# Domain Aware Multi-Task Pretraining of 3D Swin Transformer for T1-weighted Brain MRI

Jonghun Kim[1,2⋆], Mansu Kim[3⋆], and Hyunjin Park[1,2†]

[1] Department of Electrical and Computer Engineering,
Sungkyunkwan University, Suwon, Republic of Korea
[2] Center for Neuroscience Imaging Research,
Institute for Basic Science, Suwon, Republic of Korea
[3] AI Graduate School, Gwanju Institute of Science and Technology,
Gwangju, Republic of Korea
{iproj2,hyunjinp}@skku.edu, mansu.kim@gist.ac.kr

**Abstract.** The scarcity of annotated medical images is a major bottleneck in developing learning models for medical image analysis. Hence, recent studies have focused on pretrained models with fewer annotation requirements that can be fine-tuned for various downstream tasks. However, existing approaches are mainly 3D adaptions of 2D approaches ill-suited for 3D medical imaging data. Motivated by this gap, we propose novel domain-aware multi-task learning tasks to pretrain a 3D Swin Transformer for brain magnetic resonance imaging (MRI). Our method considers the domain knowledge in brain MRI by incorporating brain anatomy and morphology as well as standard pretext tasks adapted for 3D imaging in a contrastive learning setting. We pretrain our model using large-scale brain MRI data of 13,687 samples spanning several large-scale databases. Our method outperforms existing supervised and self-supervised methods in three downstream tasks of Alzheimer's disease classification, Parkinson's disease classification, and age prediction tasks. The ablation study of the proposed pretext tasks shows the effectiveness of our pretext tasks. Our code is available at github.com/jongdory/DAMT.
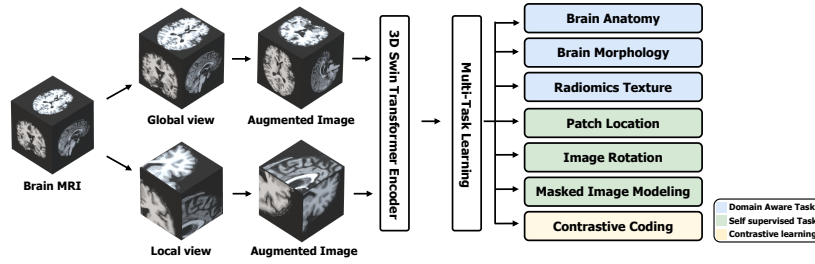
**Keywords:** Self supervised learning · Magnetic Resonance Imaging · Swin Transformer, 3D Medical Image Analysis

## 1 Introduction

The recent success of neural networks in computer vision has prompted researchers to explore new network models for medical imaging tasks. For example, tasks involving tumor region segmentation and disease classification have been performed using various modalities, including magnetic resonance imaging (MRI) and computed tomography (CT), typically using supervised learning methods [31, 35, 69, 89, 92]. One prominent model in this domain is the Vision

---

⋆ Equal Contribution
† Corresponding Author

**Fig. 1:** Overview of our proposed multi-task pretraining framework. The original MR image is divided into global and local views. Augmentation is then performed by applying masking and rotation, followed by feeding into the Swin Transformer. The process shows that the encoder learns features through seven pretext tasks.

Transformer (ViT), as introduced by [23], which has revolutionized the fields of computer vision and hence medical image analysis. ViTs particularly excel at learning pretext tasks, providing scalability for large-scale training [91], and enabling efficient gathering of both global and local information. Unlike convolutional neural networks (CNNs), which have limited receptive fields, ViTs encode visual representations from a sequence of patches and employ self-attention mechanisms to model long-range global information [47, 48, 70]. The ViT architecture has been extended to accommodate three-dimensional (3D) images, including those used in medical imaging and video analysis [76, 77]. However, there are a few key parameters to consider before applying ViT to 3D images, optimizing the trade-off between patch size and sequence length is important. For instance, in scenarios with larger patch sizes, the information capacity of each patch might become suboptimal, resulting in potential information loss. Conversely, using smaller patch size settings leads to a cubic increase in the sequence length, resulting in an exponential increase in computational cost. To address these challenges, the Swin transformer, proposed by [53], offers an efficient solution for the processing of 3D data. It is known for boosting robustness to varying patch sizes and sequence lengths, besides the reduction in computational costs and improved space efficiency.

Pretraining strategies are widely employed in both natural and medical image analyses to enhance model performance. Given the time-consuming and expensive process of annotating medical images, training with a limited number of labeled samples is essential for effective medical image analyses. The conventional approach to pretraining involves supervised pretraining using large labeled datasets of natural images, such as ImageNet [38, 39]. However, applying two-dimensional (2D) image-based neural network models to 3D medical imaging poses several challenges. This is due to the significant domain gap between natural images and medical imaging modalities, such as MRI and CT. Additionally, the absence of cross-plane contextual information in 3D images further complicates the process. Consequently, selecting suitable supervised tasks and developing domain-specific pretraining tasks remain significant challenges when effectively training models for 3D medical imaging [52]. Self-supervised learning (SSL) tasks represent effective approaches for learning useful representations

from unlabeled data. These tasks have proven successful in computer vision for pretraining models capable of learning general features applicable to a broad range of downstream tasks [14, 22, 62, 63, 66, 95]. Nevertheless, depending solely on these methods can result in the acquisition of irrelevant features. Consequently, domain-aware tasks are necessary during pretraining to facilitate the learning of crucial image features. SSL tasks can be easily trained on natural images owing to the wide availability of large databases. However, despite the scarcity of such databases in the medical imaging domain, we harnessed most of the available large-scale brain MRI databases, totaling 13,687 scans, to empower our approach. In this study, we propose domain-specific self-supervised tasks that leverage expertise in brain imaging and apply them to pretrain a Swin transformer. Inspired by previous research, [9, 83], we have designed transformations suitable for 3D medical images and applied them for pretraining. The proposed self-supervised tasks encourage the model to learn representations related to the general brain anatomy and morphological characteristics.
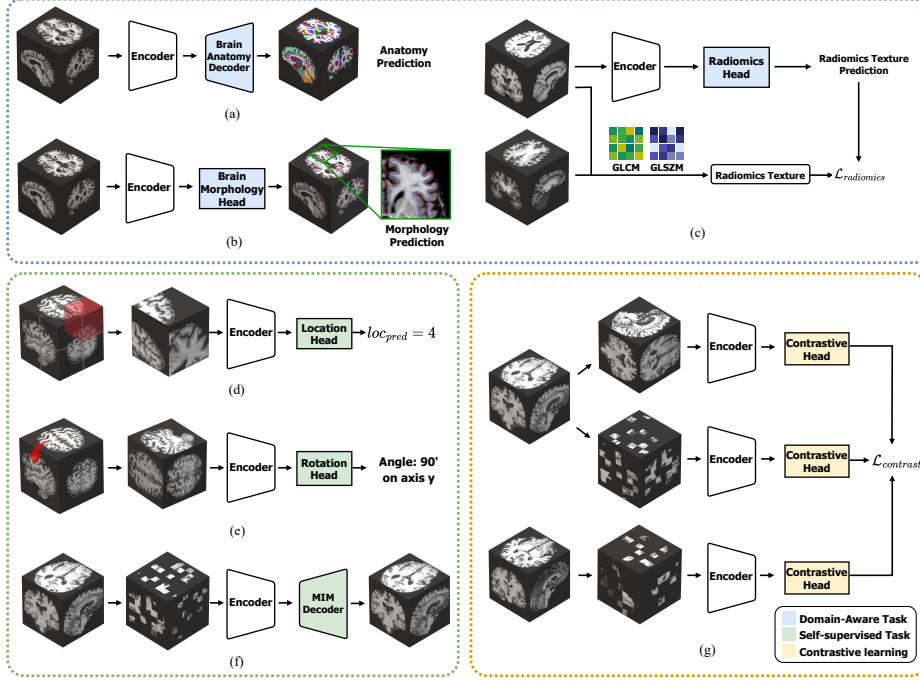
**Contribution**:

- We present a novel multi-task pretraining framework that leverages domain-specific knowledge of brain anatomy and related morphological features. This framework incorporates several self-supervised tasks, including image rotation, patch location, and masked image modeling within the contrastive learning setup.
- We successfully pretrain a Swin transformer on 3D brain T1-weighted MRI images using the proposed pretext tasks. We perform experiments on large-scale brain MRI dataset (n = 13,687) to demonstrate improvement of our pretraining strategy over the competing methods.
- We demonstrate the clinical benefits of our pretrained model, such as accurate diagnosis of Alzheimer's disease (AD) and Parkinson's disease (PD), as well as predicting the chronological age.

## 2 Related Work

### 2.1 Self-Supervised Learning

A general representation can be learned in an embedding space derived from a high-dimensional input. The objective is to enhance the similarity between semantically related data samples and increase the distance between dissimilar data samples. SSL leverages pretext tasks, also known as proxy tasks, such as solving jigsaw puzzles, memorizing the spatial context from images, predicting image rotation, colorization, and restoring images to learn feature representations [22, 62, 63, 66, 95]. Other studies have employed contrastive learning approaches, such as the simple framework for contrastive learning of visual representations (SimCLR) [14], momentum contrast (MoCo) [15, 37], bootstrap your own latent (BYOL) [33], and self-distillation with no labels (DINO) [9] to learn the representation effectively. Recent research has focused on SSL by masking and restoring random patches. Masked autoencoder (MAE) [36] pretrains ViT

**Fig. 2:** Detailed illustration of each pretext task in our proposed approach. (a) Brain Anatomy: predicting the parcellation of the input brain image. (b) Brain Morphology: predicting morphology, such as thickness or curvature, of the input brain image. (c) Radiomics Texture Prediction: predicting radiomics texture in the white matter, gray matter, and CSF regions. (d) Patch Location: identifying the position of the patch in the local view. (e) Image Rotation: rotating the original image and determining the corresponding rotation. (f) Masked Image Modeling: the original image is cut out and reconstructed back to its original form. (g) Contrastive Learning: different augmentations applied to the same patch are pulled closer as positive pairs and inputs from different images are pushed away as negative pairs.

by randomly masking and restoring images, which leads to masked image modeling (MIM). SimMIM [90] is a simplified version of MIM that can be applied to Swin transformers.

