

BrainMVP: Multi-modal Vision Pre-training for Brain Image Analysis using Multi-parametric MRI

Shaohao Rui, Lingzhi Chen, Zhenyu Tang, Lilong Wang, Mianxin Liu, Shaoting Zhang, Xiaosong Wang

Abstract—Accurate diagnosis of brain abnormalities is greatly enhanced by the inclusion of complementary multi-parametric MRI imaging data. There is significant potential to develop a universal pre-training model that can be quickly adapted for image modalities and various clinical scenarios. However, current models often rely on uni-modal image data, neglecting the cross-modal correlations among different image modalities or struggling to scale up pre-training in the presence of missing modality data. In this paper, we propose BrainMVP, a multi-modal vision pre-training framework for brain image analysis using multi-parametric MRI scans. First, we collect 16,022 brain MRI scans (over 2.4 million images), encompassing eight MRI modalities sourced from a diverse range of centers and devices. Then, a novel pre-training paradigm is proposed for the multi-modal MRI data, addressing the issue of missing modalities and achieving multi-modal information fusion. Cross-modal reconstruction is explored to learn distinctive brain image embeddings and efficient modality fusion capabilities. A modality-wise data distillation module is proposed to extract the essence representation of each MR image modality for both the pre-training and downstream application purposes. Furthermore, we introduce a modality-aware contrastive learning module to enhance the cross-modality association within a study. Extensive experiments on downstream tasks demonstrate superior performance compared to state-of-the-art pre-training methods in the medical domain, with Dice Score improvement of 0.28%-14.47% across six segmentation benchmarks and a consistent accuracy improvement of 0.65%-18.07% in four individual classification tasks.

Index Terms—Multi-parametric MRI, Data distillation, Cross-modal reconstruction, Contrastive learning

I. INTRODUCTION

MULTI-PARAMETRIC MRI (mpMRI) images combine various imaging modalities to comprehensively depict the structural and pathological features of the brain [42]. This approach substantially enhances diagnostic accuracy and thoroughness [40]. For instance, in the subregional segmentation of brain tumors, different MRI modalities reveal distinct lesion characteristics. Areas with high signal intensity in T1Gd

images compared to T1 images and healthy white matter identify the enhancing tumor (ET), whereas the tumor core, including the ET and necrotic portion, appears with low signal intensity in T1Gd images relative to T1 images. The whole tumor encompasses the tumor core and the surrounding edematous or invaded tissue, characterized by the abnormally high signal intensity in T2-FLAIR images [4].

Current methods for lesion delineation and disease classification using mpMRI heavily rely on supervised models trained on specific datasets, which limits their applicability across different datasets and tasks. Developing a multi-parametric brain MRI foundation model capable of addressing cross-modal representation and improved performance is crucial for advancing medical applications.

However, challenges like the high cost of medical image annotation, the limited accessibility of large amounts of multi-modal data [45], and the lack of dedicated multi-modal pre-training paradigm with medical image domain expertise [3] complicate the construction of such a generalized mpMRI foundation model.

Obtaining a comprehensive set of modalities in mpMRI scans can be challenging due to the complexity associated with acquisition protocol and limitations in equipment capabilities. This often leads to mismatched modality data, e.g., missing modalities cross datasets, especially when the scale of the data amount increases dramatically. Furthermore, current approaches to dealing with missing modalities primarily focus on particular downstream tasks, e.g., BraTS, and have not undergone extensive investigation in large-scale cross-modal pre-training and downstream tasks [57, 14, 48, 25, 39].

Existing research predominantly focuses on monotone image modality and utilizes cross-modal prompt learning by aligning text-based information with medical images [11, 47, 54]. However, these approaches do not directly address the challenge of effectively fusing multi-modal image information for mpMRI images. While some studies have begun exploring this issue in cross-modal mpMRI scenarios, their methodologies are often restricted to a small number of modalities and small datasets in either pre-training and downstream tasks [52, 45, 31, 41], thus leading to limited model generalizability.

Maximizing the capability of pre-training models in downstream tasks holds paramount importance ([52, 44]) and is closely related to the designed pre-training tasks. Frequently, pre-trained proxy tasks lack direct correlations to downstream applications, e.g., masked image modeling, resulting in sub-optimal performance when blindly applied. By strategically

Shaohao Rui and Lingzhi Chen contributed equally to this work. Corresponding authors: Xiaosong Wang.

The authors are with Shanghai AI Laboratory, Shanghai 200030, China. Shaohao Rui and Zhenyu Tang are interns at Shanghai AI Laboratory. (email: {ruishaohao, chenlingzhi, tangzhenyu1, wanglilong, liumianxin, zhangshaoting, wangxiaosong}@pjlab.org.cn)

Shaohao Rui and Zhenyu Tang are with the School of Electronic Information and Electrical Engineering, Shanghai Jiao Tong University, Shanghai 200240, China

incorporating these links, we can guide and enhance the efficacy of the pre-training process, ensuring that learned representations align more closely with the requirements of multi-modal fusion into the unified diagnosis.

In this paper, we introduce BrainMVP, a novel multi-modal vision pre-training framework for the multi-parametric MRI images of the brain that demonstrates distinctive and generalizable cross-modal image representation. Initially, we gather a dataset of 16,022 publicly available brain mpMRI scans from various multi-center, multi-device sources. The dataset covers different types of brain imaging modalities, including diseased and healthy brains.

To address the issue of insufficient scalability (due to mismatched or missing modalities), we propose using single-modal MRI image inputs instead of fixed modality numbers in the pre-training stage. This allows for the inclusion of arbitrary numbers of modalities in the pre-training, significantly expanding the magnitude of available pre-training data. Importantly, we propose cross-modal image reconstruction via mask modeling. A key aspect of this design is the observation that different MRI modalities for the same subject often exhibit significant similarity in anatomy. By employing cross-modal reconstruction, we encourage the model to learn the disentanglement across modalities while mining the modality-invariant representations.

Toward a more generalizable pre-training model for downstream tasks, we extract condensed representations of different modality structures using modality-wise data distillation. Our approach is inspired by the technique of data distillation, which involves learning a small synthetic dataset. The performance achieved by the model training on this synthetic dataset can rival that achieved on the original large-scale datasets [49, 58, 55]. The learned synthesized dataset indeed encapsulates dense representations of the original dataset. In a similar idea, we optimize a set of learnable modality templates tailored for each individual modality. Intuitively, the distilled modality templates retain rich structural and statistical information about a specific modality while avoiding privacy leakage concerns associated with individual patients. Moreover, the learned distilled modality templates can serve as a linkage of data between pre-training and downstream tasks, i.e., as a form of information to carry and adapt between the data domains for downstream applications.

In summary, our contributions are three fold:

- To the best of our knowledge, BrainMVP is the first multi-modal vision pre-training paradigm that aligns the features across modalities, targeting distinctive modality-aware representations. We collect a dataset of 16,022 mpMRI scans (over 2.4 million images) to facilitate the pre-training, covering a wide range of MRI brain image sequences in both diseased and healthy populations.
- We design two novel proxy task settings for the multi-modal vision pre-training, i.e., cross-modal reconstruction and cross-modality contrastive learning. To improve the generalization for downstream tasks, we also introduce modality-wise data distillation to extract the template of each modality, benefiting both the pre-training and downstream tasks.

- We demonstrate the superior performance gain and the enhanced generalizability by utilizing our BrainMVP pre-trained models on ten public segmentation and classification benchmarks, compared to state-of-the-art methods.

II. RELATED WORK

A. Multi-modal Pre-training for Natural Image Analysis

In the pursuit of acquiring knowledge, humans typically engage with data from multiple modalities. These diverse data sources, obtained from various perspectives, complement one another, enabling a more comprehensive understanding and facilitating the completion of more advanced semantic tasks. Recently, research in visual-language pre-training has seen significant advancements, primarily aiming to enhance the performance of various downstream related to vision through the alignment of different modal data. CLIP [36] pioneered large-scale image-text feature alignment by employing contrastive learning to maximize the mutual information between matched image-text pairs while minimizing it for mismatched pairs. Subsequent improvements, as noted in works like [35, 15, 30], have demonstrated robust generalization and zero-shot reasoning capabilities achieved through cross-modal knowledge alignment.

