Foundation Models for Neuroimage Segmentation

Joseph Cox¹, Peng Liu¹, Skylar E. Stolte¹, Yunchao Yang¹, Kang Liu¹, Kyle B. See¹, Huiwen Ju², and Ruogu Fang^{1,*}

¹ University of Florida, Gainesville, FL 32611, USA

{cox.j,plui1,skylastolte444,yunchaoyang,
kang.lui,kylebsee,ruogu.fang}@ufl.edu

² NVIDIA, Santa Clara, CA 95051, USA

hju@nvidia.com

Abstract. The burgeoning field of brain health research increasingly leverages artificial intelligence (AI) to analyze and interpret neuroimaging data. Medical foundation models have shown promise of superior performance with better sample efficiency. This work introduces a novel approach towards creating 3-dimensional (3D) medical foundation models for multimodal neuroimage segmentation through self-supervised training. Our approach involves a novel two-stage pretraining approach using vision transformers. The first stage encodes anatomical structures in generally healthy brains from the large-scale unlabeled neuroimage dataset of multimodal brain magnetic resonance imaging (MRI) images from 41,400 participants. This stage of pretraining focuses on identifying key features such as shapes and sizes of different brain structures. The second pretraining stage identifies disease-specific attributes, such as geometric shapes of tumors and lesions and spatial placements within the brain. This dual-phase methodology significantly reduces the extensive data requirements usually necessary for AI model training in neuroimage segmentation with the flexibility to adapt to various imaging modalities. We rigorously evaluate our model, BrainSegFounder, using the Brain Tumor Segmentation (BraTS) challenge and Anatomical Tracings of Lesions After Stroke v2.0 (ATLAS v2.0) datasets. BrainSegFounder demonstrates a significant performance gain, surpassing the achievements of the previous winning solutions using fully supervised learning. Our findings underscore the impact of scaling up both the model complexity and the volume of unlabeled training data derived from generally healthy brains. Both of these factors enhance the accuracy and predictive capabilities of the model in neuroimage segmentation tasks. Our pretrained models and code are at <https://github.com/lab-smile/BrainSegFounder>.

Keywords: Neuroimaging · 3D Foundation Model · Self-Supervised Learning · Brain Tumor Segmentation · Multi-modal MRI

1 Introduction

The fusion of artificial intelligence (AI) with neuroimaging analysis, particularly multimodal MRI, is forging a pivotal role in advancing brain health ([1], [2], [3], [4], [5], [6] [7]). The complexity of the human brain, with its elaborate anatomy and intricate functions, poses significant challenges in neuroimaging analysis ([8], [9],[10], [2], [6]). AI’s capability to interpret complex neurological data has the potential to enhance diagnostic precision and deepen our understanding of brain pathology. Numerous studies have aimed to develop AI models for specific brain health analyses, each contributing to the growing body of neuroimaging research.

Traditionally, neuroimaging AI models require extensive fine-tuning through supervised learning to address a specific downstream task. Modifications of the nnU-Net ([11]), DeepScan ([12]), and DeepMedic ([13]) architectures have performed well on a host of medical computer vision challenges such as the Brain Tumor Segmentation (BraTS) challenge ([14]), Medical Segmentation Decathlon (MSD) ([15]), and A tumor and liver automatic segmentation challenge (ATLAS) ([16]). Many of these advances stem from utilizing self-supervised pretraining methods on large, unlabeled datasets to transfer weights for model encoders and decoders to the smaller datasets present in the challenge ([17], [18]). Complementary to these pretraining modifications, there has been a recent push towards developing massive medical datasets ([19], [20], [21]) to

* Corresponding Author

aid in the creation of these models. However, medical image analysis has yet to benefit from the recent advances in natural image analysis and language processing through models like the Segment Anything Model (SAM) ([22]) and LLaMA ([23]).

In medical language processing, models like MI-Zero ([24]) and BioViL-T ([25]) utilize contrastive learning to make significant advancements in representational analysis and zero-shot transfer learning in medical image recognition. By leveraging different learning objectives, similar image-text pairs are pulled closer in the latent space while dissimilar pairs are pushed further apart. Such models have pushed the boundary of histopathology research and combined text-based analysis with computer vision. Yet, they rely on having text-based prompts accompanying their training images ([26]).

With SAM’s demonstrated success on few-shot segmentation tasks of natural images, recent research into medical image segmentation models has primarily modified the SAM architecture. Models like MedSAM ([27]), MedLSAM ([28]), and SAM-Med2D ([29]) focus on bridging the gap between SAM’s generalizability on real-world images and its performance on medical tasks. They accomplish this by adapting the SAM architecture to these medical tasks. [27] crafted MedSAM for image segmentation by constructing a massive dataset of image-mask pairs derived from sizeable medical image databases. MedLSAM further refined upon MedSAM by including landmark localization. SAM-Med2D further improved segmentation results by increasing the dataset to multiple modalities and increasing prompt density. However, these models function in 2-dimensional space, requiring 3-dimensional (3D) modalities to be sub-sampled or solved in slices ([9]). Not only is this computationally inefficient, but the most dense and often most valuable information is found in 3D modalities like CT or MRI. [30] aimed to address this discrepancy by adapting the SAM models to 3D space using a visual sampler and a mask decoder to aggregate layers. Their model, dubbed 3DSAM-adapter, outperformed leading segmentation models in various tasks while still utilizing an algorithm that functions in 2D space. These results indicated that these models would benefit from the critical anatomical and spatial information found from being fully capable of functioning in 3D space.

Despite the progress in medical imaging, holistically analyzing the vast amount of data generated by brain MRIs remains a formidable challenge ([9]). The intricate structure and function of the brain necessitate advancements in MRI analysis due to their critical impact on patient outcomes, especially in the early detection and treatment of brain disorders ([10]). Existing AI models in neuroimaging are hampered by their need for extensive supervised learning and their limited ability to generalize across different tasks without substantial retraining, revealing a gap for a robust, adaptable model that functions in 3D space ([9], [10]).

This study presents BrainSegFounder, a 3D foundational framework for multimodal neuroimage segmentation. BrainSegFounder is designed to pave the way towards setting new standards for the accuracy and efficiency of medical AI models. We focus our study on two essential tasks - brain tumor segmentation and brain lesion segmentation. A primary obstacle in creating AI models for brain tumor and brain lesion analysis is the scarcity of brain tumors within the general population. This scarcity significantly hampers the compilation of large-scale diseased patient datasets, which are essential for the supervised training of AI models. In response, the development process of BrainSegFounder incorporates a multi-stage approach to feature learning, specifically engineered to mitigate the challenges posed by data scarcity.

In its initial phase, BrainSegFounder leverages an extensive dataset from brain scans of 41,400 participants from the United Kingdom. This foundational step enables the framework to effectively encode generally healthy brain tissue structures, creating a detailed baseline of anatomical features from a predominantly healthy population. Subsequently, the framework’s training shifts focus towards identifying disease-specific attributes, such as geometric shapes of tumors and lesions and spatial placements within the brain. This dual-phase methodology significantly diminishes the extensive data requirements usually necessary for AI model training in tumor detection. Moreover, it naturally expands the dataset available for the AI to learn from efficiently and straightforwardly, sidestepping the need for generating synthetic images. This approach mirrors the analytical techniques used by radiologists and has undergone thorough validation against the BraTS challenge and ATLAS 2.0 datasets, showcasing significant improvements over current models.

BrainSegFounder, our novel framework, represents a pivotal advance in neuroimaging analysis by laying the groundwork for a future of comprehensive foundation models in this field. BrainSegFounder is designed to be adaptable for various neurological tasks, including brain tumor segmentation, stroke localization, brain region segmentation, and the diagnosis of Alzheimer’s disease. By utilizing a large dataset of brain imaging from a generally healthy population, BrainSegFounder sets the stage for transforming clinical workflows, aiming to enhance the speed and accuracy of diagnoses across a spectrum of neurological conditions. The

	UK Biobank	BraTS	ATLAS
Number of Subjects	41,400	1,251	655
Modalities	T1w, T2-FLAIR	T1w, T1-ce, T2w, T2-FLAIR	T1-ce
Number of Images	82,800	5,004	655
Diseases	Generally Healthy	Malignant Brain Neoplasms	Stroke

Table 1. A summary of the data used in this study.

The BrainSegFounder framework introduces a deep learning training scheme tailored for diverse applications by showcasing a distinct approach to self-supervised pretraining followed by precise fine-tuning. This section offers a detailed examination of the framework’s architecture and its procedural pipeline. It highlights the multi-stage self-supervised pretraining, termed Stage 1 and Stage 2, before proceeding to fine-tuning for downstream tasks. Figure 1 illustrates BrainSegFounder’s architecture. Central to BrainSegFounder is a vision transformer-based encoder that employs a series of self-attention mechanisms. This encoder is linked with an up-sampling decoder tailored for segmentation tasks. The architecture is adapted from the SwinUNETR architecture [31] with modified input channels and input hyperparameters. BrainSegFounder pioneers a novel dual-phase self-supervised pretraining method, integrating self-supervised learning components within its structure. Stage 1 pretraining exposes the framework to a wide-ranging dataset of brain MRIs from the UK Biobank dataset, predominantly consisting of healthy individuals. This initial stage equips the model with a thorough comprehension of standard brain anatomy, utilizing self-supervised learning to enhance prediction capabilities. Stage 2 of pretraining advances the model’s proficiency by introducing it to a specialized MRI dataset geared toward the downstream task. This phase leverages the architecture’s refined anomaly detection skills, focusing on distinguishing deviations in brain structure.