## 2.2  Pretraining for Medical Image

The aforementioned pretext tasks for SSL have successfully learned representations in 2D natural images, as well as in some 3D medical images in a 2D manner [3, 12]. For instance, one study employs a task involving the ordering of 2D axial slices in 3D CT and MR images, resulting in improved body part recognition [94]. Another study proposed a task that predicted the distance between 2D patches in 3D brain images, which was effective for brain tissue segmentation [80]. Tang et al. [83] pretrained a transformer on 3D medical images by simultaneously predicting rotation, inpainting reconstruction, and contrastive learning. I3D [10] is a 3D CNN model pretrained with a kinetics dataset for

action recognition, and attempts have been made to apply it to medical images [43]. However, due to the domain gap between natural and medical images, the pretrained models are not fully suitable for the medical domain. Existing studies have primarily achieved success by applying or extending 2D pretext tasks in a 3D context. However, these tasks were originally designed to address computer vision challenges in 2D natural images and may not be optimal for learning the complex anatomical and morphological properties of brain images. Hence, we focus on introducing brain imaging-specific pretasks to incorporate brain anatomy and morphological characteristics. The aim is to successfully learn high-level representations relevant to brain structure and functions.
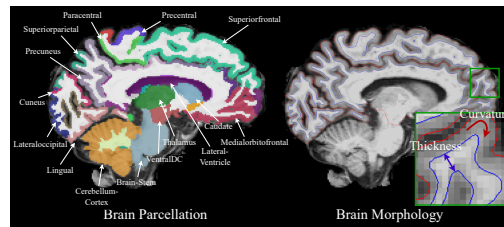
## 3   Methodology

In this section, we provide an overview of the formulations of the self-supervised pretext tasks. These tasks are designed to facilitate the learning of effective data representations $z$ in a 3D context. Furthermore, they enable the model to comprehend the complex brain anatomy and morphology from unlabeled 3D image samples during the pretraining phase. Inspired by the augmentation method introduced in a previous study [9], we incorporated augmentation into our approach. In general, we divided pretraining into two views: global and local. The global view focuses on capturing the overall image structure. In contrast, the local view is designed to facilitate the learning of localized features in the brain. The local view is a subset of the global view, and both views undergo rotations and intensity shifting. Fig. 1 depicts the data augmentation process and the proposed multi-tasking framework. Additionally, all processes and multi-tasks were applied concurrently and learned simultaneously. Fig. 2 illustrates the details of the various pretext tasks, which will be explained in the following sections.

### 3.1   Domain Aware Tasks

**Brain Anatomy Prediction.** Brain parcellation is a crucial neuroscience technique that involves dividing the whole brain into distinct, smaller regions [25, 84]. Typically, brain parcellation is derived from T1-weighted MRI. There are several ways to parcellate the brain, including atlas-based, network-based, and data-driven parcellation [30, 72, 75]. Fig. 3 provides a detailed depic-



**Fig. 3:** Illustration of brain parcellation and morphology in sagittal view of MRI. Left plot showcases 120 regions of brain parcellation using the Desikan Atlas. Right plot represents the thickness and curvature of brain morphology.

tion of the brain parcellation. Domain experts, such as radiologists and neurologists, have suggested that predicting small anatomical parcels in 3D brain images may aid in the detection and localization of brain abnormalities, such as atrophy, which might not be visible or distinguishable when considering the

whole-brain level [5, 68]. Herein, we consider the brain anatomy prediction task as a multiclass segmentation problem, aiming to generate a brain parcellation map for a given image patch. More specifically, we have divided the brain into 120 non-overlapping regions based on a pre-defined atlas (i.e., Desikan atlas) and trained a model to predict the corresponding segmentation map for a given input patch, as illustrated in Fig. 2 (a). Given that the patches are relatively small and unlikely to encompass all 120 regions, our training was limited to regions that are present within the specific patch. We employed the Dice similarity coefficient as the loss function for each input patch and minimized it between the predicted and ground truth segmentation maps ($P$ and $\hat{P}$, respectively) for a given patch.

$$\mathcal{L}_{anatomy} = \text{Dice}(P, \hat{P}) \tag{1}$$

**Brain Morphology Prediction.** Brain morphology can be assessed with structural measures of the brain, such as volume, cortical thickness, and cortical curvature [40, 81]. In neuroimage analysis, studies focused on brain morphology because it provides valuable insights for estimating age, behavioral measurements (that is, memory performance and cognitive assessments), or disease diagnosis [1, 27, 51]. Fig. 3 provides a detailed depiction of brain morphology. In this study, we specifically examined two important morphological features: cortical thickness and curvature, due to their clinical relevance and associations with various diseases [87]. The precalculated morphological features, such as the average measurement within a specific brain region, are predicted using a regression framework for a given patch, as illustrated in Fig. 2 (b). During model training, we excluded measurements from regions that were not contained within the patch. The L1 loss is employed to minimize the difference between predicted and ground truth brain morphology measurements.

$$\mathcal{L}_{morpho} = \sum_{i \in \mathcal{S}} \|v_i^{mor} - \hat{v}_i^{mor}\|_1 \tag{2}$$

Here, $\mathcal{S}$ is a set of regions in a given patch, $v_i^{mor}$ denotes the ground truth brain morphology measurement in the $i$-th region, and $\hat{v}^{mor}$ is predicted brain morphology measurements using learned representation $z$.

**Radiomics Texture Prediction.** Radiomics is a medical research field focused on extracting numerous quantitative features from medical images, offering deeper insights than that perceivable by the human eye. Among these, radiomics texture features, such as the gray-level co-occurrence matrix (GLCM) and gray-level size zone matrix (GLSZM) features, assess voxel intensity relationships in an image, providing richer perspective than mere shape and size [2, 59]. When these features are applied to the task of brain image segmentation tasks, specifically for white matter (WM), gray matter (GM), and cerebrospinal fluid (CSF), they can reveal insights into tissue microstructures and their pathologic changes. For instance, radiomics texture variations can indicate changes in the WM due to aging or injury, and those variations in CSF hint at shifts in brain ventricle size. We extract these radiomics features from WM, GM, and CSF and train our

encoder, as shown in Fig. 2 (c), using the L1 loss to ensure alignment between the extracted and predicted features.

$$\mathcal{L}_{radiomics} = \sum_{i \in \mathcal{C}} \left\| v_i^{rad} - \hat{v}_i^{rad} \right\|_1 \tag{3}$$

Here, $\mathcal{C}$ represents a set of regions (GM, WM, and CSF), $v_i^{rad}$ denotes the ground truth radiomics feature in $i$-th region, and $\hat{v}^{rad}$ is the predicted radiomics feature.

### 3.2   Self-supervised Tasks

**Patch Location.** The concept of spatial context learning was first proposed by [22] and extended to a 3D context by [82]. The task of 3D patch location estimation was applied to leverage the 3D spatial context and learn the semantic representations of the data, as depicted in Fig. 2 (d). Specifically, the patch location prediction task randomly extracts N non-overlapping 3D patches from each input 3D image and then predicts patch locations by classifying N classes. The task is optimized by minimizing the cross-entropy loss between the ground-truth location and the predicted location (i.e., $y_i^{loc}$ and $\hat{y}_i^{loc}$) defined as follows:

$$\mathcal{L}_{loc} = -\sum_{i=1}^{N} y_i^{loc} \log(\hat{y}_i^{loc}) \tag{4}$$

Here, N is the number of patches extracted from the original image and $\hat{y}_i^{loc}$ is the predicted location using learned representation $z$. To prevent the model from taking advantage of trivial solutions by exploiting edge continuity and rapidly solving the task, random gaps are introduced between adjacent 3D patches.

**Image Rotation.** The 3D image rotation prediction task, as illustrated in Fig. 2 (e), was applied to provide semantic information for the model to learn [29,82]. In this task, we randomly rotate the 3D input patches by a degree chosen from a set of R possible degrees. We consider multiples of 90° along each axis of the 3D coordinate system, resulting in 10 rotation degrees that can be classified by the model. The task was formulated as a 10-way classification problem with the model minimizing the cross-entropy $\mathcal{L}_{rot}(y_r^{rot}, \hat{y}_r^{rot})$ for each rotated image.

**Masked Image Modeling.** MIM is a prominent SSL technique that has emerged as a potent pretraining method. It operates by masking certain parts of an image and then leveraging the unmasked parts to reconstruct the masked parts. This approach effectively handles local features, and its efficacy in 3D medical image analysis was validated by [16]. We pretrain our model using the SimMIM [90] method for 3D brain images. The associated loss is defined as follows:

$$\mathcal{L}_{MIM} = \frac{1}{\Omega(\mathrm{x}_M)} \left\| \mathrm{y}_M - \mathrm{x}_M \right\|_1 \tag{5}$$

Here, x and y are the input and predicted patches respectively; $M$ denotes the set of masked pixels; $\Omega$ is the number of elements.

### 3.3   Contrastive Learning

Recently, self-supervised learning techniques in the form of contrastive learning have emerged for deriving powerful representations by contrasting sample pairs [63,82,83]. In this study, contrastive coding was utilized to effectively learn the representations for a batch of augmented patches. Contrastive coding maximizes the mutual information for positive pairs (augmented patches from the same sample) and minimizes it for negative pairs (patches from different samples within a batch). More details are in the Supplementary Material. To determine our contrastive coding's loss, we added a linear layer to the Swin transformer encoder, mapping each augmented patch to a latent representation, $z$. This process is illustrated in Fig.2 (g). The distance between the encoded representations was measured using cosine similarity. Specifically, the 3D contrastive coding loss between patch pairs $z_i$ and $z_j$ is defined as:

$$\mathcal{L}_{contrast} = -\log\frac{\exp(f_{sim}(z_i, z_j)/\tau)}{\sum_{k=1,k\neq i}^{2N}\exp(f_{sim}(z_i, z_k)/\tau)} \tag{6}$$

Here, $\tau$ is a measure of the normalized temperature scale and $f_{sim}$ denotes the dot product between normalized embeddings.

### 3.4   Loss Function

The core idea of our framework is to learn the representation that captures both the 3D context and the characteristics of brain anatomy and morphology. The domain-aware, self-supervised, and total losses are calculated as follows:

$$\mathcal{L}_{domain} = \lambda_1\mathcal{L}_{anatomy} + \lambda_2\mathcal{L}_{morpho} + \lambda_3\mathcal{L}_{radiomics} \tag{7}$$

$$\mathcal{L}_{self} = \lambda_4\mathcal{L}_{rot} + \lambda_5\mathcal{L}_{loc} + \lambda_6\mathcal{L}_{MIM} \tag{8}$$

$$\mathcal{L}_{total} = \mathcal{L}_{domain} + \mathcal{L}_{self} + \lambda_7\mathcal{L}_{contrast} \tag{9}$$

The weights were empirically set to $\lambda_2=\lambda_3=\lambda_4=\lambda_5=\lambda_6=\lambda_7 = 1$, and $\lambda_1=0.2$.

## 4   Experiments and Results

### 4.1   Experimental Setup

**Datasets.**    In this study, we collected total of 13,687 samples of T1-weighted MRI data from multi-source large-scale databases. These included the Alzheimer's Disease Neuroimaging Initiative (ADNI) [41,61,67], Human Connectome Project (HCP) [85], Information eXtraction from Images (IXI) [6], Autism Brain Imaging Data Exchange (ABIDE) [20,21], Effects of TBI & PTSD on Alzheimer's Disease in Vietnam Vets (DOD ADNI) [88], International Consortium for Brain Mapping (ICBM) [60], and Anti-Amyloid Treatment in Asymptomatic Alzheimer's (A4) [19]. Further dataset details are provided in the Supplementary Materials.