Another type of multi-modal pre-training focuses on the fusion of different modal information to enhance cross-modal data understanding and address the limitations of feature richness in uni-modal data. For instance, [2] employs a gated cross-attention mechanism to integrate features extracted from a frozen vision encoder and language model. To reduce the cost of end-to-end visual-language pre-training on large-scale datasets, BLIP-2 [28] proposes leveraging off-the-shelf frozen pre-trained image and language models. It introduces a lightweight, learnable Q-Former module to bridge the gap between modalities, facilitating image-to-text transformation through a two-stage learning process. ALBEF [27] suggests aligning visual and text features before inputting them into a multi-modal Transformer network. Additionally, to enable efficient learning from noisy web data, ALBEF [27] introduces a momentum distillation method to aid model training.

While the aforementioned works focus on cross-domain interactions, our research centers on multi-modal integration within a single domain, specifically developing pre-training methods for the fusion of different modalities in MRI images.

B. Self-Supervised Learning for Medical Image Analysis

Medical image data annotation is notably expensive, making self-supervised learning (SSL) a promising avenue for the development of efficient annotation techniques. Given the typically limited datasets available for specific medical tasks, pre-training on large-scale unlabeled data to extract highly generalizable representations is emerging as a new paradigm. Existing SSL methods related to medical images can be roughly divided into cross-domain and in-domain fashion.

1) *Cross-domain SSL*: Typical multi-modal medical image self-supervised pre-training is achieved through the joint involvement of images and text. MGCA [46] leverages the

semantic correspondence between medical images and radiology reports across three distinct levels: pathology region level, instance level, and disease level to facilitate generalized medical visual representation learning. Additionally, a multi-modal approach [42] introduces a multi-modal puzzle task designed to enhance rich representation learning from various image modalities. By obfuscating image modalities at the data level and employing the Sinkhorn operator to frame the puzzle solution as a permutation matrix inference, this method efficiently addresses multi-modal jigsaw puzzles of varying complexity. Furthermore, [10] propose a self-supervised learning paradigm for medical images and texts, named the multi-modal masked self-coder. This method acquires cross-modal domain knowledge by reconstructing missing pixels and tokens in randomly masked images and texts.

2) *In-domain SSL*: Most current SSL methods specific to medical images are based on contrastive learning and masked image modeling (MIM) to extract useful information within images. For instance, [20] introduces a geometric visual similarity pre-training framework that leverages the high topological similarity of medical images. This approach incorporates a priori information about topological invariance into the similarity metric between images and employs a proposed z-matching head to learn the similarity of semantic features at different scales. PCRLv2 [56] addresses the issue of local information loss in medical images within the contrastive learning SSL paradigm by suggesting pixel recovery and feature alignment at various scales for diverse enhancement samples. Additionally, PCRLv2 [56] recommends implementing SSL without using skip connections to avoid shortcut solutions in pixel restoration. SwinMM [50] trains several proxy tasks involving masked multi-view observation, such as image reconstruction, rotation, contrastive learning, and a novel task that exploits the consistency of multiple views to extract hidden multi-view information in 3D medical data. During the fine-tuning stage, SwinMM [50] utilizes cross-attention blocks to aggregate multi-view information. Leveraging the high structural similarity of medical images, TransVW [16] conceptualizes subregions of an image as transferable visual words and learns generic visual representations by predicting and reconstructing their region categories. Specifically, TransVW [16] identifies similar samples of the current image through nearest-neighbor classification on encoded image features and uses the four regions of these similar samples for region categorization to ensure semantic consistency. TransVW [16] also applies perturbations to transform the visual word parts, and reconstructs the transformed image to learn structured semantic representations.

Our work integrates contrastive learning and MIM based SSL methods, and we have deviated from previous MIM based approaches by discarding masked portions or filling them with noise or zero values. Instead, we use informative images from other modalities for filling, which improves the learning efficiency and highlights the correlations between modalities.

C. Data distillation

Data distillation is first inspired by knowledge distillation, proposed to distill the dataset in order to construct a core

subset of the data, the model is trained on the core subset to achieve a performance comparable to that of the complete data [49]. In this way, model training and data storage costs can be significantly reduced [55]. Several studies have already applied data distillation in the field of continual learning to achieve data replay to attenuate the significant oblivion of distributional bias on the knowledge capability of models. Inspired by this, we propose to use data distillation to preserve the structural statistical information of different models from pre-trained data, and to construct upstream-downstream linkage through modality-wise data distillation.

Overall, we propose cross-modal reconstruction and modality-aware contrastive learning as the two main proxy tasks, as well as enhancing the network’s learning for discriminative information through contrastive learning. Different from previous methods, we propose a novel masking strategy to learn efficient modal fusion capability through cross-modal reconstruction and build upstream and downstream correlations through data distillation to better adapt the pre-trained models to the downstream tasks.

III. MULTI-MODAL VISION PRE-TRAINING

As shown in Fig. 1, BrainMVP consists of three proxy tasks: cross-modal reconstruction, modality-wise data distillation, and modality-aware contrastive learning. The proposed cross-modal reconstruction (via two ways of masking and mix) module aims to achieve the disentanglement across modalities while mining the modality-invariant representations. Modality-wise data distillation is designed to learn compressed structural information for each modality from pre-trained unlabeled data while allowing the model to extract modality-wise information and learn modality-agnostic features. Furthermore, the distilled modality templates are applied to downstream tasks to establish the association of data between the pre-training and downstream domains, which helps to improve the generalization performance of the foundation model. Finally, modality-aware contrastive learning is integrated to ensure the consistency of semantic features between different masked versions of the same sample, as well as to extract discriminative information between modalities.

A. Cross-Modal Reconstruction

Problem Setting. Given an unlabeled dataset $\mathcal{D} = \{X_{im} \in \mathbb{R}^{D \times H \times W} | m \in \{1, \dots, M_i\}, i \in \{1, \dots, N\}\}$. M_i denotes the number of modalities in i -th sample and N represents total number of samples. Masked image modeling (MIM) first masks (with constant values, zeros, or noise, etc, denoted as $\Phi(\cdot)$) a large portion of X_{im} to obtain a masked input $\Phi(X_{im})$, and then reconstructs the original image from it to learn efficient representations. Specifically, let the encoder and decoder of the model be $f_{enc}(\cdot)$ and $f_{dec}(\cdot)$ respectively, where the model $f(\cdot)$ is the composition of the encoder and decoder functions, i.e., $f(\cdot) = f_{dec} \circ f_{enc}$. MIM minimizes the following reconstruction loss:

$$L_{rec} = \frac{1}{N} \sum_{i=1}^N \frac{1}{M_j} \sum_{m=1}^{M_j} \|f_{dec}(f_{enc}(\Phi(X_{im}))) - X_{im}\|_2 \quad (1)$$

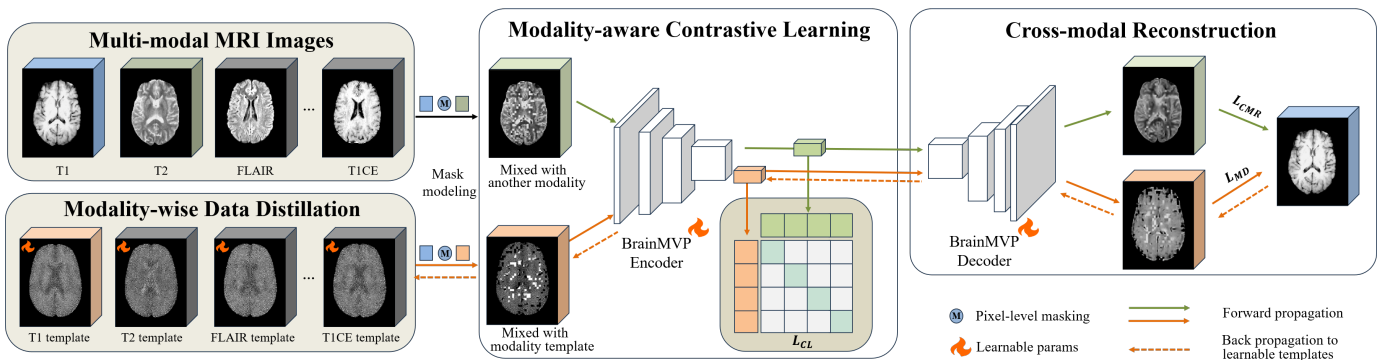


Fig. 1. Overview of the proposed BrainMVP, comprised of (a) cross-modal reconstruction module that aims at learning a mapping from images masked and mixed with another modality to the original; (b) modality-wise data distillation module that learns condensed modality templates via gradient backpropagation; and (c) modality-aware contrastive learning module for introducing case-level modality invariance to the learned features.