Following pretraining, BrainSegFounder undergoes fine-tuning on the final dataset, where transfer learning enhances the model’s encoder. As depicted in Figure 1, the fine-tuning process leverages the pretrained Swin Transformer encoder from the earlier two stages. The first pretraining stage on the UKB dataset develops a foundational understanding of normal brain anatomy. The second stage of pretraining with diseased datasets builds upon this foundation by introducing pathology, thus allowing the model to learn the distinction between healthy and pathological tissues. Transfer learning is applied after each pretraining stage to retain and refine the knowledge acquired, ensuring that the model can effectively adapt to the new dataset while preserving previously learned patterns.

The culmination of this process is the integration of the U-NET decoder, which works in concert with the pretrained encoders to generate segmentation scores that delineate tumor boundaries with precision. This hybrid approach combines the strengths of the Swin Transformer and UNETR architectures, optimizing the model for the critical task of tumor segmentation and providing an authoritative score that reflects the model’s accuracy in identifying and delineating tumor regions.

In summary, the BrainSegFounder model’s architecture and pretraining paradigm represent a comprehensive approach to understanding and segmenting brain images, with a training pipeline that methodically builds the model’s capacity to differentiate and characterize complex patterns in 3D MRI data without external annotation. Together, the self-supervised learning stages and the fine-tuning process prepare BrainSegFounder to tackle downstream tasks with high efficiency and accuracy.

1.1 Data Acquisition and Preprocessing

Throughout our pretraining and fine-tuning, we make use of the UK Biobank (UKB), Brain Tumor Segmentation (BraTS) Challenge, and Anatomical Tracings of Lesions After Stroke v2.0 (ATLAS v2.0) datasets. The following section summarizes the dataset information that is pertinent to our study. Table 1 provides an overview of this information.

UK Biobank dataset In our first stage, we utilize T1-weighted (T1w) and T2-weighted Fluid Attenuation Inversion Recovery (T2-FLAIR) from the UK Biobank (UKB) dataset ([32]). These data points were collected starting from 2014 and preprocessed by the UKB. Utilizing a comprehensive 35-minute protocol, the UKB

obtained many brain imaging modalities, including T1w and T2-FLAIR structural brain MRI images ([33]). We obtained all T1w and T2-FLAIR images available between 2014 and 2022 from 44,172 participants with neuroimaging data. Raw T1w structural volumes were processed using a processing pipeline developed by UK Biobank researchers that consisted primarily of tools from FSL and Freesurfer. The pipeline generated additional images like segmentations between different types of matter and effectively reduced the non-brain tissue interference. Volumetric measures of gray matter and internal structures were generated alongside the processed images, providing valuable insights into the characteristics of gray matter and internal structures. Each T1w structural image underwent further processing with FreeSurfer ([34]), followed by a quality control check for inclusion into the data made available by UK Biobank researchers. Additionally, T2-FLAIR images were aligned to the corresponding T1 image, resulting in two additional product images. The UK Biobank saves volumetric measures of white matter lesions as additional data with both T1w and T2-FLAIR volumes. Images were reconstructed as DICOM images and converted to the NIfTI format using dcm2niix ([35]). All imaging volumes were defaced to preserve participant anonymity and broadcast to the MNI152 template space using FNIRT ([34]). Among these 44,172 participants, 43,369 participants have both T1w and T2-FLAIR images. To build 3D foundation models of neuroimages, we selected participants who had at least 100 slices in both T1w and T2-FLAIR volumes. This criteria resulted in 41,400 participants and 82,800 imaging volumes. A CONSORT diagram depicting the data used in this study can be found in Figure 3, and demographic data for participants is summarized in Figure 2. For detailed information, see the Appendix.

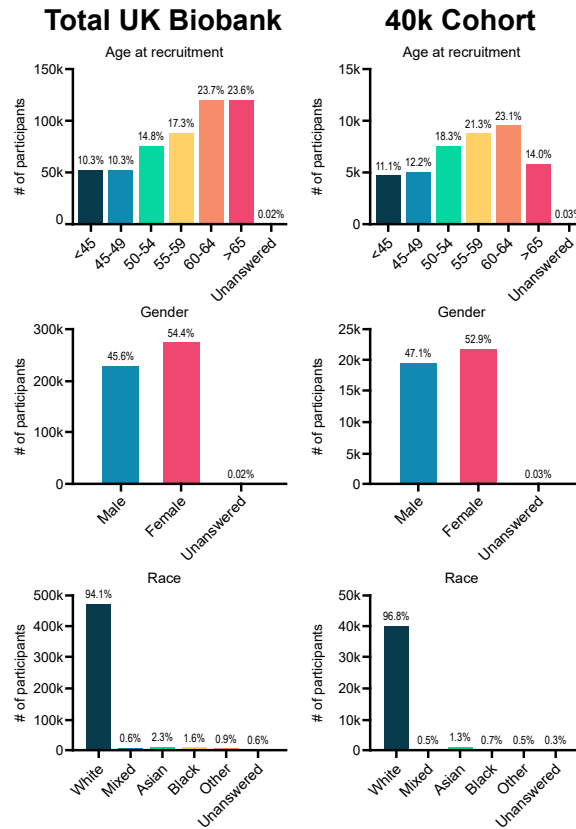


Fig. 2. Visual representation of demographic data from subjects in the UK Biobank in the study.

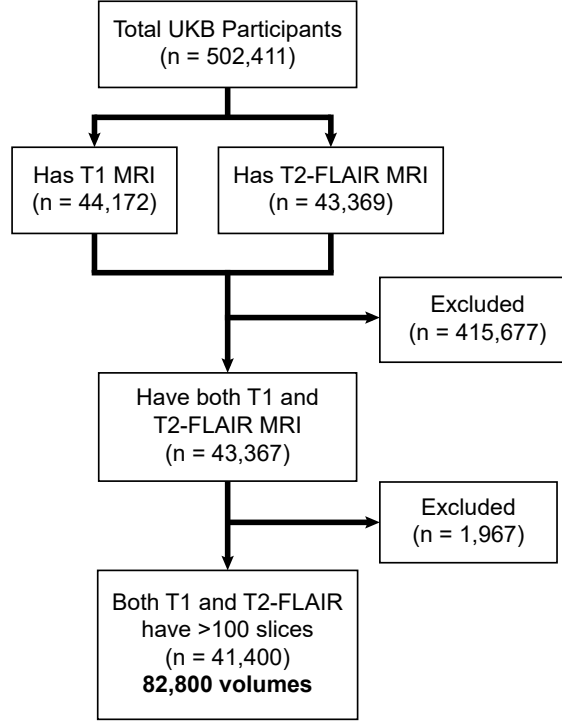


Fig. 3. CONSORT diagram of UKB data used in Stage 1 pretraining.

BRaTS dataset In Stage 2 and Stage 3, one of the datasets we used to perform self-supervised pretraining and finetuning on MRI images is from the training set of the BraTS 2021 Task 1 (Tumor Segmentation) challenge. This dataset consists of 1,251 subjects each with T1w, T1-contrast enhanced (T1-ce), T2-weighted (T2w), and T2-FLAIR images. We obtained all publicly available imaging volumes as part of the challenge. The BraTS challenge utilizes a standard preprocessing pipeline similar to the UKB dataset. First, images are converted from DICOM images to NIfTI using `dcm2niix` ([35]) and tools available from the Cancer Imaging Phenomics Toolkit (CaPTk) ([36]). Images are then co-registered to the SRI24 template and resampled to a uniform resolution of 1 mm^3 . Finally, each modality is skull stripped, defaced, and converted to NIfTI ([14]). Imaging volumes were then segmented into three tumor classes using the STAPLE algorithm ([37]) across previous BraTS winners and refined manually. These manual annotations were further verified by multiple board-certified neuro-radiologists, resulting in quality-controlled tumor segmentation labels across all four modalities in 3 classes: Gd-enhancing tumor (referred to as the whole tumor (WT), edematous tissue (ED), and necrotic tumor core (TC).