**Methods Comparison.** We compared our model with the following existing 3D-based methods: (1) four 3D-CNN based methods, including 3D ResNet50,

3D ResNet10, 3D DenseNet121, and 3D DenseNet201. These models are widely used for AD classification [24, 45, 50, 69, 74, 93]. (2) Since these methods were not designed for 3D medical images, we also employed a 3D-CNN based model for medical images [69] which was proposed for accurate AD diagnosis. (3) Additionally, we employed pretrained CNN-based models, that is, the I3D proposed by [10] and MedicalNet presented by [13]. (4) We also consider transformer-based methods, 3D ViT and 3D Swin transformers, without pretraining for comparison. Further details are provided in the Supplementary Materials.

**Model evaluation and downstream tasks.**    We conducted three different downstream tasks (i.e., AD classification, PD classification, and age prediction) and compared their performances against those of competing models. For model evaluation, we employ five-fold cross-validation to report the mean of performance metrics. First, we compared the performance of our method with existing competing methods for AD using ADNI (total: 1,869, CN: 639, MCI: 886, and AD: 344), AIBL (total: 525, CN: 434, and AD: 91), and Open Access Series of Imaging Studies (OASIS) [55, 56, 73] (total: 817; CN: 676; AD: 141). Note that the ADNI datasets employed for the downstream tasks did not overlap with those considered for pretraining and were independently separated datasets. Additionally, we used independent AD datasets, such as AIBL and OASIS, which are not utilized in the training stage. Second, we evaluated our model by comparing it with existing competing methods for PD using the Parkinson's Progression Markers Initiative (PPMI) [57, 58] dataset (total: 663, CN: 161, and PD: 502) Third, a comparative study of chronological age prediction was performed using the ADNI datasets. Performance metrics varied for each task. For AD and PD classification, we assessed accuracy and area under the curve (AUC). For the age prediction task, the mean absolute error (MAE) and $R^2$ scores were utilized to evaluate performance. Additionally, we also extended well-known SSL frameworks, such as MoCo v2, BYOL, and DINO to 3D methods compared them with our model. For detailed settings, please refer to the Supplementary Materials.

**Implementation details.** We utilized a 3D Swin transformer as our backbone framework and trained it with the proposed pretext tasks described in the Supplementary Materials. We used the AdamW [54] optimizer with an initial learning rate of 0.0005, and the pretraining process was run for 300 epochs with a linear warmup and a cosine annealing learning rate scheduler. Further information on the training hyperparameters can be found in the Supplementary Materials. We implemented our model using PyTorch [65] and MONAI [8] and trained them on four A100 80GB GPUs.

### 4.2   Experimental Results

**Alzheimer's disease classification.**    We compared the performance of our method with existing competing methods for AD using ADNI, AIBL, and OASIS. Table 1 presents the comparison results of various AD classification tasks in terms of accuracy and AUC. Overall, our model exhibited superior performance compared with the other models across downstream tasks. Despite the challenges in capturing structural changes between AD and MCI, and between MCI and

**Table 1:** Performance evaluation of the pretrained Swin Transformer (ours) and comparison models in Alzheimer's disease classification. Four tasks were performed: binary classification between AD and CN, AD and MCI, MCI and CN, and a multi-class classification of AD vs. MCI vs. CN. Bold denotes the best performance in each column.

| Task | AD vs CN | | AD vs MCI | | MCI vs CN | | AD vs MCI vs CN | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Model | Acc | AUC | Acc | AUC | Acc | AUC | Acc | AUC |
| 3D ResNet50 [38] | 0.9063 | 0.9182 | 0.7250 | 0.6890 | 0.6625 | 0.6647 | 0.6383 | 0.7126 |
| 3D ResNet101 [38] | 0.8751 | 0.8771 | 0.7171 | 0.7032 | 0.6683 | 0.6559 | 0.6317 | 0.7086 |
| 3D DenseNet121 [39] | 0.9187 | 0.9191 | 0.7364 | 0.7368 | 0.6850 | 0.7007 | 0.6518 | 0.7329 |
| 3D DenseNet201 [39] | 0.9201 | 0.9234 | 0.7385 | 0.7248 | 0.6857 | 0.6986 | 0.6483 | 0.7252 |
| 3D ViT [23] | 0.8125 | 0.8220 | 0.6801 | 0.6638 | 0.6011 | 0.5956 | 0.5694 | 0.5975 |
| I3D [10] | 0.9135 | 0.9056 | 0.7362 | 0.7295 | 0.6929 | 0.6654 | 0.6409 | 0.7202 |
| MedicalNet [13] | 0.9292 | 0.9309 | 0.7338 | 0.7337 | 0.6824 | 0.7042 | 0.6492 | 0.7291 |
| Qiu et al. [69] | 0.9286 | 0.9438 | 0.7472 | 0.7442 | 0.6902 | 0.7063 | 0.6562 | 0.7358 |
| 3D Swin Tr (scratch) [53] | 0.9227 | 0.9204 | 0.7457 | 0.7496 | 0.6832 | 0.7020 | 0.6551 | 0.7342 |
| 3D Swin Tr (ours) | **0.9462** | **0.9623** | **0.7721** | **0.7796** | **0.7037** | **0.7275** | **0.6761** | **0.7521** |

CN, our model using the proposed method demonstrated successful classification. These results support the idea that various structural changes in the brain (such as atrophy of the cerebral cortex, enlargement of the ventricular areas, and shrinkage of the hippocampal volume) are progressing during AD [1, 51, 69] and these changes are visible on MRI to distinguish among the three AD classes. Additionally, we validated the generalizability of our model by comparing the AD classification performance on independent datasets, including the AIBL and OASIS datasets. For AD/CN classification, our model outperformed all competing methods on all datasets. Our model demonstrated the best performance on ADNI, AIBL, and OASIS, respectively, as presented in Tables 1 and 2. Moreover, compared to the 3D Swin Transformer without pretraining, our model showcased a significant improvement in prediction performance on ADNI, AIBL, and OASIS. These results suggest that our model effectively captures structural changes in the brain and consistently delivers high performance. The natural progression of AD starts from the CN, then to the MCI, and finally to AD. Therefore, tasks that distinguish between AD and CN are relatively easy, because the two classes occupy the two extreme ends of the spectrum. Thus, the tasks differentiating AD and MCI, as well as MCI and CN, are relatively difficult, owing to their relative proximity in the spectrum. Our results empirically confirmed these challenges.

**Parkinson's diseases classification.** We compared the performance of our method with existing competing methods for PD using the PPMI dataset. Table 2 presents quantitative comparisons of the PD classification tasks in terms of accuracy and AUC. The performance of PD classification is relatively low compared to that of AD, owing to subtle structural differences between the brains of patients with PD and healthy individuals [4, 78]. For instance, minor volume reductions in the substantia nigra or other relevant brain areas have been reported. However, these variations are potentially influenced by individual differences or other factors, making it difficult to diagnose PD conclusively [7, 28, 44]. Overall, despite the inherent challenge of distinguishing between PD and CN using only structural MRI, our proposed model demonstrated a significant performance improvement by detecting subtle structural changes.

**Table 2:** Comparison of model performance across various datasets and downstream tasks. AD classification refers to the binary classification between AD and CN, while PD classification denotes the binary classification between PD and CN. Bold denotes the best performance in each column.

| Task | AD classification | | | | PD classification | | Age prediction | |
|---|---|---|---|---|---|---|---|---|
| Dataset | AIBL | | OASIS | | PPMI | | ADNI | |
| Model | Acc | AUC | Acc | AUC | Acc | AUC | MAE | $R^2$ |
| 3D ResNet50 [38] | 0.8872 | 0.8728 | 0.8536 | 0.8273 | 0.7028 | 0.6128 | 4.6429 | 0.7368 |
| 3D ResNet101 [38] | 0.8631 | 0.8652 | 0.8624 | 0.8186 | 0.7087 | 0.6024 | 4.8207 | 0.7135 |
| 3D DenseNet121 [39] | 0.9164 | 0.9287 | 0.8595 | 0.8736 | 0.7356 | 0.6527 | 4.4230 | 0.7487 |
| 3D DenseNet201 [39] | 0.9267 | 0.9317 | 0.8551 | 0.8663 | 0.7294 | 0.6493 | 4.5148 | 0.7378 |
| 3D ViT [23] | 0.8768 | 0.8416 | 0.7837 | 0.7716 | 0.6465 | 0.5520 | 5.6476 | 0.6472 |
| I3D [10] | 0.9181 | 0.9253 | 0.8483 | 0.8569 | 0.7062 | 0.6346 | 4.6544 | 0.7335 |
| MedicalNet [13] | 0.9251 | 0.9358 | 0.8679 | 0.8832 | 0.7211 | 0.6471 | 4.6932 | 0.7379 |
| Qiu et al. [69] | 0.9237 | 0.9339 | 0.8465 | 0.8682 | 0.7509 | 0.6708 | 4.3750 | 0.7574 |
| 3D Swin Tr (scratch) [53] | 0.9201 | 0.9214 | 0.8660 | 0.8722 | 0.7323 | 0.6589 | 4.4803 | 0.7670 |
| 3D Swin Tr (ours) | **0.9372** | **0.9531** | **0.8809** | **0.8915** | **0.7586** | **0.6782** | **3.9138** | **0.7886** |

**Age prediction.** We compared the performance of our method with existing competing approaches for the task of chronological age prediction tasks using the ADNI dataset. Table 2 presents the comparison results of the MAE and $R^2$ score. For chronological age prediction, our model performed the best. These results suggest that the general capability of our method in effectively discerning brain's structural nuances and age-related variations, such as reductions in overall brain volume and regional shrinkage [27], which manifests as volume atrophy of the frontal and temporal lobes starting in post middle-age and a notable enlargement in the central ventricles [71, 79].

**Comparision with other SSL frameworks.** We also compared our model with well-known SSL frameworks, that is, MoCo v2, BYOL, and DINO by extending them with 3D methods. Our model trained with pretext multi-tasks exhibited the highest performance, while the closest baseline displayed a relatively lower performance, as illustrated in Table 3. These results suggest that the competing methods are specialized for 2D natural images, where it is relatively straightforward and easy to distinguish between instances and learn the features. However, 3D medical images have complex structures with similar morphologies, which makes it difficult to distinguish between instances. The model pretrained with SimMIM alone displayed the second-highest performance. Previous research has demonstrated that MIM can significantly enhance 3D medical image analysis [90]. We believe that MIM is a potent pretraining strategy; thus, incorporating MIM has led to even more performance improvements in our approach.