The core idea of our proposed reconstruction proxy tasks, which are elaborated in Sections III-A and III-B, is to obtain meaningful representations via exploiting different forms of $\Phi(\cdot)$ function.

Pixel-level cross-modal masking. Given a uni-modal input volume X_{im} sampled from a mpMRI case (with M_i modalities), cross-modal masking aims to mask a large region of X_{im} with another modality image X_{in} (also sampled from X_i , $n \neq m$). Specifically, we first randomly mask a region of size $r \times r \times r$ in X_{im} , where r denotes the size of each dimension of input 3D volumes. Then, we fill the masked region with a patch cropped with the same location and size on another modality of the sampled case. Finally, we repeat the above masking-filling operation until the proportion of masked pixels over the total input volume (X_{im}) pixels arrives p^* . Referring to [31], we empirically set $r = 8$ and $p^* = 0.8$ to learn useful representations. Details can be seen in Algorithm 1.

Algorithm 1 Pixel-level cross-modal masking.

```

Sample randomly  $X_{im}$  from  $X_i$ 
Sample randomly  $X_{in}(n \neq m)$  from  $X_i$ 
 $p_{total} \leftarrow H \times W \times D$ 
 $p_{mask} \leftarrow 0$ 
while  $p_{mask} < p_{total} \times p^*$  do
  Select randomly  $(x, y, z)$  in  $X_{im}$ 
  Mask an area of size  $r \times r \times r$  centered at  $(x, y, z)$ 
  Fill with corresponding data from  $X_{in}$ 
   $p_{mask} \leftarrow p_{mask} + r \times r \times r$ 
end while
return modified  $X_{im}$ 

```

Cross-modal reconstruction. Let our proposed cross-modal masking strategy be $\Phi_{modal}(\cdot)$. Given that the masking operation masks a large portion of the image, the resulting masked input volume $\Phi_{modal}(X_{im})$ will contain information predominantly from X_{in} . The extracted representation $f_{enc}(\Phi_{modal}(X_{im}))$ will thus encode a significant amount of semantic information from X_{in} . Since we do not introduce skip connections between the encoder and decoder, we only reconstruct X_{im} from the latent representation $f_{enc}(\Phi_{modal}(X_{im}))$, which is a challenging task for natural

images. However, due to the high structural similarity between different modalities in mpMRI images, with strong contrasts only in certain regions, the cross-modal reconstruction can encourage the model to learn cross-modal representations and explore the correlations between different modalities. Furthermore, $\Phi_{modal}(X_{im})$ still contains $(1-p^*)$ proportion of information about X_{im} , and reconstructing this part will retain some original modality information, which can help reduce the difficulty of pure cross-modal reconstruction and extract the semantic information of the X_{im} image itself. Formally, the cross-modal reconstruction loss can be expressed as:

$$L_{CMR} = \frac{1}{N} \sum_{i=1}^N \frac{1}{M_j} \sum_{m=1}^{M_j} \|f_{dec}(f_{enc}(\Phi_{modal}(X_{im}))) - X_{im}\|_2 \quad (2)$$

B. Modality-wise Data Distillation

The primary objective of the foundation model is to extract highly generalizable latent representations. However, the proxy tasks currently used in pre-training models are often unrelated to the downstream application tasks. We attempt to introduce certain bridging components during the model pre-training stage that can guide the pre-training process to acquire the necessary specific representations. Simultaneously, we hope that these bridging components can facilitate the feature expression of the pre-trained model when applied to downstream tasks. As shown in Fig. 1, the proposed modality-wise data distillation is in conjunction with the cross-modal reconstruction process. Specifically, in the cross-modal reconstruction part, we use data either from another modality image X_{in} to fill in the masked region in X_{im} or from the corresponding learnable modality template.

Specifically, the learnable modality templates $T = \{T_m\}_{m=1}^S$ sized $S \times H \times W \times D$ are initialized with zero, where S represents the number of modalities in the pre-training datasets. Similar to cross-modal reconstruction, the image needed for filling X_{im} is T_m (m represents the corresponding modality) instead of another modality in the modality-wise data distillation process. The remaining steps are consistent with the cross-modal reconstruction process. An example of learned modality templates is shown in the Results section

(Fig. 4), which shows a compact representation of the structural information for each modality in the pre-training datasets.

It is worth noting that the random masked regions are used in the two processes for the same input, which further accelerates the learning. Let us denote the masking strategy for modality-wise data distillation as $\Phi_{distill}$, and the corresponding loss can be expressed as:

$$L_{MD} = \frac{1}{N} \sum_{i=1}^N \frac{1}{M_j} \sum_{m=1}^{M_j} \|f_{dec}(f_{enc}(\Phi_{distill}(X_{im})) - X_{im})\|_2 \quad (3)$$

Cross-modal reconstruction and modality-wise data distillation are performed simultaneously. The model needs to learn not only the structural information of a specific modality to form the distilled modality templates but also the transformation relationship between modalities. The representations learned by our pre-trained model are considered modality-agnostic and contain fused representations of different modalities.

C. Modality-aware Contrastive Learning

During the pre-training, we also need to consider the data distribution bias across different datasets due to the variances in MRI imaging equipment, acquisition protocols, etc. It could largely enhance the generalization of pre-trained models in the downstream tasks when unseen data are applied.

Inspired by the use of contrastive learning for aligning the paired features in multi-modal pre-training schemes, e.g., InfoNCE loss in CLIP, we first compose the positive and negative pairs. A sample masked and mixed with another modality together with another sample masked and mixed with the corresponding modality template forms a positive pair. Similarly, a sample masked and mixed with another modality together with another sample (from different datasets) masked and mixed with the corresponding modality template forms a negative pair. In such a way, the model tends to learn both dataset- and modality-independent features.

We denote the sets obtained by encoding N data samples using the aforementioned two masking strategies as $\{\Phi_{modal}(X_{im})\}_{i=1}^N$ and $\{\Phi_{distill}(X_{im})\}_{i=1}^N$, respectively. As shown in Fig. 1, We use contrastive loss [36] to bring the distance between features of positive pairs closer while repelling the distance between features of negative pairs. This can be formalized as:

$$L_{CL} = -\frac{1}{2N} \sum_{i=1}^N \left(\log \frac{e^{f_i \cdot g_i / \tau}}{\sum_{j=0}^k e^{f_i \cdot g_j / \tau}} + \log \frac{e^{g_i \cdot f_i / \tau}}{\sum_{j=0}^k e^{g_i \cdot f_j / \tau}} \right) \quad (4)$$

where $f_i = f_{enc}(\Phi_{modal}(X_{im}))$ represents current modality image masked and mixed with another modality image in the same sample, $g_i = f_{enc}(\Phi_{distill}(X_{im}))$ represents current modality image masked and mixed with corresponding distilled modality template. $g_{j, j \neq i}$ represents modality images from other samples (datasets) masked and mixed with corresponding modality templates. N is the number of samples and τ is the distillation temperature.

Overall loss function. In summary, the total loss function is a combination of L_{CMR} , L_{MD} , and L_{CL} :

$$L = L_{CMR} + \lambda_{MD} \cdot L_{MD} + \lambda_{CL} \cdot L_{CL} \quad (5)$$

where λ_{MD} as well as λ_{CL} are used to balance the corresponding loss term contributions, which are both set to 1.0 in the experiments for equal treatment.

D. Distilled Modality Template for Downstream Tasks

In this section, we will elaborate on how the distilled modality templates obtained from pre-training can be applied in downstream tasks, as a form of data augmentation. As shown in Fig. 2, in the downstream fine-tuning stage, the distilled modality templates are frozen. Let $\mathcal{D}_{ds} = \{(X_i, Y_i)\}_{i=1}^M$ denote the downstream dataset, where M represents the number of annotated samples. X_i is the multi-modal MRI input volume, and Y_i represents the corresponding label, which can be a segmentation map for segmentation tasks or a one-hot vector for classification tasks. Specifically, we randomly select n modalities in X_i and replace them with the corresponding modalities from $\{T_m\}_{m=1}^S$, obtaining two augmented copies X'_i and X''_i . The encoded features of these two copies are $f_{enc}(X'_i)$ and $f_{enc}(X''_i)$, respectively. Since the two embeddings are representations of the same sample with different numbers of replaced modalities, we use the L2 norm to maintain semantic consistency in the feature space.

$$L_{cons} = \frac{1}{N} \sum_{i=1}^N \|f_{enc}(X'_i) - f_{enc}(X''_i)\|_2 \quad (6)$$

Subsequently, the features of the two copies are decoded to the output space to calculate supervision loss with the ground-truth annotations. The overall loss is:

$$L_{total} = \frac{1}{N} \sum_{i=1}^N L_{sl}(f(X'_i), Y_i) + L_{sl}(f(X''_i), Y_i) + \lambda_{cons} * L_{cons} \quad (7)$$

where λ_{cons} is the weight of the consistency loss L_{cons} term and L_{sl} is the supervision loss used in segmentation or classification tasks, e.g., Dice Loss in segmentation or CrossEntropyLoss in classification.