ATLAS v2.0 Dataset Additionally, we perform self-supervised Stage 2 pretraining and fine-tuning on MRI images from the training set of the Anatomical Tracings of Lesions After Stroke (ATLAS) v2.0 Dataset [38]. This dataset consists of 655 T1-ce MRIs aggregated from 44 research cohorts. Each MRI is from one subject, and time points range from <24 hours to >180 days after stroke onset. The standard labeling pipeline for the ATLAS dataset consists of (1) manual quality control to exclude significant motion artifacts, (2) manual lesion tracing in ITK-SNAP [39] [40], and (3) lesion mask review by two independent raters. The data then went through a similar preprocessing pipeline to BraTS. The MR images were intensity-normalized and registered to the MNI-152 template using the MINC toolkit (<https://github.com/BIC-MNI/minc-toolkit>). Finally, FreeSurfer’s MRI deface functionality was used to deface the scans. Images were reviewed again in a final quality check at the end of the pipeline before being included in the dataset. Segmentations are evaluated on four metrics - Dice coefficient for the final segmentation (Dice), the difference between true total lesion volume and predicted total lesion volume (Volume Difference), the difference in the number of lesions

		BrainSegFounder-Tiny (62M)		BrainSegFounder-Small (64M)		BrainSegFounder-Big (69M)	
# of encoder parameters		19,097,191		20,982,103		26,636,839	
Encoder layer level	Output size	# of SSL Heads	# Swin Blocks	# of SSL Heads	# Swin Blocks	# of SSL Heads	# Swin Blocks
Level 1	48x(48x48x48)	3	2	3	2	3	2
Level 2	96x(24x24x24)	6	2	6	2	6	2
Level 3	192x(12x12x12)	12	2	12	6	12	18
Level 4	384x(6x6x6)	24	2	24	2	24	2
Feature size		48		48		48	
Bottleneck dimension		768		768		768	

Table 2. Pretraining encoder settings.

between ground truth and prediction (Lesion Count), and Lesion-wise F1 Score. The Lesion-wise F1 Score is calculated by performing a 3D connected-component analysis to determine true positives, false positives, and false negatives. A true positive is any 3D connected component in the ground-truth image that overlaps with at least one voxel in the prediction image. Conversely, a false positive is any 3D connected component in the prediction image that does not overlap with the ground-truth image. A false negative is a connected component in the ground truth that lacks overlapping voxels in the prediction image [38].

1.2 Stage 1: Pretraining on the UKB

The initial pretraining stage involves the self-supervised learning of a transformer-based neural network model using a substantial unlabeled image dataset. For this purpose, the UKB dataset ([32]) is utilized. From our 82,800 3D volumetric images used for pretraining, the input MRI modalities are randomly cropped into $96 \times 96 \times 96$ sub-volumes and augmented with random inner cutout and rotation. These augmented images are then fed into the SwinUNETR encoder for processing.

The SwinUNETR architecture incorporates a Swin Transformer encoder that handles 3D input patches. This encoder operates with a patch size of $2 \times 2 \times 2$, a feature dimension of 8, and an embedding space of 48 dimensions. It consists of four stages, with a patch merging layer introduced between stages to reduce the feature size by half.

Adopting the methodology from [18], the SwinUNETR encoder is pretrained through three distinct proxy tasks that serve as self-supervised fine-tuning mechanisms: masked volume inpainting, 3D image rotation, and contrastive coding. The primary objective of pretraining is to minimize the total loss function. This work has developed three models to address varying complexities. These models include the foundational BrainSegFounder-Tiny with 62 million parameters, the intermediate BrainSegFounder-Small with 64 million parameters, and the advanced BrainSegFounder-Big with 69 million parameters. The primary differentiation among these models is the variation in the number of sliding window blocks within their third stage. Table 2 shows BrainSegFounder’s sliding-window encoder backbone’s parameters, number of SSL heads, and number of sliding window blocks.

For the pretraining process, 64 NVIDIA DGX A100 GPUs, distributed across 8 DGX-2 nodes, are deployed at the University of Florida’s HiPerGator-AI supercomputer. Data parallelism is implemented to optimize the efficiency of model training. Both training and validation losses are monitored to track progress. The AdamW optimizer is employed using a warm-up cosine scheduler set for 500 iterations. The training employs a batch size of 2 per GPU, using $96 \times 96 \times 96$ patches. The initial learning rate is established at 6×10^{-6} , coupled with a momentum of 0.9 and a decay of 0.1 over 15,000 iterations. These parameters are summarized in Table 3.

1.3 Training on BraTS

Stage 2: Pretraining on BraTS The pretrained models based on the UK Biobank (UKB) dataset underwent further pretraining through transfer learning on the Brain Tumor Segmentation (BraTS) dataset. 1,251 subjects were employed for a 5-fold cross-validation process. To ensure consistent performance evaluation, the data splits for these 5 folds were kept identical to those used in the baseline SwinUNETR model. During training, four of the folds were utilized for training purposes, and the remaining folds served for validation.

Given that the BraTS dataset comprises four modalities, but only two (specifically, T1w and T2-FLAIR) were available for pretraining in the initial stage, the first layer of the pretrained network on UKB was

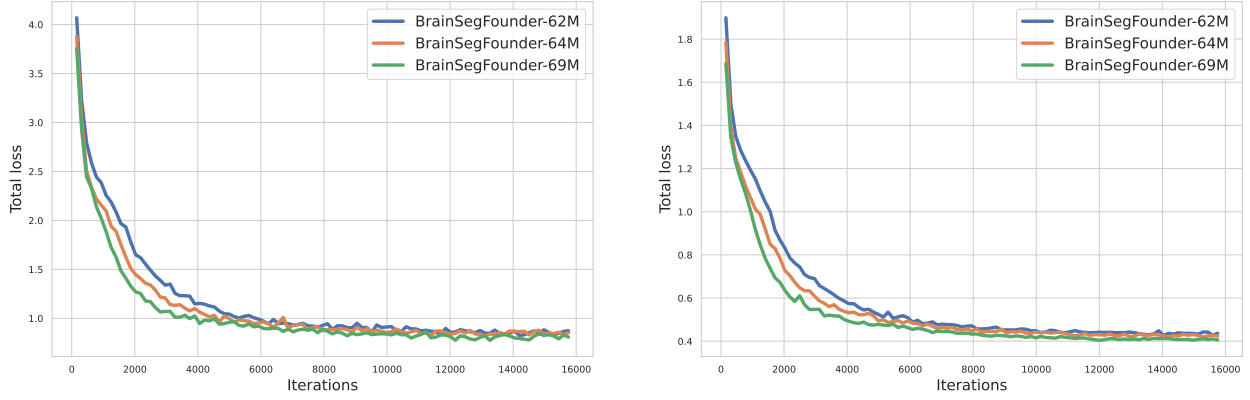


Fig. 4. Training (left) and validation (right) loss of Stage 1-pretraining three different scale of BrainSegFounder models on UKB.

Table 3. Hardware and training parameters.

	Stage	Data	No. Subjects	GPU	Batch size	Learning rate	No. steps	No. input channel
Encoder only	Stage 1 Pretraining	UKB	43369	64 × A100	128	1×10^{-6}	2×10^5	2
Encoder + Decoder	Stage 2 Pretraining & Stage 3 Fine-tuning	BraTS	1251	2 × A100	2	1×10^{-4}	50000	4
Encoder + Decoder	Stage 2 Pretraining & Stage 3 Fine-tuning	ATLAS v2.0	655	4 × A100	4	3×10^{-3}	600	1

modified. This modification involved expanding the number of input channels by adding two new channels, whose weights were randomly initialized using the Kaiming initialization method ([41]).

Hyperparameter settings for Stage-2 Pretraining can be found in Table 3. For pretraining on BraTS, two NVIDIA A100 GPUs, each with 32 GB of memory, were utilized. Depending on the model size, the BrainSegFounder models require between 48 to 72 hours for training. The batch size and learning rate were uniformly set at 2 and 1×10^{-4} for all models during this pretraining phase.

Stage 3: Fine-tuning on BraTS In the final fine-tuning stage we attach the pretrained encoder from the previous stage to a UNet decoder. This model is then finetuned directly on the BraTS dataset. We used the same hyper-parameter settings as those used in the Stage 2 pretraining phase on BraTS (in Table 3): The batch size remained at 2, mirroring the encoder-only stage, and the learning rate remained at 1×10^{-4} . The number of steps for this phase was set to 50,000, with the input data having 4 channels, which indicates the typical inclusion of multi-modal MRI scans in the BraTS dataset.

Few-shot Learning on BraTS To investigate our model’s performance using limited training data, we conducted a systematic comparison between BrainSegFounder and the baseline model, SwinUNETR, utilizing a descending percentage training approach in the context of the BraTS challenge. Using both our BrainSegFounder pretrained model and SwinUNETR, we finetuned on 40% of the BraTS training dataset, with subsequent incremental reductions in data availability, decreasing to a final 5% of the original dataset. Due to the potential high-variability on the input dataset at such small percentages of input data, we trained on 5 different randomly sampled subsets of the input training data and calculated the average performance across these subsets. This method aimed to explore the impact of training data scarcity on model performance and adaptability. Performance evaluations were carried out on the BraTS test set after each training step and evaluated with the Dice coefficient to assess segmentation accuracy.

Modality Restriction and Flexibility in Training and Inference The proposed method is designed to be adaptable across various data modalities in downstream tasks. In the case that fewer modalities are available in the downstream task, Stage 1 pretrained model using both T1- and T2-weighted MRI can be adapted and fine-tuned on fewer modalities (e.g., T1- or T2- weighted MRI alone) during Stage 2 pretraining, Stage 3 supervised training, and the inference without requiring any modifications to the network structure.