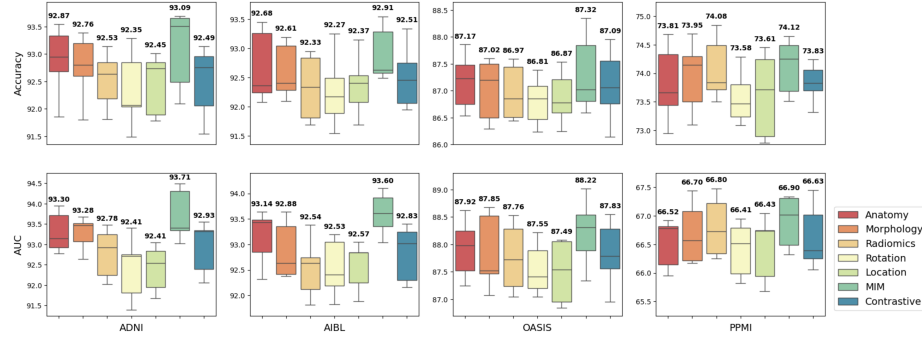### 4.3 Effectiveness of Each Self-Supervised Task

To evaluate the impact of each pretext task separately, we pretrain models using a single task. We tested them on various downstream tasks, as described above. The performance was assessed using a five-fold cross-validation, as illustrated in Fig. 4. We conduct binary classification between AD and CN for the ADNI, AIBL, and OASIS datasets. Overall, the performance of the pretrained

**Table 3:** Comparison of AD classification performance on ADNI dataset with other self-supervised methods.

| Model | Method | ACC. | AUC. |
|-------|--------|------|------|
| Swin Transformer | Scratch | 0.9227 | 0.9204 |
| | MoCo v2 [15] | 0.9246 | 0.9256 |
| | BYOL [33] | 0.9215 | 0.9218 |
| | DINO [9] | 0.9267 | 0.9282 |
| | SimMIM [90] | 0.9309 | 0.9371 |
| | Ours | **0.9462** | **0.9623** |

**Table 4:** Ablation study to evaluate the AD classification performance of various task combinations on ADNI dataset.

| Tasks | | | Metric | |
|-------|-------|-------|--------|------|
| $\mathcal{L}_{domain}$ | $\mathcal{L}_{self}$ | $\mathcal{L}_{contrast}$ | ACC. | AUC. |
| ✓ | | | 0.9348 | 0.9468 |
| | ✓ | | 0.9315 | 0.9424 |
| | | ✓ | 0.9249 | 0.9293 |
| | ✓ | ✓ | 0.9362 | 0.9480 |
| ✓ | | ✓ | 0.9397 | 0.9508 |
| ✓ | ✓ | | 0.9431 | 0.9588 |
| ✓ | ✓ | ✓ | **0.9462** | **0.9623** |



**Fig. 4:** The comparison of downstream tasks performance with varying pretext tasks for pretraining. The average accuracy (top) and AUC (bottom) for five-fold cross-validation are reported in each box plot.

model with the MIM task alone was the highest for all downstream tasks. From these results, we showcase that MIM effectively captures the structural context of images during pretraining. In addition, the novel tasks we proposed effectively learned structural information compared with MIM. Although most performance differences resembled those observed in AD classification, it is noteworthy that the radiomics texture prediction task excelled in PD classification. This suggests that subtle texture variations within brain regions can be crucial for PD classification, given the less-pronounced structural changes associated with PD [49].
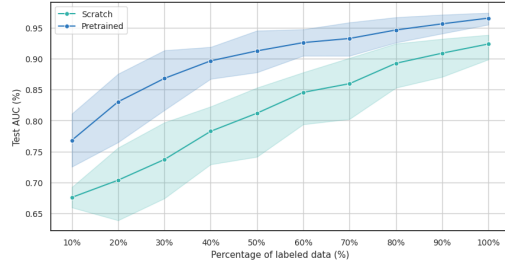
### 4.4 Effectiveness of Multi Self-Supervised Tasks

To evaluate the impact of combining different pretext tasks, we conducted experiments using various multi-task setups, as listed in Table 4. Overall, our approach of integrating self-supervised tasks with contrastive learning exhibited a significantly improved performance over models pre-trained on a single task. Specifically, the model pretrained with domain-aware tasks alone showcased a significant performance improvement compared with the best single-task pretrained model. The self-supervised tasks alone, such as rotation prediction, location prediction, and MIM, demonstrated a slight improvement over using MIM alone. Through this, we believe that by conducting various tasks in a multitasking manner for pretraining, it is possible to learn structural information in a more diverse manner, leading to performance improvement. It also highlights the

importance of effectively capturing both semantic and local information when learning about brain structural features. Additionally, combining our domain-aware tasks with other tasks confirmed that we could effectively grasp the overall context and structural information of the brain, which aids in AD classification.

### 4.5 Reducing the amount of manual labeling

The data size is important for downstream tasks. We compared the pretrained Swin transformer model and the Swin transformer without pretraining for a downstream task (AD classification). The comparison was conducted by adjusting the size of the fine-tuning data, ranging from 10% to 100% of the labeled dataset. As presented in Fig. 5, the pretrained model outperformed the scratch model across all data sizes



**Fig. 5:** The AUC graphs of the scratch model and pretrained model of the Swin Transformer according to the percentage of labeled data for the AD classification.

in terms of the average AUC. Specifically, the pretrained model displays a significantly improved performance even when 10% of the labeled data is adopted These results indicate that our model learns the data more effectively at a faster rate than the scratch model.

## 5 Visualization Explanations for AD Classification

AD is a complex debilitating neurodegenerative disorder. It is characterized by the accumulation of beta-amyloid and tau proteins in the brain, leading to neuronal injury, synaptic dysfunction, and cognitive impairment. The pathological hallmarks of AD are widespread and manifest as cerebral cortex atrophy, ventricular enlargement, and hippocampal volume loss. The natural course of AD involves gradual progression from CN to MCI, and finally to AD. MCI is a transitional state between the CN and AD, where brain alterations that occur in MCI are heterogeneous and can range from mild to severe, affecting different regions and functions of the brain. Therefore, distinguishing between AD and CN is relatively straightforward, whereas differentiating between AD and MCI, and between MCI and CN, is more challenging due to their close proximity to the disease spectrum. We compared the activated areas of our pretrained model and a scratch-learned model on the AD classification task using GradCAM++ and M3d-CAM [11,32]. From these results, we visually interpret imaging patterns as follows: (1) Our findings reveal that the pretrained model successfully detected both the hippocampus and corpus callosum in early MCI stages, as depicted in the bottom left of Fig. 6. However, the scratch model detected only the hippocampus, as presented in the top left of Fig. 6. (2) Our approach successfully

Fig. 6: **Visualization of the network's attention map.** The dashed line indicates various regions of interest. Top: scratch, bottom: pretrained.

identified the precuneus and prefrontal cortex, which are important for cognitive function, and are strongly associated with AD in the transition from the MCI to AD stages [64] as presented in the bottom middle of Fig. 6. (3) In the AD stage, our pretrained method detects not only the ventricular region but also the right temporal lobe reduction, which is a well-known AD biomarker [46] as indicated in the bottom right of Fig. 6. In contrast, the scratch model detected only the ventricular region, as illustrated in the top right of Fig. 6. These findings suggest that our pretrained model is capable of capturing the most common AD progression patterns and is more interpretable than the scratch model.

## 6    Conclusion

In this paper, we propose a pretraining method that integrates a novel domain recognition task with self-supervised task adapted to brain MRI data. The model was pre-trained using a substantial dataset of 13,687 brain MRI samples obtained from several large databases. We evaluated our pretrained method on three downstream tasks: AD classification, PD classification, and age prediction. The experimental results show the effectiveness of the multitasking approach in learning structural properties of the brain. The study also highlights that pre-training models with tasks specifically designed for structural MRI images of the brain can be used as a powerful pretraining tool to capture structural changes.

# References

1. Apostolova, L.G., Green, A.E., Babakchanian, S., Hwang, K.S., Chou, Y.Y., Toga, A.W., Thompson, P.M.: Hippocampal atrophy and ventricular enlargement in normal aging, mild cognitive impairment and alzheimer's disease. Alzheimer disease and associated disorders **26**(1), 17 (2012) 6, 10
2. Ardakani, A.A., Bureau, N.J., Ciaccio, E.J., Acharya, U.R.: Interpretation of radiomics features–a pictorial review. Computer Methods and Programs in Biomedicine **215**, 106609 (2022) 6
3. Azizi, S., Mustafa, B., Ryan, F., Beaver, Z., Freyberg, J., Deaton, J., Loh, A., Karthikesalingam, A., Kornblith, S., Chen, T., et al.: Big self-supervised models advance medical image classification. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 3478–3488 (2021) 4
4. Betrouni, N., Moreau, C., Rolland, A.S., Carrière, N., Chupin, M., Kuchcinski, G., Lopes, R., Viard, R., Defebvre, L., Devos, D.: Texture-based markers from structural imaging correlate with motor handicap in parkinson's disease. Scientific Reports **11**(1), 2724 (2021) 10
5. Betzel, R.F., Bassett, D.S.: Multi-scale brain networks. Neuroimage **160**, 73–83 (2017) 6
6. Biomedical Image Analysis Group Imperial College London Centre for the Developing Brain King's College London: Information eXtraction From Images. https://brain-development.org/ixi-dataset/ (2018), accessed: December 15, 2022 8, 22
7. Burton, E.J., McKeith, I.G., Burn, D.J., Williams, E.D., O'Brien, J.T.: Cerebral atrophy in parkinson's disease with and without dementia: a comparison with alzheimer's disease, dementia with lewy bodies and controls. Brain **127**(4), 791–800 (2004) 10
8. Cardoso, M.J., Li, W., Brown, R., Ma, N., Kerfoot, E., Wang, Y., Murrey, B., Myronenko, A., Zhao, C., Yang, D., et al.: Monai: An open-source framework for deep learning in healthcare. arXiv preprint arXiv:2211.02701 (2022) 9
9. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 9650–9660 (2021) 3, 5, 12, 25, 26
10. Carreira, J., Zisserman, A.: Quo vadis, action recognition? a new model and the kinetics dataset. In: proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 6299–6308 (2017) 4, 9, 10, 11
11. Chattopadhay, A., Sarkar, A., Howlader, P., Balasubramanian, V.N.: Gradcam++: Generalized gradient-based visual explanations for deep convolutional networks. In: 2018 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 839–847. IEEE (2018) 13
12. Chen, L., Bentley, P., Mori, K., Misawa, K., Fujiwara, M., Rueckert, D.: Self-supervised learning for medical image analysis using image context restoration. Medical image analysis **58**, 101539 (2019) 4
13. Chen, S., Ma, K., Zheng, Y.: Med3d: Transfer learning for 3d medical image analysis. arXiv preprint arXiv:1904.00625 (2019) 9, 10, 11
14. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: III, H.D., Singh, A. (eds.) Proceedings of the 37th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 119, pp. 1597–1607. PMLR (13–18 Jul 2020), https://proceedings.mlr.press/v119/chen20j.html 3