For the uni-modal input scenario, instead of replacing the selected modalities with distilled modality templates, we perform a partially masking strategy like Algorithm 1 where X_{im} is replaced with the corresponding distilled modality template. Then we randomly mask the uni-modal input volume twice to obtain two augmented copies of X_{im} , and the remaining procedures are the same as the aforementioned multi-modal scenario.

IV. EXPERIMENTS

A. Datasets

Pre-training Datasets We collect five publicly available mpMRI datasets for pre-training, spanning 8 modalities with a total of 3,755 cases and 16,022 scans. The task types and number of modalities for each dataset are summarized in Table I.

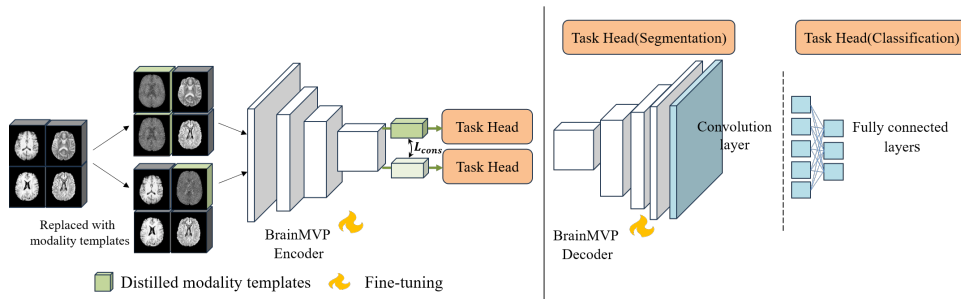


Fig. 2. Modality-wise data distillation for **downstream tasks**. The input multi-modal MRI images are randomly selected to replace a certain number of modalities with the corresponding modality templates. Then L2 norm is used to ensure feature consistency between the two replacement copies. Finally, the task head is replaced with corresponding modules based on the task type.

TABLE I
DETAILS OF DATASETS USED IN OUR WORK.

Dataset	Task type	Modality type	cases
Pre-training			
BraTS2021 [4]	-	T1,T1CE,T2,FLAIR	3755
BraTS2023-SSA [1]	-	T1,T1CE,T2,FLAIR	75
BraTS2023-MEN [26]	-	T1,T1CE,T2,FLAIR	1141
BrainAtlas [22]	-	T1,T2,MRA,PD	568
UCSF-PDGM [7]	-	T1,T1CE,T2,FLAIR,DWI,ADC	501
Downstream			
BraTS-PED [24]	seg. (pediatric tumor)	T1,T1CE,T2,FLAIR	99
BraTS2023-MET [34]	seg. (brain metastases)	T1,T1CE,T2,FLAIR	237
ISLES22 [21]	seg. (ischemic stroke lesion)	FLAIR,DWI,ADC	238
MRBrainS13 [33]	seg. (CF,GM,WM)	T1,T1CE,FLAIR	20
UPENN-GBM [6]	seg. (glioblastoma)	T1,T1CE,T2,FLAIR	127
VSseg [38]	seg. (vestibular schwannoma)	T1	242
BraTS2018 [5]	cls. (HGG and LGG)	T1,T1CE,T2,FLAIR	285
ADNI [23]	cls. (MCI and NC)	T1	1348
ADHD-200 [12]	cls. (ADHD and NC)	T1	767
ABIDE-I [13]	cls. (ASD and NC)	T1	819

seg.: segmentation; cls.: classification; CF: Cerebrospinal Fluid; GM: Gray Matter; WM: White Matter; HGG: Higher Grade Glioma; LGG: Lower Grade Glioma; MCI: Mild Cognitive Impairment; NC: Normal Control; ADHD: Attention Deficit Hyperactivity Disorder; ASD: Autism Spectrum Disorder.

Downstream Datasets Our pre-trained models are evaluated on ten downstream tasks (in a full fine-tuning setting), including six segmentation tasks and four classification applications, as detailed in Table I.

B. Experiments setting

Data Pre-processing. During the pre-training stage, data pre-processing is performed sequentially in Python based on MONAI 1.3.0¹ library. The orientation of the mpMRI images is first unified to the RAS axcodes and co-registered to the same anatomical template. Subsequently, each MRI scan is resampled to an isotropic voxel spacing of $1.0mm \times 1.0mm \times 1.0mm$ using bilinear interpolation, and skull-stripping is performed as well. We linearly clip the pixel values between the 1st and 99th percentiles and re-scale them to $[0, 1]$. The images are then cropped into $96 \times 96 \times 96$ voxel patches centered on either foreground or background areas, to ensure that the modality-wise data distillation is learned sufficiently. We do not apply any other data augmentation techniques.

Implementation details. We use UniFormer [29] for pre-training and downstream tasks, benefiting from its natural multi-modal fusion capabilities. In addition, we have conducted experiments related to the UNET3D [37] architecture. We conduct all experiments using the PyTorch framework on

8 NVIDIA GeForce RTX 4090 GPUs. During the training process, we utilize the AdamW [32] optimizer with a momentum of 0.9 and the weight decay is $1e-5$. We train the model for 1,500 epochs with a batch size of 3 and introduce the modality-aware contrastive learning module at epoch 900. The initial learning rate is set to $3e-4$ and we employ a cosine learning rate decay strategy. Detailed hyperparameters for downstream experiments can be found in the Appendix I.

Comparison methods. We compare our BrainMVP against three different types of approaches, i.e., training from scratch, general domain SSL methods, and medical domain SSL methods. There are three mainstream medical image segmentation networks for training from scratch: UNETR [18], UNET3D [37] [37], and Swin-UNETR [17]. UniFormer [29] is a novel 3D medical image segmentation network initially developed in the field of video object detection and extensive experiments have been conducted to verify its effectiveness. The subsequent SSL methods are pre-trained on the above architectures, allowing for a fair comparison of the impact of different network architectures on the final performance. The baseline SSL methods include MAE3D [19, 9], MIM-based SimMIM [53], and contrastive learning related MoCoV3 [8] for general domain, and MG [59], TransVW [16], GVSL [20], Swin-UNETR [43], and VoCo [51] for medical domain. Specifically, two MIM-based methods in medical domain, namely, DAE [45] and M³AE [31], are also taken as comparisons.

Evaluation metrics. For segmentation tasks, we use the Dice Score and Hausdorff distance at 95th percentile (HD95) as evaluation metrics. For classification tasks, we report accuracy (ACC), area under the curve (AUC), and F1 score for comprehensive assessment with higher metric values indicating better classification performance.

Label efficiency experiments. To validate if our BrainMVP, pre-trained on large-scale mpMRI image datasets, can significantly reduce annotation workload in clinical practices, particularly for handling label-deficient segmentation tasks (which incur higher annotation costs), we conduct label efficiency experiments on five segmentation and one classification datasets. Specifically, we randomly split the training labeled samples into five partitions and gradually increase the training set size by one partition at a time until reaching the full dataset size. The resulted experiments are configured with 20%, 40%, 60%, 80%, and 100% of the total training data.

¹<https://monai.io/>

TABLE II

EXPERIMENTAL RESULTS ON SIX DOWNSTREAM SEGMENTATION DATASETS. WE REPORT THE MEAN DICE SCORE (%) ON EACH DATASET AND THE BEST RESULTS ARE BOLDED. THE SECOND BEST RESULTS ARE UNDERLINED.