This is achieved by simply configuring the two input channels to process the same type of data (either T1- or T2-weighted).

In the case that more modalities are available in the downstream task, our model can also accommodate this by increasing the number of input channels. The pre-trained weights are then loaded into the corresponding layers of the network.

To investigate the efficacy of this method, we performed ablation testing on BraTS by restricting the modalities available to the model. The Stage 1 model was given only either T1w or T2w images rather than utilizing both modalities available from the UKBiobank. These models, pretrained on only one modality, were then trained with our Stage 2 pretraining pipeline on the BraTS data with only the modality trained on in Stage 1 instead of the 4 modalities available from BraTS. Finally, the model was finetuned on BraTS with the same single modality. Hyperparameters and number of GPUs were kept the same as in our earlier pretraining steps on BraTS as summarized in Table 3.

1.4 Training on ATLAS v2.0

Stage 2: Pretraining on ATLAS v2.0 Similarly, the pretrained models based on the UK Biobank (UKB) underwent further self-supervised pretraining through transfer learning on the ATLAS v2.0 dataset. In this stage, a total of 665 MR images were included in the training set, intentionally avoiding cross-validation to align our methodology with that employed by the submissions in the challenge leaderboard. This approach allows for direct performance comparisons under similar training conditions, providing a robust test of our models against established benchmarks.

Since the ATLAS dataset has only one modality, T1-ce, the first layers of the Stage-1 pretrained model were modified by dropping the channel corresponding to the T2w modality present in the UKB. For pre-training, four NVIDIA A100 GPUs, each with 32 GB of memory, were utilized. The stage 2 model took 35 hours to train, with a batch size and learning rate set to 4 and 5×10^{-3} , respectively.

Stage 3: Fine-tuning on ATLAS v2.0 Upon completion of pretraining, the model was fine-tuned further on the ATLAS v2.0 dataset to adapt to the specific challenges of lesion detection in stroke patients. The fine-tuning employed a cosine-annealing learning rate scheduler, starting with an initial learning rate of 1×10^{-4} . Batch size was set to 4, and the model was trained for a total of 600 epochs.

Training was conducted using 2 NVIDIA A100 GPUs with 32GB of RAM accessible to each GPU. We applied data augmentation techniques of random cropping, rotation, and to improve model robustness against variations in real-world data. The loss function used was Dice-loss, suited for addressing imbalance between classes. Dropout at a rate of 10% was applied to prevent overfitting.

Model performance was periodically assessed using a held-out set of 100 images from the training dataset, ensuring that the model’s improvements were generalizable and aiding with tuning the hyperparameters effectively and tracking training progress without the use of a separate validation dataset.

2 Results

2.1 Pretraining

The pretraining of our BrainSegFounder models, which varied in size based on the number of parameters, took between 3 to 6 days. This process utilized a computational setup ranging from 8 to 64 NVIDIA A100 GPUs, each with 80 GB capacity. Figure 4 illustrates the validation loss during the pretraining phase across different BrainSegFounder model sizes.

2.2 Evaluation on BraTS Challenge Dataset

Comparison to State-of-the-Art Methods Table 4 summarizes our BrainSegFounder’s best performing model against published results on our validation splits from other state-of-the-art models on the BraTS challenge: the baseline SwinUNETR model [31], nnU-Net [11], TransBTS [42], SegResNet [43], and MONAI’s Model-Zoo [44]. SegResNet and nnU-Net are both winning methodologies in previous BraTS challenges.

Table 4. A comparison of BrainSegFounder (BSF) models’ performance in terms of average Dice coefficient on the BraTS challenge. BSF-S indicates our best performing BrainSegFounder model (small, 64M parameters). BrainSegFounder models were pretrained with SSL on T1- and T2-weighted MRI 3D volumes and finetuned with supervised learning using all four modalities present in BraTS. BSF-1-S indicates this model with only the Stage 1 (SSL) pertaining on UKB and without the Stage 2 pretraining step. SwinU models are models using the SwinUNETR architecture trained on BraTS via supervised learning. SwinU-MRI is the model trained directly using supervised learning on BraTS published on GitHub (<https://github.com/Project-MONAI/research-contributions/tree/main/SwinUNETR/BRAATS21>), SwinU-Res is pretrained with SSL on only T1w and T2w and finetuned on BraTS, and SwinU-CT pretrained using CT data and finetuned with supervised learning on BraTS. nnU-Net and SegResNet are former BraTS challenge winners trained using supervised learning on our folds. TransBTS is a vision-transformer based segmentation algorithm optimized for brain segmentation. Model-Zoo is a bundle of models published by MONAI that can perform BraTS segmentation out of the box using their "Brats mri segmentation" [sic] model found at <https://monai.io/model-zoo.html>

	BSF-S (64M)	BSF-1-S (64M)	SwinU-MRI	SwinU-Res	SwinU-CT	nnU-Net	SegResNet	TransBTS	Model-Zoo
Fold 1	0.9032	0.8994	0.8854	0.895	0.894	0.896	0.899	0.883	0.857
Fold 2	0.9182	0.9055	0.9059	0.899	0.902	0.917	0.916	0.902	0.879
Fold 3	0.9121	0.9125	0.8981	0.894	0.898	0.910	0.909	0.889	0.820
Fold 4	0.9100	0.9133	0.8924	0.890	0.893	0.909	0.908	0.893	0.889
Fold 5	0.9139	0.9114	0.9035	0.903	0.902	0.909	0.906	0.892	0.893
Average	0.9115	0.9110	0.8971	0.896	0.898	0.908	0.907	0.891	0.868

TransBTS is a vision-transformer based approach tailored for brain tumor segmentation. MONAI’s Model-Zoo is a bundle of medical imaging models capable of performing a wide variety of tasks, including BraTS segmentation. In addition, we include comparison to the corresponding single-stage pretrained model that was pretrained only on the UKB. Table 8 in the appendix present a comparative analysis of all trained BrainSegFounder models (of varying sizes) against the current leading model in this field broken down by model size.

These results show that pretraining on a large scale of healthy brain MRI data from UKB can significantly improve performance. Other observations can be made as below:

- Across all folds, the BrainSegFounder-Small (BSF-S) framework consistently outperformed the SwinUNETR model. This indicates that the additional training steps taken within the BrainSegFounder framework play a significant role in enhancing its effectiveness in brain tumor segmentation tasks.
- The Small (64M parameters) version of BrainSegFounder achieved higher Dice coefficients on average than the 62M and 69M versions. We believe this indicates that there is an optimal range of model complexity that maximizes performance given certain training data size. Simply increasing the number of parameters does not necessarily lead to better results if the training data does not scale up.
- The one-stage Tiny (62M parameters) model performed comparably to the two-stage BrainSegFounder-Tiny (62M parameters) model, which is notable, implying it did not benefit considerably from the second stage pretraining on the BraTS. This might imply that the UKB dataset alone provides enough variability for effective training. Further study should be made to verify whether the benefit of pretraining on the target datasets can be found using large-scale networks.

Few-shot learning Our experimental results demonstrate the performance capabilities of BrainSegFounder relative to the baseline model, SwinUNETR, under constrained training data conditions. As depicted in Figure 5, BrainSegFounder consistently matched the performance of SwinUNETR across higher levels of available training data and outperformed SwinUNETR when training data was constrained to lower input percentages. As the percentages of training data approached 40% of the input data, both models achieved nearly equivalent accuracy. However, as the amount of training data decreased, BrainSegFounder exhibited superior robustness and adaptability. Notably, at all data availability levels, BrainSegFounder maintained higher mean segmentation accuracy. The results presented in Figure 5 for these input levels are an average of 5 independent subsets of input data to account for variability in small datasets. Supplemental Table 9 contains the results for each of these subsets.