15. Chen, X., Fan, H., Girshick, R., He, K.: Improved baselines with momentum contrastive learning. arXiv preprint arXiv:2003.04297 (2020) 3, 12, 26
16. Chen, Z., Agarwal, D., Aggarwal, K., Safta, W., Balan, M.M., Brown, K.: Masked image modeling advances 3d medical image analysis. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 1970–1980 (2023) 7
17. Cruces, R.R., Royer, J., Herholz, P., Larivière, S., Vos de Wael, R., Paquola, C., Benkarim, O., yong Park, B., Degré-Pelletier, J., Nelson, M.C., DeKraker, J., Leppert, I.R., Tardif, C., Poline, J.B., Concha, L., Bernhardt, B.C.: Micapipe: A pipeline for multimodal neuroimaging and connectome analysis. NeuroImage **263**, 119612 (2022). https://doi.org/https://doi.org/10.1016/j.neuroimage.2022.119612, https://www.sciencedirect.com/science/article/pii/S1053811922007273 22, 24
18. Desikan, R.S., Ségonne, F., Fischl, B., Quinn, B.T., Dickerson, B.C., Blacker, D., Buckner, R.L., Dale, A.M., Maguire, R.P., Hyman, B.T., et al.: An automated labeling system for subdividing the human cerebral cortex on mri scans into gyral based regions of interest. Neuroimage **31**(3), 968–980 (2006) 24
19. Deters, K.D., Napolioni, V., Sperling, R.A., Greicius, M.D., Mayeux, R., Hohman, T., Mormino, E.C.: Amyloid pet imaging in self-identified non-hispanic black participants of the anti-amyloid in asymptomatic alzheimer's disease (a4) study. Neurology **96**(11), e1491–e1500 (2021) 8, 23
20. Di Martino, A., O'connor, D., Chen, B., Alaerts, K., Anderson, J.S., Assaf, M., Balsters, J.H., Baxter, L., Beggiato, A., Bernaerts, S., et al.: Enhancing studies of the connectome in autism using the autism brain imaging data exchange ii. Scientific data **4**(1), 1–15 (2017) 8, 23
21. Di Martino, A., Yan, C.G., Li, Q., Denio, E., Castellanos, F.X., Alaerts, K., Anderson, J.S., Assaf, M., Bookheimer, S.Y., Dapretto, M., et al.: The autism brain imaging data exchange: towards a large-scale evaluation of the intrinsic brain architecture in autism. Molecular psychiatry **19**(6), 659–667 (2014) 8, 23
22. Doersch, C., Gupta, A., Efros, A.A.: Unsupervised visual representation learning by context prediction. In: Proceedings of the IEEE international conference on computer vision. pp. 1422–1430 (2015) 3, 7
23. 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. In: International Conference on Learning Representations (2021), https://openreview.net/forum?id=YicbFdNTTy 2, 10, 11
24. Ebrahimi, A., Luo, S., Chiong, R.: Introducing transfer learning to 3d resnet-18 for alzheimer's disease detection on mri images. In: 2020 35th international conference on image and vision computing New Zealand (IVCNZ). pp. 1–6. IEEE (2020) 9
25. Fan, L., Li, H., Zhuo, J., Zhang, Y., Wang, J., Chen, L., Yang, Z., Chu, C., Xie, S., Laird, A.R., et al.: The human brainnetome atlas: a new brain atlas based on connectional architecture. Cerebral cortex **26**(8), 3508–3526 (2016) 5
26. Fischl, B.: Freesurfer. Neuroimage **62**(2), 774–781 (2012) 24
27. Fjell, A.M., Walhovd, K.B.: Structural brain changes in aging: courses, causes and cognitive consequences. Reviews in the Neurosciences **21**(3), 187–222 (2010) 6, 11
28. de la Fuente-Fernández, R.: Role of datscan and clinical diagnosis in parkinson disease. Neurology **78**(10), 696–701 (2012) 10
29. Gidaris, S., Singh, P., Komodakis, N.: Unsupervised Representation Learning by Predicting Image Rotations. In: ICLR 2018. Vancouver, Canada (Apr 2018), https://hal-enpc.archives-ouvertes.fr/hal-01864755 7

30. Glasser, M.F., Coalson, T.S., Robinson, E.C., Hacker, C.D., Harwell, J., Yacoub, E., Ugurbil, K., Andersson, J., Beckmann, C.F., Jenkinson, M., et al.: A multi-modal parcellation of human cerebral cortex. Nature **536**(7615), 171–178 (2016) 5

31. Gordillo, N., Montseny, E., Sobrevilla, P.: State of the art survey on mri brain tumor segmentation. Magnetic resonance imaging **31**(8), 1426–1438 (2013) 1

32. Gotkowski, K., Gonzalez, C., Bucher, A., Mukhopadhyay, A.: M3d-cam: A pytorch library to generate 3d attention maps for medical deep learning. In: Bildverarbeitung für die Medizin 2021: Proceedings, German Workshop on Medical Image Computing, Regensburg, March 7-9, 2021. pp. 217–222. Springer (2021) 13

33. Grill, J.B., Strub, F., Altché, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., et al.: Bootstrap your own latent-a new approach to self-supervised learning. Advances in neural information processing systems **33**, 21271–21284 (2020) 3, 12, 26

34. Hatamizadeh, A., Nath, V., Tang, Y., Yang, D., Roth, H.R., Xu, D.: Swin unetr: Swin transformers for semantic segmentation of brain tumors in mri images. In: Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: 7th International Workshop, BrainLes 2021, Held in Conjunction with MICCAI 2021, Virtual Event, September 27, 2021, Revised Selected Papers, Part I. pp. 272–284. Springer (2022) 26

35. Havaei, M., Davy, A., Warde-Farley, D., Biard, A., Courville, A., Bengio, Y., Pal, C., Jodoin, P.M., Larochelle, H.: Brain tumor segmentation with deep neural networks. Medical image analysis **35**, 18–31 (2017) 1

36. He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable vision learners. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16000–16009 (2022) 3

37. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 9729–9738 (2020) 3, 26

38. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016) 2, 10, 11

39. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4700–4708 (2017) 2, 10, 11

40. Im, K., Lee, J.M., Lyttelton, O., Kim, S.H., Evans, A.C., Kim, S.I.: Brain size and cortical structure in the adult human brain. Cerebral cortex **18**(9), 2181–2191 (2008) 6

41. Jack Jr, C.R., Bernstein, M.A., Fox, N.C., Thompson, P., Alexander, G., Harvey, D., Borowski, B., Britson, P.J., L. Whitwell, J., Ward, C., et al.: The alzheimer's disease neuroimaging initiative (adni): Mri methods. Journal of Magnetic Resonance Imaging: An Official Journal of the International Society for Magnetic Resonance in Medicine **27**(4), 685–691 (2008) 8, 22

42. Jenkinson, M., Beckmann, C.F., Behrens, T.E., Woolrich, M.W., Smith, S.M.: Fsl. Neuroimage **62**(2), 782–790 (2012) 24

43. Jun, E., Jeong, S., Heo, D.W., Suk, H.I.: Medical transformer: Universal brain encoder for 3d mri analysis. arXiv preprint arXiv:2104.13633 (2021) 5

44. Kalia, L.V., Lang, A.E.: Parkinson's disease. The Lancet **386**(9996), 896–912 (2015) 10

45. Karasawa, H., Liu, C.L., Ohwada, H.: Deep 3d convolutional neural network architectures for alzheimer's disease diagnosis. In: Intelligent Information and Database Systems: 10th Asian Conference, ACIIDS 2018, Dong Hoi City, Vietnam, March 19-21, 2018, Proceedings, Part I 10. pp. 287–296. Springer (2018) 9

46. Killiany, R.J., Moss, M.B., Albert, M.S., Sandor, T., Tieman, J., Jolesz, F.: Temporal lobe regions on magnetic resonance imaging identify patients with early alzheimer's disease. Archives of neurology **50**(9), 949–954 (1993) 14

47. Kim, J., Park, H.: Multi-modal cross attention network for predicting pathological complete response in breast cancer mri. In: 2023 5th International Conference on Control and Robotics (ICCR). pp. 250–254. IEEE (2023) 2

48. Kim, J., Park, H.: Radiomics-guided multimodal self-attention network for predicting pathological complete response in breast mri. In: 2024 IEEE International Symposium on Biomedical Imaging (ISBI). pp. 1–5 (2024). https://doi.org/10.1109/ISBI56570.2024.10635671 2

49. Korda, A.I., Andreou, C., Rogg, H.V., Avram, M., Ruef, A., Davatzikos, C., Koutsouleris, N., Borgwardt, S.: Identification of texture mri brain abnormalities on first-episode psychosis and clinical high-risk subjects using explainable artificial intelligence. Translational Psychiatry **12**(1), 481 (2022) 12

50. Korolev, S., Safiullin, A., Belyaev, M., Dodonova, Y.: Residual and plain convolutional neural networks for 3d brain mri classification. In: 2017 IEEE 14th international symposium on biomedical imaging (ISBI 2017). pp. 835–838. IEEE (2017) 9

51. Ledig, C., Schuh, A., Guerrero, R., Heckemann, R.A., Rueckert, D.: Structural brain imaging in alzheimer's disease and mild cognitive impairment: biomarker analysis and shared morphometry database. Scientific reports **8**(1), 11258 (2018) 6, 10

52. Litjens, G., Kooi, T., Bejnordi, B.E., Setio, A.A.A., Ciompi, F., Ghafoorian, M., Van Der Laak, J.A., Van Ginneken, B., Sánchez, C.I.: A survey on deep learning in medical image analysis. Medical image analysis **42**, 60–88 (2017) 2

53. 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) 2, 10, 11

54. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: International Conference on Learning Representations 9

55. 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) 9, 24

56. 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) 9, 24

57. Marek, K., Chowdhury, S., Siderowf, A., Lasch, S., Coffey, C.S., Caspell-Garcia, C., Simuni, T., Jennings, D., Tanner, C.M., Trojanowski, J.Q., et al.: The parkinson's progression markers initiative (ppmi)–establishing a pd biomarker cohort. Annals of clinical and translational neurology **5**(12), 1460–1477 (2018) 9, 24

58. 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) 9, 24

59. Mayerhoefer, M.E., Materka, A., Langs, G., Häggström, I., Szczypiński, P., Gibbs, P., Cook, G.: Introduction to radiomics. Journal of Nuclear Medicine **61**(4), 488–495 (2020) 6

60. Mazziotta, J., Toga, A., Evans, A., Fox, P., Lancaster, J., Zilles, K., Woods, R., Paus, T., Simpson, G., Pike, B., et al.: A probabilistic atlas and reference system for the human brain: International consortium for brain mapping (icbm). Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences **356**(1412), 1293–1322 (2001) 8, 23

61. Mueller, S.G., Weiner, M.W., Thal, L.J., Petersen, R.C., Jack, C.R., Jagust, W., Trojanowski, J.Q., Toga, A.W., Beckett, L.: Ways toward an early diagnosis in alzheimer's disease: the alzheimer's disease neuroimaging initiative (adni). Alzheimer's & Dementia **1**(1), 55–66 (2005) 8, 22