Method	Modality	Network	BraTS-PED [24]				BraTS-MET [34]				ISLES22 [21]	MRBrainS13 [33]				VSseg [38]	UPENN-GBM [6]			
			ET	TC	WT	AVG	ET	TC	WT	AVG	IS	CF	GM	WM	AVG	VS	ET	TC	WT	AVG
<i>From Scratch</i>																				
UNETR [18]	-	-	46.46	76.43	78.66	67.19	54.01	54.87	59.44	56.11	74.65	67.55	78.73	83.69	76.66	70.28	83.10	80.88	81.98	81.99
UNET3D [37]	-	-	47.12	81.60	83.94	70.89	56.44	58.75	62.76	59.32	80.94	70.47	73.93	82.96	75.78	69.43	85.65	88.76	86.27	86.89
UniFormer [29]	-	-	46.73	83.87	86.97	72.52	67.22	<u>72.74</u>	<u>70.78</u>	70.25	84.97	77.66	74.09	75.60	75.78	<u>80.33</u>	<u>87.93</u>	91.86	88.81	89.53
Swin-UNETR [17]	-	-	49.66	81.10	84.13	71.63	63.84	67.08	68.58	66.50	75.88	70.35	<u>81.66</u>	84.65	78.89	76.82	87.60	91.15	87.34	88.70
<i>With General SSL</i>																				
MAE3D [19, 9]	Natural	UNETR	46.55	77.08	79.32	67.65	57.45	59.19	62.06	59.57	70.43	68.30	80.57	84.69	77.86	69.57	83.66	80.42	81.86	81.98
SimMIM [53]	Natural	UNETR	45.14	76.59	78.61	66.78	54.46	55.84	58.89	56.40	69.94	68.11	80.49	<u>84.76</u>	77.79	69.08	83.70	81.68	82.44	82.61
MoCov3 [8]	Natural	UNETR	45.66	77.37	79.88	67.64	55.84	56.77	61.62	58.07	70.32	67.97	79.64	84.36	77.32	69.83	83.02	80.54	81.77	81.78
<i>With Medical SSL</i>																				
MG [59]	CXR, CT	UNET3D	47.99	86.69	88.41	<u>74.36</u>	60.11	64.05	65.43	63.19	83.53	71.40	74.71	80.41	75.51	76.33	86.64	90.58	87.03	88.08
TransVW [16]	CT	UNET3D	46.38	80.05	81.98	69.47	56.10	58.69	62.81	59.20	80.24	68.92	80.53	83.70	77.72	71.76	85.95	89.51	86.91	87.46
GVSL [20]	CT	UNET3D	49.05	84.47	86.81	73.45	62.46	66.81	67.26	65.51	80.05	69.34	75.07	82.85	75.75	72.21	87.09	91.75	87.53	88.79
Swin-UNETR* [43]	MRI	Swin-UNETR	49.07	81.74	84.13	71.65	60.60	64.56	64.53	63.23	79.55	69.67	82.09	86.13	79.30	75.55	87.24	91.46	87.28	88.66
VoCo [51]	MRI	Swin-UNETR	48.66	82.26	84.64	71.85	57.49	59.33	63.59	60.13	77.58	71.29	76.43	81.40	76.37	76.45	86.65	90.54	87.34	88.18
DAE [45]	MRI	Swin-UNETR	49.30	82.12	84.78	72.07	62.27	65.99	64.85	64.37	73.92	71.37	78.50	83.20	77.69	74.51	86.90	90.83	87.32	88.35
M ³ AE [31]	MRI	UNET3D	46.77	85.67	86.89	73.11	66.01	70.92	70.18	69.04	83.85	71.32	69.56	79.28	73.39	75.96	87.15	91.90	88.44	89.16
M ³ AE [31]	MRI	UniFormer	<u>50.77</u>	84.95	86.70	74.14	<u>68.08</u>	72.35	70.74	70.39	<u>86.32</u>	<u>78.23</u>	77.20	76.43	77.29	79.31	87.75	<u>92.43</u>	<u>88.72</u>	<u>89.63</u>
BrainMVP	MRI	UNET3D	47.75	85.99	88.46	74.07	67.24	71.27	68.63	69.05	83.31	68.88	74.60	82.66	75.38	76.02	87.30	91.87	<u>88.98</u>	89.38
BrainMVP	MRI	UniFormer	55.45	<u>86.54</u>	<u>88.41</u>	76.80	70.70	75.80	74.52	73.67	86.60	81.04	78.17	81.61	80.27	83.64	88.49	92.48	89.07	90.01

CXR: Chest X-Ray; ET: enhancing tumor; TC: tumor core; WT: whole tumor; AVG: average; IS: Ischemic Stroke; CF: Cerebrospinal Fluid; GM: Gray matter; WM: White matter; VS: Vestibular schwannoma.

The validation and test sets are kept the same for a fair comparison. For the comparison methods, we select representative approaches for each pre-training data modality (natural, CT, and MRI), including MAE3D [19, 9], GVSL [20], MG [59], and VoCo [51]. Notably, we observe that MG [59] exhibits strong generalization performance across many datasets, so we include it in the label efficiency experiments to verify whether our method has superior performance.

V. RESULTS

A. Experimental Results on Ten Downstream Tasks

1) *Segmentation Results: Superior performance on tumor segmentation datasets.* To evaluate the improvement of downstream segmentation performance after pre-training, we select two brain gliomas subregion segmentation datasets, BraTS-PED [24] and UPENN-GBM [6]. BraTS-PED, dedicated to pediatric glioma, comprises only 99 annotated cases, making it a challenging testbed for assessing the generalization capability of pre-trained foundation models. Consequently, initial comparative experiments are conducted on this dataset. As depicted in Table II, SSL methods tailored for medical image domains consistently outperform general SSL methods. It is worth noting that models pre-trained on natural image demonstrate poorer generalization on medical image domains. Specifically, the best average Dice Score achieved by general SSL methods based on MIM is 67.65%, which is 9.15% lower than BrainMVP’s best result of 76.80%. Also, MoCoV3 [8] performs less effectively, achieving 9.16% lower in Dice Score compared to BrainMVP. This disparity arises because typical pre-training methods developed primarily for 2D image tasks often require full images or large patches as input, which is usually impractical for 3D medical images. Our BrainMVP also outperforms medical SSL methods based on mask modeling, such as M³AE [31] (76.80% vs. 74.14%) and DAE [45] (76.80% vs. 72.07%). We further validate the effectiveness of BrainMVP on UPENN-GBM [6], as shown in Table II.

BrainMVP achieves an average Dice Score of 90.01% and outperforms state-of-the-art methods.

Performance improvement on brain structure segmentation dataset. We utilize the MRBrainS13 [33] dataset for the segmentation of normal brain structures to assess the efficacy of BrainMVP in scenarios with limited normal brain structure samples during pre-training. As detailed in Table II, our BrainMVP achieves an average Dice Score of 80.27%. In contrast, MG [59], employing multiple proxy tasks, attains 75.51%, and VoCo [51], leveraging contrastive learning, achieves 76.37%. Based on the UniFormer [29] architecture, BrainMVP surpasses all previous methods and demonstrates a notable 4.49% average Dice Score improvement over training from scratch. This underscores its robust capability to effectively improve downstream tasks performance, even under constraints with limited normal brain structure data samples in pre-training.

Strong generalization performance on Unseen datasets. Given that our pre-training datasets primarily include normal brain structures and those afflicted with glioma, we aim to verify the generalization capabilities of BrainMVP on other types of samples. To assess this, we evaluate our BrainMVP on three datasets: BraTS-MET [34], ISLES22 [21], and VSseg [38]. For the BraTS-MET [34] dataset focusing on brain metastasis subregion segmentation, as seen in Table II, our BrainMVP achieves an average Dice Score of 73.67%. Further, BrainMVP notably outperforms existing state-of-the-art methods in medical applications, including MG [59] (63.19%), and Swin-UNETR* [43] (63.23%). In the context of the ISLES22 [21] ischemic stroke segmentation task, which involves abnormalities distinct from tumors targeted in pre-training, BrainMVP achieves substantial improvement compared to MG [59] (86.60% vs. 83.53%) and GVSL [20] (86.60% vs. 80.05%). For the VSseg [38] dataset focusing on vestibular schwannoma segmentation task, in previous methods, M³AE [31] achieves the best performance with 79.31% Dice Score, while our BrainMVP outperforms all previous methods with 83.64% Dice Score, proving the effectiveness of BrainMVP.

TABLE III

EXPERIMENTAL RESULTS ON FOUR DOWNSTREAM **CLASSIFICATION** DATASETS. WE REPORT THE OVERALL ACCURACY (ACC), AREA UNDER THE CURVE (AUC) AND F1 SCORE ON EACH DATASET. THE BEST RESULTS ARE BOLDED AND THE SECOND BEST RESULTS ARE UNDERLINED.