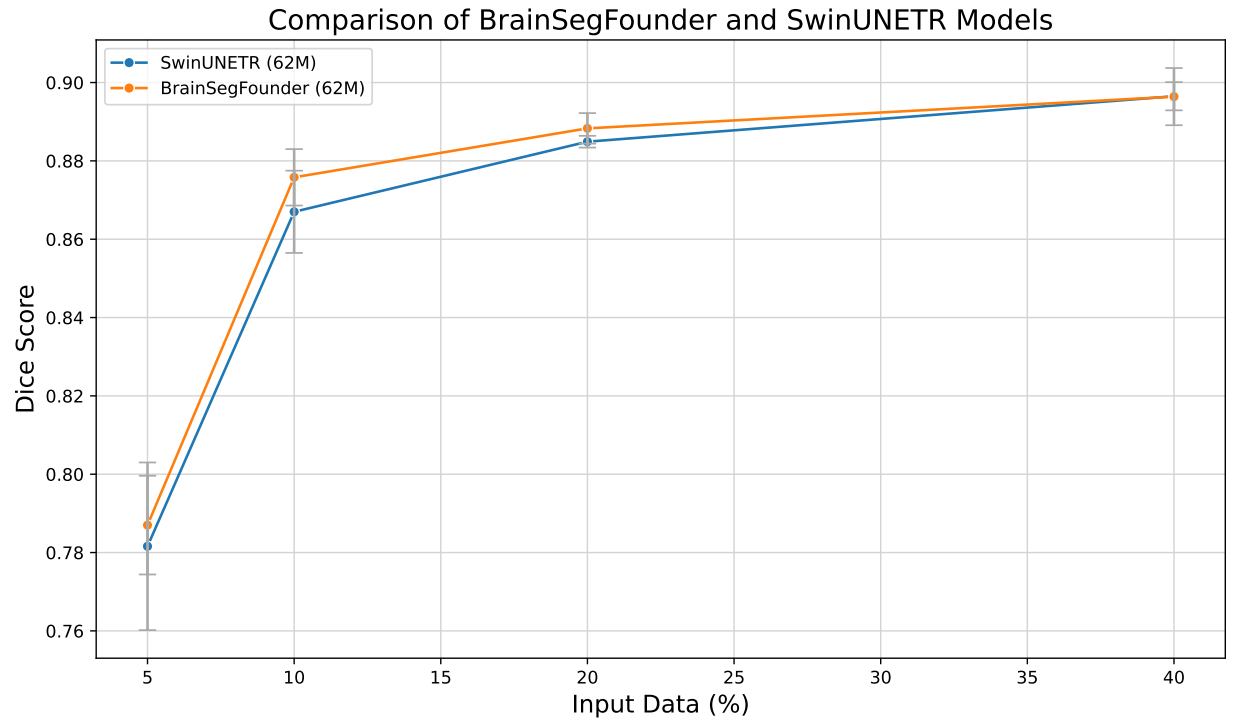


Fig. 5. Dice coefficients for baseline (SwinUNETR) and our model across different levels of training data availability. All models were trained 5 times to account for variability in the input data randomly selected. Error bars represent \pm one standard deviation.

Overall, the BrainSegFounder (64M) model provides the best balance between complexity and performance, as evidenced by its leading average Dice coefficient. These results demonstrate the potential benefits of pretraining on the data from a large number of health subjects.

Modality Restriction Our model can effectively adapt with fewer or single modality data in the downstream task. Table 5 shows the five-fold cross-validation and average Dice scores comparing SSL-pretrained BrainSegFounder and fully supervised SwinUNETR, both using T1-weighted MRI only for all training and inference stages. Our multi-stage pretraining demonstrated a DICE score improvement of 0.04 (6%), indicating the substantial benefit of multi-stage pretraining when the downstream task has limited data modality.

	BrainSegFounder (T1w only)	SwinUNETR (T1w only)
Fold 1	0.718	0.700
Fold 2	0.721	0.672
Fold 3	0.707	0.657
Fold 4	0.731	0.686
Fold 5	0.725	0.676
Average	0.721	0.678

Table 5. Performance comparison of modality restricted models on the BraTS dataset. SwinUNETR is fully supervised learning on T1-weighted MRI without pretraining, while BrainSegFounder uses our multi-stage pretraining on UKB and BraTS T1-weighted MRI and is then finetuned on BraTS T1-weighted MRI.

2.3 Evaluation on ATLAS Challenge Dataset

Our model’s performance on the ATLAS v2.0 dataset was compared against the top-performing models listed on the challenge leaderboard. The results, as summarized in Table 6, demonstrate that our model achieved a Dice score of 0.712, a lesion-wise F1-score of 0.711, simple lesion count of 3.421, and volume difference of 8993.85. These scores would place our model within the top 3 models in the training set leaderboard. Worth noting again is that our training protocol did not include 5-fold cross-validation due to the lack of predetermined folds for which ATLAS can be evaluated in the training set. Top-performing models on the leaderboard³ (CTRL, HeRN, and POBOTRI) were not available for independent validation. Instead, our approach focused on maximizing the comparability with the leaderboard conditions by adhering closely to their reported training setups.

Metric	BrainSegFounder	SwinUNETR	CTRL(*†)	HeRN(*‡)	POBOTRI(*)
Dice (↑)	0.712	0.703	0.663	0.718	0.663
Lesion-wise F1 Score (↑)	0.711	0.703	0.556	0.724	0.559
Simple Lesion Count (↓)	3.421	3.677	4.657	2.750	4.500
Volume Difference (↓)	8993.85	9165.18	8804.91	6162.00	9535.23

Table 6. Performance comparison of segmentation models on the ATLAS v2.0 dataset. All metrics from the challenge (Dice coefficient, Lesion-wise F1 Score, Simple Lesion Count, and Volume Difference) are included for each model for each model. Scores for models marked with an asterisk (*) are sourced directly from the official challenge leaderboard and pertain to their performance on the training set. All of these models utilize ensemble learning methods. CTRL (denoted with †) is the official challenge winner, while HeRN (denoted with ‡) leads on the training set. Top-performing models on the leaderboard (CTRL, HeRN, and POBOTRI) were unavailable for independent validation.

³ <https://atlas.grand-challenge.org/evaluation/lesion-segmentation-hidden-test-set/leaderboard/>

3 Discussion

The findings from our work with BrainSegFounder, especially the "Small" model comprising 64 million parameters, signify a noteworthy progression in 3D foundation models for neuroimage segmentation. The framework’s novel two-stage pretraining strategy—initially utilizing a broad dataset of multi-modal neuroimages from the generally healthy population found in the UK Biobank, followed by training on diseased brain MRI volumes from the BraTS dataset—has demonstrated substantial efficacy. The first stage of pretraining enables the framework to capture the latent representation of normal brain anatomy, a vital aspect for the precise identification of anomalies, like those encountered in brain tumor segmentation tasks. The second stage of pretraining learns the spatial distribution and texture representations of lesions presented in different brain disorders.

Table 4 demonstrates superior performance of BrainSegFounder compared to state-of-the-art approaches including SwinUNETR, nn-UNET, SegResNet, TransBTS, and the brain segmentation foundation model published in Model-Zoo. Each of these models represents a significant advancement in brain segmentation techniques. SwinUNETR is the most direct comparison - its architecture is identical to ours, and outperforming SwinUNETR on these tasks indicates that our multi-stage, self-supervised pretraining method is effective in improving segmentation performance. nn-UNET and SegResNet both were former BraTS challenge winners. As such, they are empirically validated models that excel in the BraTS challenge. By outperforming these models, BrainSegFounder demonstrates its capability to perform segmentation tasks with exceptionally high accuracy and precision. TransBTS utilizes both convolutional and transformer-based architectures to SwinOTH local and global context; our superior results compared to TransBTS demonstrate that our model’s performance increase is not merely due to the inclusion of a transformer based architecture and further validate our pipeline. The brain segmentation foundation model from Model-Zoo serves as a comprehensive pre-trained model designed specifically for brain imaging tasks. By demonstrating increased results compared to this novel foundation model, we show that our methodology could streamline the creation of more effective medical foundation models for segmentation.

One of the key findings is the superior performance of the BrainSegFounder-Small (64M) model over other BrainSegFounder variants. Based on our limited explored range of parameters, our model performs best with an intermediate number of parameters. This suggests that an optimal balance of model complexity and training data is crucial. It is also indicative of the importance of large-scale datasets in training 3D vision foundation models for medical imaging, as even the one-stage pretraining model showed significant effectiveness. However, there is still a possibility to see higher performance using a higher number of parameters and large-scale diverse training data that we did not explore in this work.

Further, the comparable performance of the one-stage 62M model with the two-stage approach indicates that extensive pretraining on a large and diverse dataset like UKB might be sufficient for effective model training, reducing the need for additional pretraining on targeted datasets. This insight could streamline future 3D foundation model development for medical imaging, especially in scenarios where specific pathological datasets are limited or hard to acquire.

Our results from our study limiting training data in few-shot learning indicate that BrainSegFounder’s training methods potentially offer better generalization from limited data, a crucial factor for practical applications in medical imaging where annotated data can be scarce. Though only a slight improvement, the BrainSegFounder consistently outperforms the baseline model at lower levels of input data (see Fig 5 and Supplemental Table 9). Even with incredibly limited data, our Stage 2 self-supervised pretraining serves as a meaningful inclusion in the training pipeline. These findings suggest that the enhancements integrated into BrainSegFounder are effective in optimizing performance under varying data constraints, thereby affirming its suitability for real-world deployment in medical imaging contexts.

In our modality restriction experiment, our model sees a significant reduction in quality when training with fewer modalities. This drop indicates that the multiple modalities present in BraTS contain important information not present in just T1-weighted MRI images about tumor segmentation. However, BrainSegFounder’s better performance under these more challenging scenarios with limited modality input when compared to the base SwinUNETR model validates the feasibility of our extensible approach to handling varying numbers of modalities. When fewer modalities are present, our training scheme still provides valuable information and performance improvements by keeping only the layers trained on the modality present. Similarly, our results using all four modalities present in BraTS suggest our method effectively utilizes infor-

mation given in the pretraining step when presented with additional modalities. Therefore, we conclude that the pretraining steps have a positive effect even when the model is provided with more or less information than is present in the original pretraining stage. Moreover, our results on ATLAS (discussed below) further support our method of handling multiple modalities.