62. Noroozi, M., Favaro, P.: Unsupervised learning of visual representations by solving jigsaw puzzles. In: Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VI. pp. 69–84. Springer (2016) 3

63. Oord, A.v.d., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. arXiv preprint arXiv:1807.03748 (2018) 3, 8

64. Ossenkoppele, R., Cohn-Sheehy, B.I., La Joie, R., Vogel, J.W., Möller, C., Lehmann, M., van Berckel, B.N., Seeley, W.W., Pijnenburg, Y.A., Gorno-Tempini, M.L., et al.: Atrophy patterns in early clinical stages across distinct phenotypes of a lzheimer's disease. Human brain mapping **36**(11), 4421–4437 (2015) 14

65. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. Advances in neural information processing systems **32** (2019) 9

66. Pathak, D., Krahenbuhl, P., Donahue, J., Darrell, T., Efros, A.A.: Context encoders: Feature learning by inpainting. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2536–2544 (2016) 3

67. Petersen, R.C., Aisen, P.S., Beckett, L.A., Donohue, M.C., Gamst, A.C., Harvey, D.J., Jack, C.R., Jagust, W.J., Shaw, L.M., Toga, A.W., et al.: Alzheimer's disease neuroimaging initiative (adni): clinical characterization. Neurology **74**(3), 201–209 (2010) 8, 22

68. Plassard, A.J., McHugo, M., Heckers, S., Landman, B.A.: Multi-scale hippocampal parcellation improves atlas-based segmentation accuracy. In: Medical Imaging 2017: Image Processing. vol. 10133, pp. 666–672. SPIE (2017) 6

69. Qiu, S., Joshi, P.S., Miller, M.I., Xue, C., Zhou, X., Karjadi, C., Chang, G.H., Joshi, A.S., Dwyer, B., Zhu, S., et al.: Development and validation of an interpretable deep learning framework for alzheimer's disease classification. Brain **143**(6), 1920–1933 (2020) 1, 9, 10, 11

70. Raghu, M., Unterthiner, T., Kornblith, S., Zhang, C., Dosovitskiy, A.: Do vision transformers see like convolutional neural networks? Advances in Neural Information Processing Systems **34**, 12116–12128 (2021) 2

71. Raz, N., Rodrigue, K.M.: Differential aging of the brain: patterns, cognitive correlates and modifiers. Neuroscience & Biobehavioral Reviews **30**(6), 730–748 (2006) 11

72. Rolls, E.T., Huang, C.C., Lin, C.P., Feng, J., Joliot, M.: Automated anatomical labelling atlas 3. Neuroimage **206**, 116189 (2020) 5

73. Rowe, C.C., Ellis, K.A., Rimajova, M., Bourgeat, P., Pike, K.E., Jones, G., Fripp, J., Tochon-Danguy, H., Morandeau, L., O'Keefe, G., et al.: Amyloid imaging re-

sults from the australian imaging, biomarkers and lifestyle (aibl) study of aging. Neurobiology of aging **31**(8), 1275–1283 (2010) 9, 23

74. Ruiz, J., Mahmud, M., Modasshir, M., Shamim Kaiser, M., Alzheimer's Disease Neuroimaging Initiative, f.t.: 3d densenet ensemble in 4-way classification of alzheimer's disease. In: Brain Informatics: 13th International Conference, BI 2020, Padua, Italy, September 19, 2020, Proceedings 13. pp. 85–96. Springer (2020) 9

75. Schaefer, A., Kong, R., Gordon, E.M., Laumann, T.O., Zuo, X.N., Holmes, A.J., Eickhoff, S.B., Yeo, B.T.: Local-global parcellation of the human cerebral cortex from intrinsic functional connectivity mri. Cerebral cortex **28**(9), 3095–3114 (2018) 5

76. Selva, J., Johansen, A.S., Escalera, S., Nasrollahi, K., Moeslund, T.B., Clapés, A.: Video transformers: A survey. arXiv preprint arXiv:2201.05991 (2022) 2

77. Shamshad, F., Khan, S., Zamir, S.W., Khan, M.H., Hayat, M., Khan, F.S., Fu, H.: Transformers in medical imaging: A survey. arXiv preprint arXiv:2201.09873 (2022) 2

78. Sikiö, M., Holli-Helenius, K.K., Harrison, L.C., Ryymin, P., Ruottinen, H., Saunamäki, T., Eskola, H.J., Elovaara, I., Dastidar, P.: Mr image texture in parkinson's disease: A longitudinal study. Acta Radiologica **56**(1), 97–104 (2015) 10

79. Sowell, E.R., Peterson, B.S., Thompson, P.M., Welcome, S.E., Henkenius, A.L., Toga, A.W.: Mapping cortical change across the human life span. Nature neuroscience **6**(3), 309–315 (2003) 11

80. Spitzer, H., Kiwitz, K., Amunts, K., Harmeling, S., Dickscheid, T.: Improving cytoarchitectonic segmentation of human brain areas with self-supervised siamese networks. In: Medical Image Computing and Computer Assisted Intervention–MICCAI 2018: 21st International Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part III 11. pp. 663–671. Springer (2018) 4

81. Sporns, O., Tononi, G., Kötter, R.: The human connectome: a structural description of the human brain. PLoS computational biology **1**(4), e42 (2005) 6

82. Taleb, A., Loetzsch, W., Danz, N., Severin, J., Gaertner, T., Bergner, B., Lippert, C.: 3d self-supervised methods for medical imaging. Advances in neural information processing systems **33**, 18158–18172 (2020) 7, 8

83. Tang, Y., Yang, D., Li, W., Roth, H.R., Landman, B., Xu, D., Nath, V., Hatamizadeh, A.: Self-supervised pre-training of swin transformers for 3d medical image analysis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 20730–20740 (2022) 3, 4, 8

84. Tzourio-Mazoyer, N., Landeau, B., Papathanassiou, D., Crivello, F., Etard, O., Delcroix, N., Mazoyer, B., Joliot, M.: Automated anatomical labeling of activations in spm using a macroscopic anatomical parcellation of the mni mri single-subject brain. Neuroimage **15**(1), 273–289 (2002) 5

85. Van Essen, D.C., Ugurbil, K., Auerbach, E., Barch, D., Behrens, T.E., Bucholz, R., Chang, A., Chen, L., Corbetta, M., Curtiss, S.W., et al.: The human connectome project: a data acquisition perspective. Neuroimage **62**(4), 2222–2231 (2012) 8, 22

86. Van Griethuysen, J.J., Fedorov, A., Parmar, C., Hosny, A., Aucoin, N., Narayan, V., Beets-Tan, R.G., Fillion-Robin, J.C., Pieper, S., Aerts, H.J.: Computational radiomics system to decode the radiographic phenotype. Cancer research **77**(21), e104–e107 (2017) 27

87. Wang, J., Li, W., Miao, W., Dai, D., Hua, J., He, H.: Age estimation using cortical surface pattern combining thickness with curvatures. Medical & biological engineering & computing **52**, 331–341 (2014) 6

88. Weiner, M.W., Harvey, D., Hayes, J., Landau, S.M., Aisen, P.S., Petersen, R.C., Tosun, D., Veitch, D.P., Jack Jr, C.R., Decarli, C., et al.: Effects of traumatic brain injury and posttraumatic stress disorder on development of alzheimer's disease in vietnam veterans using the alzheimer's disease neuroimaging initiative: preliminary report. Alzheimer's & Dementia: Translational Research & Clinical Interventions **3**(2), 177–188 (2017) 8, 23
89. Wen, J., Thibeau-Sutre, E., Diaz-Melo, M., Samper-González, J., Routier, A., Bottani, S., Dormont, D., Durrleman, S., Burgos, N., Colliot, O., et al.: Convolutional neural networks for classification of alzheimer's disease: Overview and reproducible evaluation. Medical image analysis **63**, 101694 (2020) 1
90. Xie, Z., Zhang, Z., Cao, Y., Lin, Y., Bao, J., Yao, Z., Dai, Q., Hu, H.: Simmim: A simple framework for masked image modeling. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9653–9663 (2022) 4, 7, 11, 12
91. Zhai, X., Kolesnikov, A., Houlsby, N., Beyer, L.: Scaling vision transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12104–12113 (2022) 2
92. Zhang, D., Wang, Y., Zhou, L., Yuan, H., Shen, D., Initiative, A.D.N., et al.: Multimodal classification of alzheimer's disease and mild cognitive impairment. Neuroimage **55**(3), 856–867 (2011) 1
93. Zhang, J., Zheng, B., Gao, A., Feng, X., Liang, D., Long, X.: A 3d densely connected convolution neural network with connection-wise attention mechanism for alzheimer's disease classification. Magnetic Resonance Imaging **78**, 119–126 (2021) 9
94. Zhang, P., Wang, F., Zheng, Y.: Self supervised deep representation learning for fine-grained body part recognition. In: 2017 IEEE 14th international symposium on biomedical imaging (ISBI 2017). pp. 578–582. IEEE (2017) 4
95. Zhang, R., Isola, P., Efros, A.A.: Colorful image colorization. In: Leibe, B., Matas, J., Sebe, N., Welling, M. (eds.) Computer Vision – ECCV 2016. pp. 649–666. Springer International Publishing, Cham (2016) 3

# Supplementary Materials of Domain Aware Multi-Task Pretraining of 3D Swin Transformer for T1-weighted Brain MRI

This Supplementary Materials provide additional details not included in the main paper. In Sec. A, we provide details about the several datasets we employed. Sec. B includes information on training details and network hyperparameters. Sec. C details the pretraining task. Finally, Sec. D details the results of pretraining tasks.

## A   Datasets

To pretrain and evaluate our proposed methods, we utilized 13,687 samples from several large-scale T1 structural MRI databases, including ADNI, HCP, IXI, ABIDE, DOD ADNI, ICBM, and A4. We further employed four datasets for the model assessment: ADNI, AIBL, OASIS, and PPMI. These datasets are independent of the datasets utilized for pretraining and were considered solely for evaluation.

**Alzheimer's Disease Neuroimaging Initiative (ADNI)**   ADNI [41,61,67] is a research database dedicated to collecting multi-modal neuroimaging (MRI, fMRI, PET, and DTI) and non-imaging data (clinical outcome and genotyping data) related to AD. In our study, we obtained a total of 10,169 T1-weighted MR images from ADNI. These images encompass longitudinal data, different field strengths (1.5T and 3T), and scans from various manufacturers (Philips, Siemens, and GE). Of these images, we employed 8,300 images for pretraining phase and reserved 1,869 images for model evaluation. These data underwent preprocessing [17] and were employed to pretrain the model. For downstream tasks, we utilized a dataset of 1,869 samples, comprising CN: 639, MCI: 886, and AD: 344, for model evaluation.