Method	Modality	Network	BraTS2018 [5]			ADNI [23]			ADHD-200 [12]			ABIDE-I [13]		
			ACC	AUC	F1	ACC	AUC	F1	ACC	AUC	F1	ACC	AUC	F1
<i>From Scratch</i>														
UNETR [18]	-	-	0.7895	0.7817	0.6621	0.5672	0.6066	0.5645	0.6688	0.6523	0.6204	0.6121	0.5478	0.5507
UNET3D [37]	-	-	0.7368	0.7373	0.4242	0.5756	0.4966	0.3653	0.6494	0.6798	0.4265	0.6061	0.5059	0.4591
UniFormer [29]	-	-	0.7762	0.7719	0.6994	0.5546	0.6343	0.5526	0.6039	0.6387	0.5796	0.5879	0.4433	0.4292
Swin-UNETR [17]	-	-	0.7018	0.7143	0.6069	0.5672	0.5853	0.5650	0.6494	0.6950	0.6240	0.6121	0.5530	0.5596
<i>With General SSL</i>														
MAE3D [19, 9]	Natural	UNETR	0.7018	0.6754	0.5645	0.5756	0.5414	0.5651	0.6169	0.6489	0.5906	0.6061	0.4983	0.4591
SimMM [53]	Natural	UNETR	0.7368	0.8349	0.7077	0.6218	0.6026	0.5446	0.6234	0.6567	0.5790	0.5394	0.5819	0.5318
MoCov3 [8]	Natural	UNETR	0.7368	0.8135	<u>0.7304</u>	0.6092	0.5769	0.5996	0.6104	0.6265	0.6007	0.5939	<u>0.6284</u>	0.5890
<i>With Medical SSL</i>														
MG [59]	CXR, CT	UNET3D	0.7368	<u>0.9286</u>	0.4242	0.5756	0.5496	0.3653	0.6169	0.6980	0.6141	0.6121	0.6266	0.5892
TransVW [16]	CT	UNET3D	0.7368	0.7222	0.4242	0.4958	0.6661	0.4450	0.6818	0.7228	0.6271	<u>0.6424</u>	0.5292	0.5003
GVSL [20]	CT	UNET3D	0.7895	0.8516	0.7286	0.5966	0.6661	0.5959	0.6623	0.7309	0.6565	0.6242	0.5244	0.4701
Swin-UNETR* [43]	MRI	Swin-UNETR	0.7368	0.5032	0.4242	0.5462	0.5517	0.5461	0.6299	0.6437	0.5953	0.6303	0.4993	0.3866
VoCo [51]	MRI	Swin-UNETR	0.7368	0.5135	0.4242	0.5210	0.5740	0.5207	0.6558	0.6971	0.6413	0.5818	0.5626	0.5466
DAE [45]	MRI	Swin-UNETR	0.7719	0.8151	0.7120	0.5294	0.5666	0.5294	0.6688	0.7129	0.6548	0.6061	0.5173	0.5548
M ³ AE [31]	MRI	UNET3D	0.7370	0.6984	0.5915	0.6008	0.6338	0.6003	0.6364	0.7049	0.6177	0.6061	0.5453	0.4769
M ³ AE [31]	MRI	UniFormer	<u>0.7895</u>	0.8659	0.7159	0.6092	0.5352	0.5756	0.6169	0.6597	0.6028	0.5636	0.4682	0.4500
BrainMVP	MRI	UNET3D	0.7895	0.7746	0.6621	<u>0.6555</u>	<u>0.6669</u>	<u>0.6421</u>	<u>0.6818</u>	0.7245	<u>0.6665</u>	0.6970	0.5817	0.6327
BrainMVP	MRI	UniFormer	0.8596	0.9452	0.8324	0.6765	0.6964	0.6609	0.6883	<u>0.7249</u>	0.6723	0.6182	0.6329	<u>0.5890</u>

2) *Classification Results*: We select four distinct classification tasks to assess the generalizability of BrainMVP across diverse domains. Initially, experiments are conducted on the BraTS2018 [5] glioblastoma grading task. As depicted in Table III, M³AE [31] achieves the best ACC compared with other methods with 0.7895, while our BrainMVP achieves an outstanding 0.8596, surpassing the state-of-the-art methods by a large margin. For example, M³AE [31] achieves accuracies of 0.7370 and 0.7895, VoCo [51] achieves an accuracy of 0.7368, and GVSL [20] achieves an accuracy of 0.7895. Additionally, BrainMVP exhibits superior performance in F1 score and AUC compared to prior SSL methods, underscoring its efficacy.

Subsequently, BrainMVP is validated on the ADNI [23] dataset to assess its ability to differentiate between healthy and diseased states, a task not represented in the pre-training datasets. Notably, the pre-training datasets comprise only a small fraction of normal brain data (12.5%). Experimental results reveal that BrainMVP displays robust generalization capabilities, achieving the highest accuracy of 0.6765 with the UNET3D [37] network, compared to an accuracy of 0.5756 for training from scratch. In terms of AUC, BrainMVP achieves 0.6964, a 3.03% improvement over the best-performing method GVSL [20] (0.6661). Similarly, BrainMVP outperforms the state-of-the-art methods in F1 score, demonstrating strong generalizability despite disparities between pre-training and downstream tasks.

Moreover, to further investigate BrainMVP’s generalization across diverse tasks and domains, experiments are conducted on the ADHD-200 [12] and ABIDE-I [13] datasets. Results indicate that BrainMVP consistently outperforms the state-of-the-art SSL methods. Specifically, on the ADHD-200 [12] dataset, BrainMVP achieves an accuracy of 0.6883 while the best result achieved in the previous method is 0.6818. On the ABIDE-I [13] dataset, BrainMVP demonstrates 5.46% accuracy improvement, 0.45% AUC improvement, and 4.35% F1 score improvement, establishing BrainMVP as a robust approach surpassing existing SSL methods.

It is important to note that our BrainMVP leverages pre-training on partially relevant normal brain mpMRI images, further validating the strong generalizability and superior performance of our BrainMVP.

B. Results of Label Efficiency Experiments

Fig. 3 illustrates the results of the label efficiency experiments. It can be observed that when BrainMVP is fine-tuned on downstream segmentation and classification tasks with varying ratios of labeled training data, BrainMVP consistently exhibits superior performance compared to the representative methods. As the labeled data increase from 20% to 40%, BrainMVP shows a significant performance improvement on multiple datasets, such as BraTS-PED [24] (Dice Score 66.41% to 70.46%), BraTS-MET [34] (Dice Score 60.45% to 70.12%), and ISLES22 [21] (Dice Score 73.27% to 84.03%). Similarly, on the classification task of BraTS2018 [5], the AUC also shows a substantial increase (0.6833 to 0.8008).

It is noteworthy that with only 40% of the labeled data, BrainMVP can achieve performance on par with or even surpassing those of other methods using the fully labeled data. With just 20% of the labeled data, BrainMVP can attain 66.41% Dice Score on the BraTS-PED [24] dataset, 70.39% Dice Score on the VSseg [38] dataset, and 86.82% Dice Score on the UPENN-GBM [6] dataset, while the corresponding best-performing methods achieve 59.50%, 52.31%, and 80.97%, respectively. This demonstrates the excellent efficiency of our BrainMVP, which can lift the annotation requirements in clinical practices.

C. Ablation Study

We conduct comprehensive ablation experiments on three key components of the proposed BrainMVP: cross-modal reconstruction, modality-wise data distillation, and modality-aware contrastive learning on the BraTS-PED [24], BraTS2018 [5], and ADNI [23] datasets. The results are shown in Table IV.