The performance of BrainSegFounder on the ATLAS dataset indicates that its training scheme is generalizable and effective at more than just tumor segmentation, a trait desirable for foundation models. While methods specifically adapted to optimizing results on this dataset do outperform ours, we still maintain third place in the leaderboard. Remarkably, our results were achieved without the use of ensemble learning techniques, which are commonly employed to boost performance by leveraging the strengths of multiple models. The fact that our single-model approach is competitive with ensemble models underscores the robustness and efficiency of our model in managing the intricacies of medical image analysis. We believe that the methodology used for BrainSegFounder can be refined and extended to move towards a Medical Foundation Model for neuroimages.

In addition, BrainSegFounder’s training scheme and model provide a clear advantage over SAM and MedSAM, two powerful existing segmentation foundation models. (1) While SAM is restricted to 2D RGB images, BrainSegFounder is designed to handle 3D medical images with any number of channels as input, providing greater versatility in medical imaging applications. (2) MedSAM requires bounding-box input prompts and its 3D functionality is limited to manually uploading each image to a plugin for prompting and slice-by-slice annotation. Both methods require manual input. In contrast, our model eliminates the need for such manual interventions once trained, streamlining the segmentation process. (3) Although SAM is capable of automated segmentation without input, it lacks the ability to specify a fixed number of classes and instead generates an arbitrary number of classes; this property leads sub-optimally on medical images with specific segmentation tasks (e.g., lesion detection), and cannot be used without additional human input. (4) Neither SAM nor MedSAM efficiently process multimodal data as they generate predictions from a single scan, whereas BrainSegFounder is designed to integrate multiple scans from the same individual.

However, it’s important to note that while BrainSegFounder shows promise in brain tumor segmentation and brain region segmentation, its application in other neuroimaging tasks remains to be explored. One such task is brain tissue segmentation - a common task in automated analysis. Future research should investigate its adaptability to other neurological conditions, its performance in different clinical environments, and its usefulness in additional common analysis tasks.

In conclusion, BrainSegFounder is a significant step forward 3D foundation models for medical image segmentation and analysis, particularly for multi-modal neuroimaging. Its development underscores the potential of AI and foundation models in enhancing diagnostic accuracy and efficiency, paving the way for more advanced, adaptable, and robust AI tools in healthcare.

4 Acknowledgments

This work was partially supported by the National Science Foundation (1908299, 2318984, 2123809) and the National Institute of Aging of the National Institutes of Health (RF1AG071469). This research has been conducted using the UK Biobank Resource under application number 48388. We extend our appreciation to the contributors of the UK Biobank, BraTS, and ALTAS datasets, whose extensive data collection efforts have been fundamental to the success of this research. We gratefully acknowledge NVIDIA AI Technology (NVAITC)’s support for this research project, especially on GPU parallelization technology. Special thanks to the University of Florida’s HiPerGator-AI supercomputer team for providing the computational resources essential in training the BrainSegFounder models. We appreciate the support from Ying Zhang (UFIT Research Computing) and Kaleb Smith (NVAITC) in providing the necessary computing resources and technology support.

References

- [1] Y. Chen, C. B. Schonlieb, P. Lio, *et al.*, “AI-Based Reconstruction for Fast MRI-A Systematic Review and Meta-Analysis,” *Proceedings of the IEEE*, vol. 110, no. 2, pp. 224–245, Feb. 2022, ISSN: 0018-9219. DOI: [10.1109/JPROC.2022.3141367](https://doi.org/10.1109/JPROC.2022.3141367). [Online]. Available: <http://www.scopus.com/inward/record.url?scp=85124598976&partnerID=8YFLogxK> (visited on 01/14/2024).

- [2] A. Segato, A. Marzullo, F. Calimeri, and E. De Momi, “Artificial intelligence for brain diseases: A systematic review,” *APL Bioengineering*, vol. 4, no. 4, p. 041503, Oct. 2020, ISSN: 2473-2877. DOI: [10.1063/5.0011697](https://doi.org/10.1063/5.0011697). [Online]. Available: <https://doi.org/10.1063/5.0011697> (visited on 01/14/2024).
- [3] R. P. N. Rao, “Brain Co-processors: Using AI to Restore and Augment Brain Function,” en, in *Handbook of Neuroengineering*, N. V. Thakor, Ed., Singapore: Springer Nature, 2023, pp. 1225–1260, ISBN: 9789811655401. DOI: [10.1007/978-981-16-5540-1_32](https://doi.org/10.1007/978-981-16-5540-1_32). [Online]. Available: https://doi.org/10.1007/978-981-16-5540-1_32 (visited on 01/14/2024).
- [4] M. O. Owolabi, M. Leonardi, C. Bassetti, *et al.*, “Global synergistic actions to improve brain health for human development,” en, *Nature Reviews Neurology*, vol. 19, no. 6, pp. 371–383, Jun. 2023, Number: 6 Publisher: Nature Publishing Group, ISSN: 1759-4766. DOI: [10.1038/s41582-023-00808-z](https://doi.org/10.1038/s41582-023-00808-z). [Online]. Available: <https://www.nature.com/articles/s41582-023-00808-z> (visited on 01/14/2024).
- [5] D. Moreno-Blanco, J. Solana-Sanchez, P. Sanchez-Gonzalez, *et al.*, “Technologies for Monitoring Lifestyle Habits Related to Brain Health: A Systematic Review,” en, *Sensors*, vol. 19, no. 19, p. 4183, Jan. 2019, Number: 19 Publisher: Multidisciplinary Digital Publishing Institute, ISSN: 1424-8220. DOI: [10.3390/s19194183](https://doi.org/10.3390/s19194183). [Online]. Available: <https://www.mdpi.com/1424-8220/19/19/4183> (visited on 01/14/2024).
- [6] P. Rajpurkar, E. Chen, O. Banerjee, and E. J. Topol, “AI in health and medicine,” en, *Nature Medicine*, vol. 28, no. 1, pp. 31–38, Jan. 2022, Number: 1 Publisher: Nature Publishing Group, ISSN: 1546-170X. DOI: [10.1038/s41591-021-01614-0](https://doi.org/10.1038/s41591-021-01614-0). [Online]. Available: <https://www.nature.com/articles/s41591-021-01614-0> (visited on 01/14/2024).
- [7] A. S. Khachaturian, A. Dengel, V. Dočkal, P. Hroboň, and M. Tolar, “Accelerating Innovations for Enhanced Brain Health. Can Artificial Intelligence Advance New Pathways for Drug Discovery for Alzheimer’s and other Neurodegenerative Disorders?” en, *The Journal of Prevention of Alzheimer’s Disease*, vol. 10, no. 1, pp. 1–4, Jan. 2023, ISSN: 2426-0266. DOI: [10.14283/jpad.2023.1](https://doi.org/10.14283/jpad.2023.1). [Online]. Available: <https://doi.org/10.14283/jpad.2023.1> (visited on 01/14/2024).
- [8] M. Moor, O. Banerjee, Z. S. H. Abad, *et al.*, “Foundation models for generalist medical artificial intelligence,” en, *Nature*, vol. 616, no. 7956, pp. 259–265, Apr. 2023, Number: 7956 Publisher: Nature Publishing Group, ISSN: 1476-4687. DOI: [10.1038/s41586-023-05881-4](https://doi.org/10.1038/s41586-023-05881-4). [Online]. Available: <https://www.nature.com/articles/s41586-023-05881-4> (visited on 01/14/2024).
- [9] B. Azad, R. Azad, S. Eskandari, *et al.*, *Foundational Models in Medical Imaging: A Comprehensive Survey and Future Vision*, arXiv:2310.18689 [cs], Oct. 2023. DOI: [10.48550/arXiv.2310.18689](https://doi.org/10.48550/arXiv.2310.18689). [Online]. Available: <http://arxiv.org/abs/2310.18689> (visited on 01/14/2024).
- [10] S. Zhang and D. Metaxas, “On the challenges and perspectives of foundation models for medical image analysis,” eng, *Medical Image Analysis*, vol. 91, p. 102996, Jan. 2024, ISSN: 1361-8423. DOI: [10.1016/j.media.2023.102996](https://doi.org/10.1016/j.media.2023.102996).
- [11] F. Isensee, P. F. Jaeger, S. A. A. Kohl, J. Petersen, and K. H. Maier-Hein, “nnU-Net: A self-configuring method for deep learning-based biomedical image segmentation,” en, *Nature Methods*, vol. 18, no. 2, pp. 203–211, Feb. 2021, Number: 2 Publisher: Nature Publishing Group, ISSN: 1548-7105. DOI: [10.1038/s41592-020-01008-z](https://doi.org/10.1038/s41592-020-01008-z). [Online]. Available: <https://www.nature.com/articles/s41592-020-01008-z> (visited on 01/14/2024).
- [12] R. McKinley, R. Meier, and R. Wiest, “Ensembles of Densely-Connected CNNs with Label-Uncertainty for Brain Tumor Segmentation,” en, in *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries*, A. Crimi, S. Bakas, H. Kuijff, F. Keyvan, M. Reyes, and T. van Walsum, Eds., ser. Lecture Notes in Computer Science, Cham: Springer International Publishing, 2019, pp. 456–465, ISBN: 978-3-030-11726-9. DOI: [10.1007/978-3-030-11726-9_40](https://doi.org/10.1007/978-3-030-11726-9_40).
- [13] K. Kamnitsas, C. Ledig, V. F. J. Newcombe, *et al.*, “Efficient multi-scale 3D CNN with fully connected CRF for accurate brain lesion segmentation,” eng, *Medical Image Analysis*, vol. 36, pp. 61–78, Feb. 2017, ISSN: 1361-8423. DOI: [10.1016/j.media.2016.10.004](https://doi.org/10.1016/j.media.2016.10.004).
- [14] U. Baid, S. Ghodasara, S. Mohan, *et al.*, *The RSNA-ASNR-MICCAI BraTS 2021 Benchmark on Brain Tumor Segmentation and Radiogenomic Classification*, arXiv:2107.02314 [cs], Sep. 2021. DOI: [10.48550/arXiv.2107.02314](https://doi.org/10.48550/arXiv.2107.02314). [Online]. Available: <http://arxiv.org/abs/2107.02314> (visited on 01/14/2024).