**Human Connectome Project (HCP)**   HCP [85] is a large-scale initiative aimed at comprehensively mapping the neural connections within the human brain. In our study, a total of 1,104 MR images were acquired. The following parameters were considered to acquire MR scans: manufacturer = Siemens, field strength=3T, TR = 2400 ms, TE = 2.14 ms, Flip angle = 8 degrees, FOV = $224 \times 224 mm^2$, Matrix size = $256 \times 256$, and Voxel size = $0.7 \times 0.7 \times 0.7 mm^3$. These data were used to pretrain the model.

**Information eXtraction from Images (IXI)**   IXI [6] contains 581 MR images from healthy participants. These images include various MR scan types such as T1, T2, PD-weighted, MRA, and DWI. The T1-weighted images are

available in two field strengths (1.5T and 3T), and were scanned by different manufacturers (Philips, Siemens, and GE). We employed all these 581 images for pretraining.

**Autism Brain Imaging Data Exchange (ABIDE)** ABIDE [20, 21] contains 1099 MR images from Autism Spectrum Disorder (ASD) and control. These images include various MR scan types such as T1, resting state fMRI, and DWI. The T1-weighted images are available in two field strengths (1.5T and 3T), and were scanned by different manufacturers (Philips, Siemens, and GE). We used all these 1099 images for pretraining.

**Effects of TBI & PTSD on Alzheimer's Disease in Vietnam Vets (DOD ADNI)** DOD ADNI [88] focuses on exploring potential connections between traumatic brain injury (TBI), post-traumatic stress disorder (PTSD), MR scans, including longitudinal data. These T1-weighted scans were taken at two field strengths: 1.5T and 3T. The parameters for the 3T scanner were as follows: TR/TE = 2300/2.98ms, TI = 900ms, Flip angle = 9°, with a $1 \times 1 \times 1.2mm^3$ voxel size and $256 \times 256$ matrix over 170 slices. For the 1.5T scanner, they are: TR/TE = 2400/3.16ms, TI = 1000ms, Flip angle = 8°, with a $1.25 \times 1.25 \times 1.2mm^3$ voxel size and $256 \times 256$ matrix over 170 slices.

**International Consortium for Brain Mapping (ICBM)** ICBM [60] consists of 344 MRI images. These images were acquired axially in a 3D type using a body coil. The scans were taken with a SIEMENS TrioTim 3.0 Tesla machine. Key parameters include: Field Strength of 3.0 tesla, Flip Angle of 13.0°, and a GR/IR pulse sequence. The matrix dimensions are $220 \times 320 \times 208$ voxels with voxel sizes of $0.8 \times 0.8 \times 0.8mm^3$. Other notable parameters were TE = 2.8 ms, TI = 773 ms, and TR = 2200 ms, with a T1 weighting.

**Anti-Amyloid Treatment in Asymptomatic Alzheimer's (A4)** The A4 [19] provides a unique opportunity to compare MRI findings, such as Amyloid-related imaging abnormalities (ARIA), between cognitively impaired elderly individuals with high or low brain amyloid levels. This dataset includes sequences like T1, T2, GRE, FLAIR, and DWI, captured using a 3T MRI. The specifications for the 3T scanner are: voxel size of $1 \times 1 \times 1.2mm^3$ and a $256 \times 256$ matrix over 170 slices. For the pretraining of our model, we utilized 1791 T1-weighted images from this dataset.

**Australian Imaging, Biomarkers and Lifestyle (AIBL)** The AIBL [73] aims to provide researchers with new insights into the onset and progression of Alzheimer's disease. The dataset encompasses both AD and control groups. We utilized a total of 525 T1-weighted (T1w) images from this dataset, consisting of 434 CN and 91 AD samples for model validation. The T1 scanner parameters are set as follows: a matrix size of $240 \times 240 \times 160$, voxel size of $1 \times 1 \times 1.2mm^3$, TE=3.0 ms, TI=900.0 ms and TR=2300.0 ms.

**Open Access Series of Imaging Studies (OASIS)**  We utilized the OASIS [55, 56] dataset, specifically OASIS 3, which includes sequences such as T1w, T2w, FLAIR, ASL, SWI, time of flight, resting-state BOLD, and DTI. Out of these, we used 817 T1-weighted images (comprising 676 CN and 141 AD) for model validation.

**Parkinson's Progression Markers Initiative (PPMI)**  PPMI [57,58] dataset is a collection of a variety of medical data, including demographic and clinical, genetic, and neuroimaging data (i.e., MRI, PET, and SPECT). In our study, we obtained T1-weighted MRI data from a total of 663 images, which were acquired using the following parameters: field strength = 3T, repetition time (TR) = 2300 ms, echo time (TE) = 2.98 ms, and inversion time (TI) = 900 ms. Field of view (FOV) was $256 \times 256 mm^2$, matrix size was $256 \times 256$, and voxel size was $1 \times 1 \times 1.2 mm^3$.
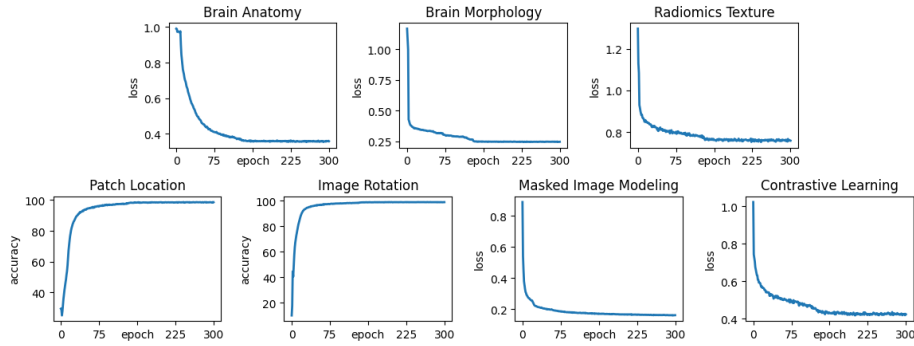
**Preprocessing**  The T1-weighted MR images used in our study were collected from various institutions, resulting in different matrix sizes, voxel spacings, and FOV. We employed the standard preprocessing steps [17], including skull stripping, bias field correction, and intensity normalization. Specifically, we skull-stripped MR using FSL-BET [42]. We resampled the voxels to $1.25 \times 1.25 \times 1.25 mm^3$. Then, we normalized the image intensities of all voxels using the zero-mean unit variance method. Brain anatomy was analyzed using the Desikan atlas, which involves dividing the whole brain into 120 regions and 17 subcortical regions, as computed by Freesurfer [18, 26]. Brain morphology measurements, cortical thickness and curvature, are also calculated using Freesurfer on Desikan atlas and 17 subcortical regions, resulting in 274 measurements.

## B   Implementation Details

**Model architecture**  We employ the Swin transformer as our backbone framework due to its efficiency on 3D data. Table A shows the model configuration. Specifically, the encoder architecture consists of four stages, each containing two transformer blocks except for the third stage, which consists of six transformer blocks, resulting in a total of $L = 24$ layers. Between stages, a patch merging layer is used to reduce the resolution by a factor of 2. In the first stage, the linear embedding layer and transformer blocks maintain the number of tokens at $\frac{H}{2} \times \frac{W}{2} \times \frac{D}{2}$. Additionally, a patch merging layer groups patches with a resolution of $2 \times 2 \times 2$ and concatenates them resulting in a 4C-dimensional feature embedding. A linear layer is then utilized to downsample the resolution by reducing the dimension to 2C. The procedure is repeated in stages 2, 3, and 4, with resolutions of $\frac{H}{4} \times \frac{W}{4} \times \frac{D}{4}$, $\frac{H}{8} \times \frac{W}{8} \times \frac{D}{8}$, and $\frac{H}{16} \times \frac{W}{16} \times \frac{D}{16}$, respectively. The patch size is set to $2 \times 2 \times 2$, with a feature dimension of $2 \times 2 \times 2 = 8$. The embedding space has a dimension of $C = 48$. The window size for multi-head self-attention is $7 \times 7 \times 7$.

**Table A:** Our swin Transformer configuration. FLOPs; floating point operations per second

| Patch Size | Window size | Feature size | Embedded Dimension |
|---|---|---|---|
| $2 \times 2 \times 2$ | $7 \times 7 \times 7$ | 48 | 768 |

| Number of Blocks | Number of Heads | Parameters | FLOPs |
|---|---|---|---|
| [2,2,18,2] | [3,6,12,24] | 57.16M | 82.38G |



**Fig. A:** The graphs represent various metrics during pretraining with multi-task learning. The y-axes of the graphs for patch location and image rotation show accuracy, which converges to nearly 100% during training. The y-axes of other tasks show the loss, which converges during 300 epochs.

**Settings of 3D ViT**    We set up the 3D patch embedding of size $16 \times 16 \times 16$ and a projection dimension of 2048. For 3D Swin transformer, we set the patch size to $2 \times 2 \times 2$, with feature dimensions of 8. The dimensions of the embedding space are $C = 48$. For multi-head self-attention, the window size was set to $7 \times 7 \times 7$.

**Data augmentation**    Two strategies were employed for data augmentation. First, we used multi-view (i.e., global and local views) augmentation inspired by DINO [9] for 3D input images. A global view was obtained by cropping and resizing the full image to remove the background to $128 \times 128 \times 128$, which included the entire brain. The local view, on the other hand, is a randomly cropped patch of size $56 \times 56 \times 56$ to focus on specific brain structures and that is further resized to $64 \times 64 \times 64$. Three local and one global views were considered for each sample. Second, we used a series of operations such as rotation and shifted intensity to augment the data. For contrast training, each view was augmented twice, yielding two enhanced views from the same sample. Furthermore, only one of these augmented views is masked, allowing for the simultaneous execution of contrastive learning and masked image modeling. All pretext tasks were applied

to both global and local views, except for the patch location prediction task, owing to the nature of the task.

**Hyperparameters**    We conducted training on four NVIDIA A100 GPUs, each with a batch size of 2. The pretraining phase involved an initial learning rate of 0.0001 for 300 epochs with a cosine annealing scheduler and linear warm-up. We utilized AdamW optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.999$.

**Settings of Other SSL Frameworks**    We tried to keep the original settings of SSL frameworks (i.e., MoCo v2 [15, 37], BYOL [33], and DINO [9]) in the comparative experiments as much as possible. However, our dataset consists of single channel 3D images and has a relatively small number of samples compared to previous studies. Therefore, we made some modifications to several hyperparameters. For the common augmentation method between ours and other SSL, we followed their implementations but replaced the color jitter with intensity scaling and shifting due to the single-channel nature of our medical images. The image size was cropped to $128 \times 128 \times 128$, and a pretrain batch size of 2 per GPU was used for 300 epochs.

**MoCov2**    We modified the default queue size to 12,288, because the total number of subjects in our dataset is 13,687.

**DINO**    We leveraged a global view of size $128 \times 128 \times 128$ and local views of size $56 \times 56 \times 56$. We used two global views and eight local views for the training process.