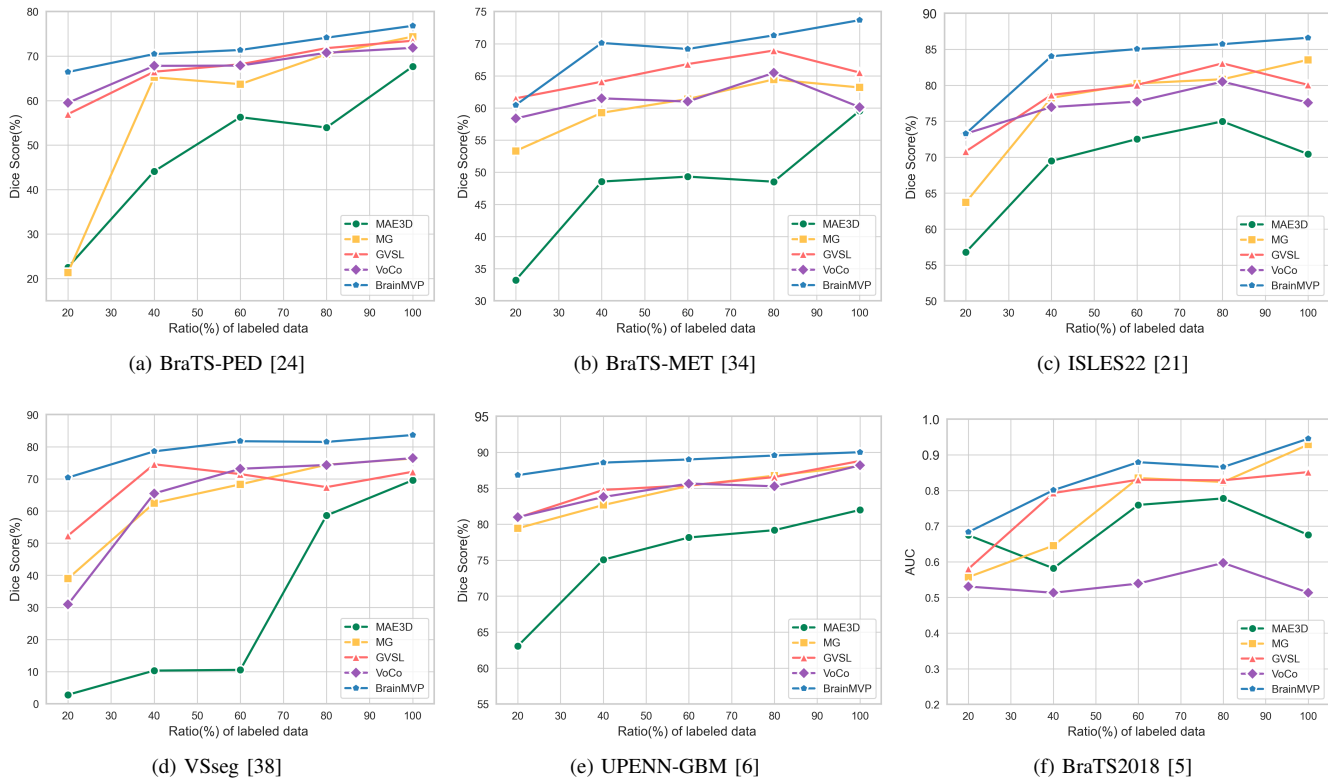


Fig. 3. Label efficiency results of the downstream segmentation and classification tasks. We report the mean Dice Score (%) in segmentation and area under the curve (AUC) in classification.

TABLE IV
ABLATION EXPERIMENTAL RESULTS ON BRA-TS-PED [24], BRA-TS2018 [5] AND ADNI [23] DATASETS.

Task			BraTS-PED [24]		BraTS2018 [5]			ADNI [23]		
CMR	MD	MCL	Dice Score (%)	ACC	AUC	F1	ACC	AUC	F1	
✗	✗	✗	72.52	0.7762	0.7719	0.6994	0.5546	0.6343	0.5526	
✓	✗	✗	75.16	0.7895	0.8056	0.7286	0.6261	0.6770	0.5552	
✓	✓	✗	75.87	0.8421	0.9032	0.8081	0.6261	0.6835	0.6187	
✓	✓	✓	76.80	0.8596	0.9452	0.8324	0.6765	0.6964	0.6609	

CMR: cross-modal reconstruction; MD: Modality-wise data distillation; MCL: modality-aware contrastive learning

Cross-modal reconstruction. As shown in Table IV, when our proposed cross-modal reconstruction is added to the pre-training, results get a notable performance improvement, specifically from Dice Score 72.52% to 75.16% on the BraTS-PED [24] dataset, AUC 0.7719 to 0.8056 on the BraTS2018 [5] and ACC 0.5546 to 0.6261 on the ADNI [23] dataset. In addition, for the BraTS-PED [24] tumor subregion segmentation task, which has a greater demand for mpMRI information, the incorporation of cross-modal reconstruction leads to the prominent improvement in performance, indicating that the proposed cross-modal reconstruction can effectively capture the associations between modalities, enabling more efficient multi-modal information fusion.

Modality-wise data distillation. Subsequently, we assess the efficacy of the modality-wise data distillation module. Importantly, integrating this module to learn aggregated data across diverse modalities necessitates employing mutual learning on downstream tasks to establish a linkage between upstream and downstream processes. As seen in Table IV,

the AUC in the BraTS2018 [5] tumor subtype classification task exhibits a noticeable improvement (from 0.8056 to 0.9032), suggesting that the distilled modality templates obtained through pre-training can effectively enrich the diversity of downstream data, thereby facilitating our BrainMVP’s capability for generalized representation.

Modality-aware contrastive learning. Finally, we explore the role of modality-aware contrastive learning. It can be observed that with the incorporation of the modality-aware contrastive learning component, the performance of BrainMVP continues to rise across multiple datasets. On the BraTS-PED [24] dataset, the average Dice Score is improved from 75.87% to 76.80%, and on the BraTS2018 [5] dataset for tumor subtype classification, the AUC increases from 0.9032 to 0.9452. On the ADNI [23] dataset, the ACC is enhanced from 0.6261 to 0.6765. Modality-aware contrastive learning relies on cross-modal reconstruction and modality-wise data distillation, and with the integration of these three components, our BrainMVP achieves the best results across multiple datasets,

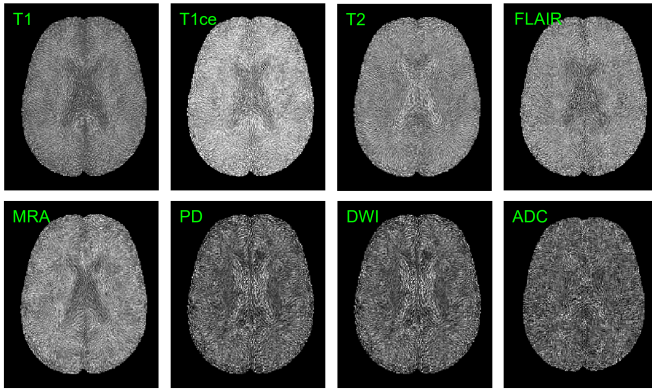


Fig. 4. Visualization of the pre-trained modality templates based on modality-wise data distillation.

demonstrating the effectiveness of the proposed BrainMVP pre-training framework.

VI. FINAL REMARKS

In this paper, we propose an efficient multi-modal vision pre-training paradigm, BrainMVP, that aligns the features across modalities for multi-parametric brain MRI image analysis. Inspired by the structural similarities between different modalities of MRI images, we design cross-modal reconstruction to learn the correlations between modalities. Meanwhile, we leverage single-channel modality image input for handling an arbitrary number of modalities of MRI datasets so as to scale up the pre-training datasets. Subsequently, we learn the condensed structural representation of the pre-trained specific modality based on modality-wise data distillation and build the association between pre-training and downstream tasks by mixing downstream input modality images with condensed modality templates. In addition, we propose modality-aware contrastive learning to ensure semantic consistency of different masking replicas while enhancing the model’s ability to extract discriminative information between different samples. Through extensive experiments on ten downstream datasets, our proposed BrainMVP demonstrates superior performance and strong generalizability across diverse tasks involving multiple types of abnormalities. The label efficiency experiment shows that we can achieve the performance of the state-of-the-art methods using only 40% of the labeled data, demonstrating the potential of BrainMVP in real-world clinical practice.