- [15] M. Antonelli, A. Reinke, S. Bakas, *et al.*, “The Medical Segmentation Decathlon,” en, *Nature Communications*, vol. 13, no. 1, p. 4128, Jul. 2022, Number: 1 Publisher: Nature Publishing Group, ISSN: 2041-1723. DOI: [10.1038/s41467-022-30695-9](https://doi.org/10.1038/s41467-022-30695-9). [Online]. Available: <https://www.nature.com/articles/s41467-022-30695-9> (visited on 01/14/2024).
- [16] F. e. a. Quinton, “A tumor and liver automatic segmentation challenge,” 2023. DOI: [10.5281/zenodo.7835370](https://doi.org/10.5281/zenodo.7835370).
- [17] Z. Zhou, V. Sodha, J. Pang, M. B. Gotway, and J. Liang, “Models Genesis,” *Medical Image Analysis*, vol. 67, p. 101840, Jan. 2021, ISSN: 1361-8415. DOI: [10.1016/j.media.2020.101840](https://doi.org/10.1016/j.media.2020.101840). [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1361841520302048> (visited on 01/14/2024).
- [18] Y. Tang, D. Yang, W. Li, *et al.*, *Self-Supervised Pre-Training of Swin Transformers for 3D Medical Image Analysis*, en, arXiv:2111.14791 [cs], Mar. 2022. [Online]. Available: <http://arxiv.org/abs/2111.14791> (visited on 01/14/2024).
- [19] X. Mei, Z. Liu, P. M. Robson, *et al.*, “RadImageNet: An Open Radiologic Deep Learning Research Dataset for Effective Transfer Learning,” *Radiology: Artificial Intelligence*, vol. 4, no. 5, e210315, Sep. 2022, Publisher: Radiological Society of North America. DOI: [10.1148/ryai.210315](https://doi.org/10.1148/ryai.210315). [Online]. Available: <https://pubs.rsna.org/doi/full/10.1148/ryai.210315> (visited on 01/14/2024).
- [20] K. Clark, B. Vendt, K. Smith, *et al.*, “The Cancer Imaging Archive (TCIA): Maintaining and Operating a Public Information Repository,” en, *Journal of Digital Imaging*, vol. 26, no. 6, pp. 1045–1057, Dec. 2013, ISSN: 1618-727X. DOI: [10.1007/s10278-013-9622-7](https://doi.org/10.1007/s10278-013-9622-7). [Online]. Available: <https://doi.org/10.1007/s10278-013-9622-7> (visited on 01/14/2024).
- [21] C. Bycroft, C. Freeman, D. Petkova, *et al.*, “The UK Biobank resource with deep phenotyping and genomic data,” en, *Nature*, vol. 562, no. 7726, pp. 203–209, Oct. 2018, Number: 7726 Publisher: Nature Publishing Group, ISSN: 1476-4687. DOI: [10.1038/s41586-018-0579-z](https://doi.org/10.1038/s41586-018-0579-z). [Online]. Available: <https://www.nature.com/articles/s41586-018-0579-z>. (visited on 01/14/2024).
- [22] A. Kirillov, E. Mintun, N. Ravi, *et al.*, *Segment Anything*, arXiv:2304.02643 [cs], Apr. 2023. DOI: [10.48550/arXiv.2304.02643](https://doi.org/10.48550/arXiv.2304.02643). [Online]. Available: <http://arxiv.org/abs/2304.02643> (visited on 01/14/2024).
- [23] H. Touvron, T. Lavril, G. Izacard, *et al.*, *LLaMA: Open and Efficient Foundation Language Models*, arXiv:2302.13971 [cs], Feb. 2023. DOI: [10.48550/arXiv.2302.13971](https://doi.org/10.48550/arXiv.2302.13971). [Online]. Available: <http://arxiv.org/abs/2302.13971> (visited on 01/14/2024).
- [24] M. Y. Lu, B. Chen, A. Zhang, *et al.*, *Visual Language Pretrained Multiple Instance Zero-Shot Transfer for Histopathology Images*, arXiv:2306.07831 [cs], Jun. 2023. DOI: [10.48550/arXiv.2306.07831](https://doi.org/10.48550/arXiv.2306.07831). [Online]. Available: <http://arxiv.org/abs/2306.07831> (visited on 01/14/2024).
- [25] S. Bannur, S. Hyland, Q. Liu, *et al.*, *Learning to Exploit Temporal Structure for Biomedical Vision-Language Processing*, arXiv:2301.04558 [cs], Mar. 2023. DOI: [10.48550/arXiv.2301.04558](https://doi.org/10.48550/arXiv.2301.04558). [Online]. Available: <http://arxiv.org/abs/2301.04558> (visited on 01/14/2024).
- [26] E. Tiu, E. Talus, P. Patel, C. P. Langlotz, A. Y. Ng, and P. Rajpurkar, “Expert-level detection of pathologies from unannotated chest X-ray images via self-supervised learning,” eng, *Nature Biomedical Engineering*, vol. 6, no. 12, pp. 1399–1406, Dec. 2022, ISSN: 2157-846X. DOI: [10.1038/s41551-022-00936-9](https://doi.org/10.1038/s41551-022-00936-9).
- [27] J. Ma, Y. He, F. Li, L. Han, C. You, and B. Wang, *Segment Anything in Medical Images*, arXiv:2304.12306 [cs, eess], Jul. 2023. DOI: [10.48550/arXiv.2304.12306](https://doi.org/10.48550/arXiv.2304.12306). [Online]. Available: <http://arxiv.org/abs/2304.12306> (visited on 01/14/2024).
- [28] W. Lei, X. Wei, X. Zhang, K. Li, and S. Zhang, *MedLSAM: Localize and Segment Anything Model for 3D CT Images*, arXiv:2306.14752 [cs], Nov. 2023. DOI: [10.48550/arXiv.2306.14752](https://doi.org/10.48550/arXiv.2306.14752). [Online]. Available: <http://arxiv.org/abs/2306.14752> (visited on 01/14/2024).
- [29] J. Cheng, J. Ye, Z. Deng, *et al.*, *SAM-Med2D*, arXiv:2308.16184 [cs], Aug. 2023. DOI: [10.48550/arXiv.2308.16184](https://doi.org/10.48550/arXiv.2308.16184). [Online]. Available: <http://arxiv.org/abs/2308.16184> (visited on 01/14/2024).
- [30] S. e. a. Gong, “3dsam-adapter: Holistic adaptation of sam from 2d to 3d for promptable medical image segmentation,” 2023. DOI: [10.48550/arXiv.2306.13465](https://doi.org/10.48550/arXiv.2306.13465).
- [31] A. Hatamizadeh, V. Nath, Y. Tang, D. Yang, H. Roth, and D. Xu, *Swin UNETR: Swin Transformers for Semantic Segmentation of Brain Tumors in MRI Images*, arXiv:2201.01266 [cs, eess], Jan. 2022.