## C    Pretraining Task Details

**Brain Anatomy Prediction**    This task involved predicting the brain parcellation of a given patch. Only the regions belonging to the patch are considered during training, and the other regions are masked out during the loss calculation. For example, we are likely to consider only a few anatomically neighboring regions in a given patch. The segmentation task was performed by adding a simple CNN decoder to form a UNet-like structure, which is based on a previous study that employed a Swin Transformer as an encoder [34]. A total of 120 regions are predicted.

**Brain Morphology Prediction**    This task involved predicting the morphological features of each brain region. We predict the average thickness and curvature in each of the 137 brain regions. Similar to the brain anatomy prediction, only the regions within the patch are considered during training and other regions are masked out during the loss calculation. The morphology values were predicted using a morphology head composed of a simple multilayer perception (MLP) consisting of two FC layers.

**Radiomics Texture Prediction**    This task aimed to predict the radiomics texture features of the gray matter, white matter, and CSF regions. For each region, 20 GLCM features and four GLSZM features are extracted, resulting in 72 features (3 regions with 24 features each). These features were extracted using Pyradiomics v3.0.1 [86]. To execute this task, representation z from the swin transformer was passed through a two-layer perceptron for regression prediction.

**Patch Location**    For the patch location task, an eight-way classification was conducted to estimate the location of $2 \times 2 \times 2$ sub-patches within the 3D images. This task is performed only locally. For patch location, the representation was trained with a single FC layer to perform an 8-way classification.

**Image Rotation**    In our 3D rotation prediction task, we randomly rotate 3D input patches by a degree chosen from a set of 12 possible degrees (i.e., 0, 90, 180, 270 degrees along each axis), then train the model to predict rotation degree in a classification manner. Since the zero-degree rotation of the x, y, and z axes were the same, only 10 possible rotation degrees were available for our classification task. In our study, we added a single FC layer as the image rotation head for 10-way classification.

**Masked Image Modeling**    In global and local perspectives, 75% of the 3D volume within the patch was masked out. We employed a patch size of 16 and randomly generated the cut-out regions. A single-layer projection with pixel shuffle served as the MIM head. During pretraining, the L1 loss was calculated between the original and reconstructed patches.
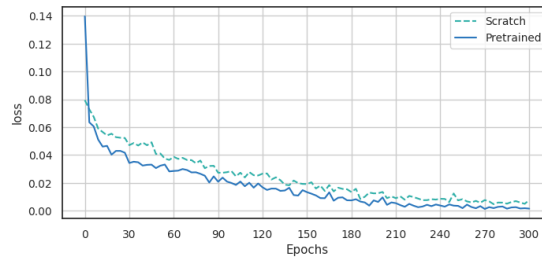
**Contrastive Learning**    To perform contrastive Learning to randomly augmente the patches to generate positive and negative pairs. Specifically, because we set the batch size to two, one positive pair and two negative pairs were available for $i$-th augmented patch. Then, we computed the latent representation z of each augmented patch using linear projection, where the dimension of the latent representation was 512. Finally, the contrastive coding loss is computed using eq.6. In our study, we applied a contrastive learning task to both global and local views to learn multiscale representations.

## D    Results of Pretraining Tasks

**Learning progress of each task** To demonstrate the effectiveness of our multi-tasking approach, we evaluated the performance of each task during the pretraining phase. Fig. A showcases the metrics for each task during the training phase. The y-axes of the graphs for patch location and image rotation represent accuracy, whereas masked image modeling and brain morphology use L1 loss, brain anatomy uses the Dice coefficient, and contrastive learning uses information noise

and contrastive estimation loss. Each task demonstrated that the learning metrics converged during training. The model shows pretraining on various aspects of the brain's structural features across seven tasks.
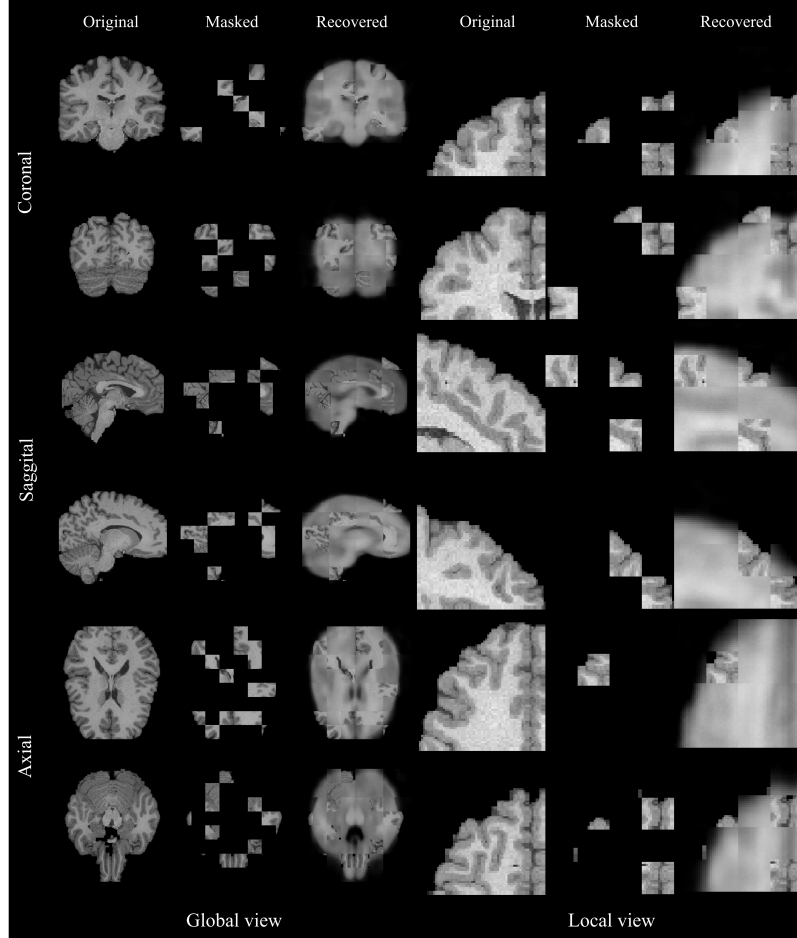
**Effectiveness of Pretraining** We compared the convergence speed of training with the scratch model and our pretrained model. Fig. B depicts the convergence graphs of the training losses for the two swin transformer models. Our pretrained model not only converges faster in the early epoch, but also has a lower loss than the scratch model at all epochs. The results demonstrate the effectiveness of our pretraining method using multi-task learning.
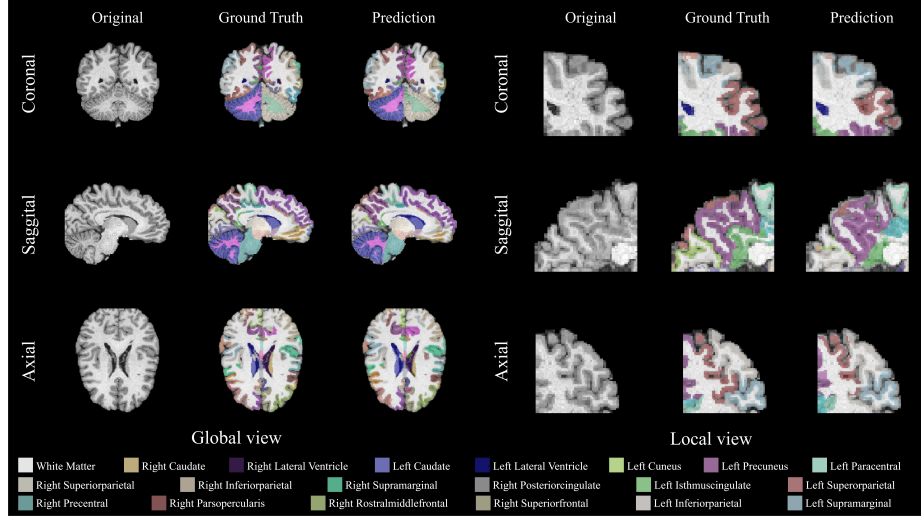


**Fig. B:** The train loss graphs of the scratch and pretrained Swin Transformer

**Results of MIM** Fig. C illustrates the reconstruction process for MIM. To pretrain the encoder, we attached a single projection layer to reconstruct the masked 3D volume. Despite performing reconstruction through a simple single projection layer, it is evident that the masked areas are effectively encoded, enabling the identification and restoration of the corresponding structure.
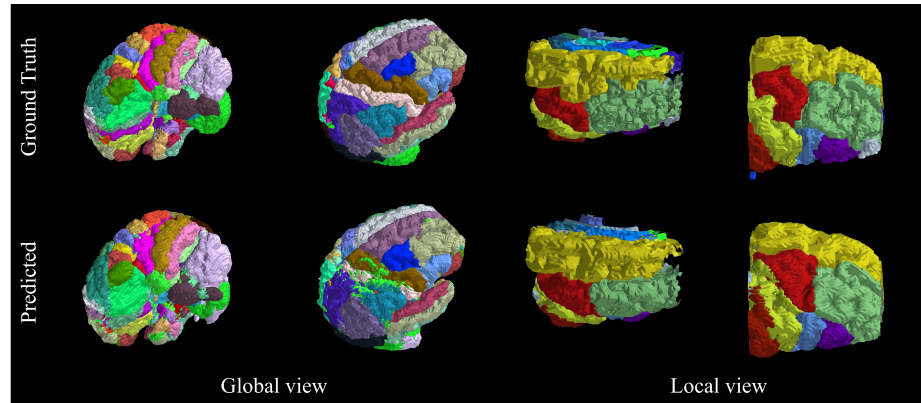
**Results of Brain Anatomy Prediction** To assess task performance, we visualized ground truth and prediction parcellation. Fig. D illustrates the process of predicting brain anatomy. Fig. E shows a 3D rendering comparing the ground truth with the predicted brain parcellation. Our objective was to train the encoder. Therefore, we utilized a lightweight CNN decoder and observed its ability to reasonably predict the locations of rough parcellations.

**Fig. C:** Illustration of the training process for the masked image modeling task. Original: source image. Masked: Image from the original with 75% masked out. Recovered: Image restored after passing the masked image through a single projection head. The model is trained using the L1 Loss between the original and recovered. Given that the input is a 3D volume, we present three planes: coronal, sagittal, and axial. Each plane displays two distinct views. Top: coronal, Middle: sagittal, Bottom: axial. Left: global view, Right: local view.

**Fig. D:** Illustration of the training process for the brain anatomy prediction task. Original: Source image. Ground Truth: Image overlayed with the ground truth brain parcellation on the original. Predicted: Image overlayed with the predicted brain parcellation on the original. It shows that the pretrained encoder with our multi-task effectively predicts brain parcellation. Left: global view, Right: local view



**Fig. E:** 3D rendering of both the ground truth and predicted brain parcellation. Two different viewing angles are presented. Left: global view, Right: local view