- DOI: [10.48550/arXiv.2201.01266](https://doi.org/10.48550/arXiv.2201.01266). [Online]. Available: <http://arxiv.org/abs/2201.01266> (visited on 01/14/2024).
- [32] T. J. Littlejohns, J. Holliday, L. M. Gibson, *et al.*, “The UK Biobank imaging enhancement of 100,000 participants: Rationale, data collection, management and future directions,” en, *Nature Communications*, vol. 11, no. 1, p. 2624, May 2020, Number: 1 Publisher: Nature Publishing Group, ISSN: 2041-1723. DOI: [10.1038/s41467-020-15948-9](https://doi.org/10.1038/s41467-020-15948-9). [Online]. Available: <https://www.nature.com/articles/s41467-020-15948-9> (visited on 01/14/2024).
 - [33] S. M. Smith, A.-A. Fidel, and M. L. Karla, *UK Biobank Brain Imaging Documentation, Version 1.9*, English, Sep. 2022. [Online]. Available: https://biobank.ctsu.ox.ac.uk/crystal/crystal/docs/brain_mri.pdf.
 - [34] M. W. Woolrich, S. Jbabdi, B. Patenaude, *et al.*, “Bayesian analysis of neuroimaging data in FSL,” eng, *NeuroImage*, vol. 45, no. 1 Suppl, S173–186, Mar. 2009, ISSN: 1095-9572. DOI: [10.1016/j.neuroimage.2008.10.055](https://doi.org/10.1016/j.neuroimage.2008.10.055).
 - [35] X. Li, P. S. Morgan, J. Ashburner, J. Smith, and C. Rorden, “The first step for neuroimaging data analysis: DICOM to NIfTI conversion,” eng, *Journal of Neuroscience Methods*, vol. 264, pp. 47–56, May 2016, ISSN: 1872-678X. DOI: [10.1016/j.jneumeth.2016.03.001](https://doi.org/10.1016/j.jneumeth.2016.03.001).
 - [36] C. Davatzikos, S. Rathore, S. Bakas, *et al.*, “Cancer imaging phenomics toolkit: Quantitative imaging analytics for precision diagnostics and predictive modeling of clinical outcome,” eng, *Journal of Medical Imaging (Bellingham, Wash.)*, vol. 5, no. 1, p. 011018, Jan. 2018, ISSN: 2329-4302. DOI: [10.1117/1.JMI.5.1.011018](https://doi.org/10.1117/1.JMI.5.1.011018).
 - [37] S. K. Warfield, K. H. Zou, and W. M. Wells, “Simultaneous truth and performance level estimation (STAPLE): An algorithm for the validation of image segmentation,” eng, *IEEE transactions on medical imaging*, vol. 23, no. 7, pp. 903–921, Jul. 2004, ISSN: 0278-0062. DOI: [10.1109/TMI.2004.828354](https://doi.org/10.1109/TMI.2004.828354).
 - [38] S.-L. Liew, B. P. Lo, M. R. Donnelly, *et al.*, “A large, curated, open-source stroke neuroimaging dataset to improve lesion segmentation algorithms,” *Scientific Data*, vol. 9, no. 1, p. 320, Jun. 2022, ISSN: 2052-4463. DOI: [10.1038/s41597-022-01401-7](https://doi.org/10.1038/s41597-022-01401-7). [Online]. Available: <https://doi.org/10.1038/s41597-022-01401-7>.
 - [39] P. A. Yushkevich, J. Piven, H. C. Hazlett, *et al.*, “User-guided 3D active contour segmentation of anatomical structures: Significantly improved efficiency and reliability,” *NeuroImage*, vol. 31, no. 3, pp. 1116–1128, Jul. 2006, ISSN: 1053-8119. DOI: [10.1016/j.neuroimage.2006.01.015](https://doi.org/10.1016/j.neuroimage.2006.01.015). [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1053811906000632> (visited on 04/18/2024).
 - [40] P. A. Yushkevich and G. Gerig, “ITK-SNAP: An Intractive Medical Image Segmentation Tool to Meet the Need for Expert-Guided Segmentation of Complex Medical Images,” *IEEE Pulse*, vol. 8, no. 4, pp. 54–57, Jul. 2017, Conference Name: IEEE Pulse, ISSN: 2154-2317. DOI: [10.1109/MPUL.2017.2701493](https://doi.org/10.1109/MPUL.2017.2701493). [Online]. Available: <https://ieeexplore.ieee.org/document/7979667> (visited on 04/18/2024).
 - [41] K. He, X. Zhang, S. Ren, and J. Sun, *Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification*, arXiv:1502.01852 [cs] version: 1, Feb. 2015. DOI: [10.48550/arXiv.1502.01852](https://doi.org/10.48550/arXiv.1502.01852). [Online]. Available: <http://arxiv.org/abs/1502.01852> (visited on 01/14/2024).
 - [42] W. Wang, C. Chen, M. Ding, J. Li, H. Yu, and S. Zha, *TransBTS: Multimodal Brain Tumor Segmentation Using Transformer*, arXiv:2103.04430 [cs], Jun. 2021. DOI: [10.48550/arXiv.2103.04430](https://doi.org/10.48550/arXiv.2103.04430). [Online]. Available: <http://arxiv.org/abs/2103.04430> (visited on 07/21/2024).
 - [43] A. Myronenko, *3D MRI brain tumor segmentation using autoencoder regularization*, arXiv:1810.11654 [cs, q-bio], Nov. 2018. DOI: [10.48550/arXiv.1810.11654](https://doi.org/10.48550/arXiv.1810.11654). [Online]. Available: <http://arxiv.org/abs/1810.11654> (visited on 07/21/2024).
 - [44] P. MONAI, *Project-MONAI/model-zoo*, original-date: 2022-05-09T03:40:16Z, Jul. 2024. [Online]. Available: <https://github.com/Project-MONAI/model-zoo> (visited on 07/21/2024).

5 Appendix

5.1 UK Biobank Data

Table 7 presents a comprehensive summary of the participants used from the UK Biobank.

	Entire UK Biobank	40K Cohort (%)
Age at Recruitment		
≤ 45	51,763 (10.3%)	4,601 (11.1%)
45-49	51,866 (10.3%)	5,031 (12.2%)
50-54	74,387 (14.8%)	7,563 (18.3%)
55-59	86,899 (17.3%)	8,819 (21.3%)
60-64	118,959 (23.7%)	9,579 (23.1%)
≥ 65	118,435 (23.6%)	5,796 (14.0%)
Unanswered	101 (0.02%)	11 (0.03%)
Gender		
Male	229,051 (45.6%)	19,497 (47.1%)
Female	273,258 (54.4%)	21,891 (52.9%)
Unanswered	101 (0.02%)	12 (0.03%)
Race		
White	472521 (94.1%)	40057 (96.8%)
Mixed	2953 (0.6%)	190 (0.5%)
Asian	11447 (2.3%)	541 (1.3%)
Black	8055 (1.6%)	268 (0.7%)
Other	4555 (0.9%)	223 (0.5%)
Unanswered	2,778 (0.6%)	121 (0.3%)
Data Information		
# Samples	502,309	41,400
# Brain Tumors	1,210	42

Table 7. UKB Data Demographic information.

	SwinUNETR (62M)	BSF-T (62M)	BSF-S (64M)	BSF-B (69M)	One-Stage (62M)	One-Stage (64M)	One-Stage (69M)
Fold 1	0.8854	0.9027	0.9032	0.9014	0.9019	0.8994	0.8999
Fold 2	0.9059	0.9181	0.9182	0.9164	0.9188	0.9186	0.9055
Fold 3	0.8981	0.9102	0.9121	0.9097	0.9119	0.9125	0.9002
Fold 4	0.8924	0.9103	0.9100	0.9070	0.9107	0.9133	0.9109
Fold 5	0.9035	0.9139	0.9141	0.9101	0.9132	0.9114	0.9103
Average	0.8971	0.9110	0.9115	0.9089	0.9112	0.9110	0.9054

Table 8. Comparison of BrainSegFounder models through 5-fold cross-validation with metric Dice coefficient on BraTS. SwinUNETR is the winning solution on BraTS challenge 2021, which is performed with fully supervised learning without UKB pretraining. BrainSegFounder is the proposed method, which is conducted with the two-stage pretraining and then finetuning on the target dataset. The one-stage means that pretraining on UKB is performed but not on the BraTS. Note: The performance results for SwinUNETR were published on the official GitHub, utilizing hyper-parameter settings similar to those in our finetuning stage but without implementing the ensembling approach that was described in the published work.

5.2 Fold-wise comparison of all models.

Table 8 provides fold-wise comparison of our BrainSegFounder models across all tested parameters.

Table 9 presents a comparison of few-shot learning Dice scores on the testing set at varying levels of input training data.

Repeat	5%		10%		20%		40%	
	BSF	SwinU	BSF	SwinU	BSF	SwinU	BSF	SwinU
1	0.7810	0.7437	0.8771	0.8594	0.8893	0.8837	0.8949	0.8885
2	0.7912	0.7899	0.8764	0.8553	0.8814	0.8834	0.8981	0.8981
3	0.7797	0.7886	0.8683	0.8812	0.8901	0.8869	0.8906	0.8956
4	0.8071	0.7966	0.8705	0.8735	0.8911	0.8844	0.9048	0.8938
5	0.7758	0.7893	0.8869	0.8656	0.8895	0.8860	0.8857	0.8938
Average	0.7870	0.7816	0.8758	0.8670	0.8883	0.8849	0.8948	0.8937

Table 9. Comparison of BrainSegFounder (BSF) and SwinUNETR (SwinU) Baseline models trained on 5 repeats of varying percentages of the input data. Data was randomly sampled from the BraTS training dataset, and models were evaluated on the testing dataset.