REFERENCES

- [1] Maruf Adewole et al. “The brain tumor segmentation (brats) challenge 2023: Glioma segmentation in sub-saharan africa patient population (brats-africa)”. In: *ArXiv* (2023).
- [2] Jean-Baptiste Alayrac et al. “Flamingo: a visual language model for few-shot learning”. In: *Advances in neural information processing systems* 35 (2022), pp. 23716–23736.
- [3] Saeid Asgari Taghanaki et al. “Deep semantic segmentation of natural and medical images: a review”. In: *Artificial Intelligence Review* 54 (2021), pp. 137–178.
- [4] Ujjwal Baid et al. “The rsna-asnr-miccai brats 2021 benchmark on brain tumor segmentation and radiogenomic classification”. In: *arXiv preprint arXiv:2107.02314* (2021).
- [5] Spyridon Bakas et al. “Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the BRATS challenge”. In: *arXiv preprint arXiv:1811.02629* (2018).
- [6] Spyridon Bakas et al. “Multi-parametric magnetic resonance imaging (mpMRI) scans for de novo Glioblastoma (GBM) patients from the University of Pennsylvania Health System (UPENN-GBM)”. In: *The Cancer Imaging Archive (TCIA) Public Access* (2021).
- [7] Evan Calabrese et al. “The University of California San Francisco preoperative diffuse glioma MRI dataset”. In: *Radiology: Artificial Intelligence* 4.6 (2022), e220058.
- [8] Xinlei Chen, Saining Xie, and Kaiming He. “An empirical study of training self-supervised vision transformers”. In: *Proceedings of the IEEE/CVF international conference on computer vision*. 2021, pp. 9640–9649.
- [9] Zekai Chen et al. “Masked image modeling advances 3d medical image analysis”. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 2023, pp. 1970–1980.
- [10] Zhihong Chen et al. “Multi-modal masked autoencoders for medical vision-and-language pre-training”. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer. 2022, pp. 679–689.
- [11] Junlong Cheng et al. “Sam-med2d”. In: *arXiv preprint arXiv:2308.16184* (2023).
- [12] ADHD-200 consortium. “The ADHD-200 consortium: a model to advance the translational potential of neuroimaging in clinical neuroscience”. In: *Frontiers in systems neuroscience* 6 (2012), p. 62.
- [13] Adriana Di Martino et al. “The autism brain imaging data exchange: towards a large-scale evaluation of the intrinsic brain architecture in autism”. In: *Molecular psychiatry* 19.6 (2014), pp. 659–667.
- [14] Yuhang Ding, Xin Yu, and Yi Yang. “RFNet: Region-aware Fusion Network for Incomplete Multi-Modal Brain Tumor Segmentation”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 3975–3984. (Visited on 02/26/2024).
- [15] Xiuye Gu et al. “Open-vocabulary object detection via vision and language knowledge distillation”. In: *arXiv preprint arXiv:2104.13921* (2021).
- [16] Fatemeh Haghighi et al. “Transferable visual words: Exploiting the semantics of anatomical patterns for self-supervised learning”. In: *IEEE transactions on medical imaging* 40.10 (2021), pp. 2857–2868.
- [17] Ali Hatamizadeh et al. “Swin unetr: Swin transformers for semantic segmentation of brain tumors in mri images”. In: *International MICCAI Brainlesion Workshop*. Springer. 2021, pp. 272–284.
- [18] Ali Hatamizadeh et al. “Unetr: Transformers for 3d medical image segmentation”. In: *Proceedings of the*

- IEEE/CVF winter conference on applications of computer vision*. 2022, pp. 574–584.
- [19] Kaiming He et al. “Masked autoencoders are scalable vision learners”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2022, pp. 16000–16009.
- [20] Yuting He et al. “Geometric visual similarity learning in 3d medical image self-supervised pre-training”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 9538–9547.
- [21] Moritz R Hernandez Petzsche et al. “ISLES 2022: A multi-center magnetic resonance imaging stroke lesion segmentation dataset”. In: *Scientific data* 9.1 (2022), p. 762.
- [22] IXI Dataset. Accessed: July 5, 2024. 2024. URL: <https://brain-development.org/ixi-dataset/>.
- [23] Clifford R Jack Jr et al. “The Alzheimer’s disease neuroimaging initiative (ADNI): MRI methods”. In: *Journal of Magnetic Resonance Imaging: An Official Journal of the International Society for Magnetic Resonance in Medicine* 27.4 (2008), pp. 685–691.
- [24] Anahita Fathi Kazerooni et al. “The brain tumor segmentation (BRATS) challenge 2023: Focus on pediatrics (CBTN-CONNECT-DIPGR-ASNR-MICCAI BraTS-PEDs)”. In: *arXiv preprint arXiv:2305.17033* (2023).
- [25] Aishik Konwer et al. “Enhancing Modality-Agnostic Representations via Meta-Learning for Brain Tumor Segmentation”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2023, pp. 21415–21425. (Visited on 02/26/2024).
- [26] Dominic LaBella et al. “The asnr-miccai brain tumor segmentation (brats) challenge 2023: Intracranial meningioma”. In: *arXiv preprint arXiv:2305.07642* (2023).
- [27] Junnan Li et al. “Align before fuse: Vision and language representation learning with momentum distillation”. In: *Advances in neural information processing systems* 34 (2021), pp. 9694–9705.
- [28] Junnan Li et al. “Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models”. In: *International conference on machine learning*. PMLR. 2023, pp. 19730–19742.
- [29] Kunchang Li et al. “Uniformer: Unifying convolution and self-attention for visual recognition”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023).
- [30] Liunian Harold Li et al. “Grounded language-image pre-training”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 10965–10975.
- [31] Hong Liu et al. “M³AE: multimodal representation learning for brain tumor segmentation with missing modalities”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 37. 2. 2023, pp. 1657–1665.
- [32] Ilya Loshchilov and Frank Hutter. “Decoupled weight decay regularization”. In: *arXiv preprint arXiv:1711.05101* (2017).
- [33] Adrienne M Mendrik et al. “MRBrainS challenge: online evaluation framework for brain image segmentation in 3T MRI scans”. In: *Computational intelligence and neuroscience* 2015.1 (2015), p. 813696.
- [34] Ahmed W Moawad et al. “The brain tumor segmentation (brats-mets) challenge 2023: Brain metastasis segmentation on pre-treatment mri”. In: *ArXiv* (2023).
- [35] Or Patashnik et al. “Styleclip: Text-driven manipulation of stylegan imagery”. In: *Proceedings of the IEEE/CVF international conference on computer vision*. 2021, pp. 2085–2094.
- [36] Alec Radford et al. “Learning transferable visual models from natural language supervision”. In: *International conference on machine learning*. PMLR. 2021, pp. 8748–8763.
- [37] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. “U-net: Convolutional networks for biomedical image segmentation”. In: *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5–9, 2015, proceedings, part III 18*. Springer. 2015, pp. 234–241.
- [38] Jonathan Shapey et al. “Segmentation of vestibular schwannoma from MRI, an open annotated dataset and baseline algorithm”. In: *Scientific Data* 8.1 (2021), p. 286.
- [39] Junjie Shi et al. “M²FTrans: Modality-Masked Fusion Transformer for Incomplete Multi-Modality Brain Tumor Segmentation”. In: *IEEE Journal of Biomedical and Health Informatics* (2023).
- [40] Aron S Talai et al. “Utility of multi-modal MRI for differentiating of Parkinson’s disease and progressive supranuclear palsy using machine learning”. In: *Frontiers in Neurology* 12 (2021), p. 648548.
- [41] Aiham Taleb et al. “Multimodal self-supervised learning for medical image analysis”. In: *International conference on information processing in medical imaging*. Springer. 2021, pp. 661–673.
- [42] Aiham Taleb et al. “Self-supervised learning for medical images by solving multimodal jigsaw puzzles”. In: *Ieee Transactions on Medical Imaging* 12729 (2017), pp. 661–673.
- [43] Yucheng Tang et al. “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*. 2022, pp. 20730–20740.
- [44] Yao-Hung Hubert Tsai et al. “Conditional contrastive learning: Removing undesirable information in self-supervised representations”. In: *arXiv e-prints* (2021), arXiv–2106.
- [45] Jeya Maria Jose Valanarasu et al. “Disruptive Autoencoders: Leveraging Low-level features for 3D Medical Image Pre-training”. In: *arXiv preprint arXiv:2307.16896* (2023).

- [46] Fuying Wang et al. “Multi-granularity cross-modal alignment for generalized medical visual representation learning”. In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 33536–33549.
- [47] Haoyu Wang et al. “Sam-med3d”. In: *arXiv preprint arXiv:2310.15161* (2023).
- [48] Hu Wang et al. “Multi-Modal Learning With Missing Modality via Shared-Specific Feature Modelling”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023, pp. 15878–15887. (Visited on 02/26/2024).
- [49] Tongzhou Wang et al. “Dataset distillation”. In: *arXiv preprint arXiv:1811.10959* (2018).
- [50] Yiqing Wang et al. “Swinmm: masked multi-view with swin transformers for 3d medical image segmentation”. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer. 2023, pp. 486–496.
- [51] Linshan Wu, Jiabin Zhuang, and Hao Chen. “VoCo: A Simple-yet-Effective Volume Contrastive Learning Framework for 3D Medical Image Analysis”. In: *arXiv preprint arXiv:2402.17300* (2024).
- [52] Yutong Xie et al. “ReFs: A hybrid pre-training paradigm for 3D medical image segmentation”. In: *Medical Image Analysis* 91 (2024), p. 103023.
- [53] Zhenda Xie et al. “Simmm: A simple framework for masked image modeling”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2022, pp. 9653–9663.
- [54] Yiwen Ye et al. “Uniseg: A prompt-driven universal segmentation model as well as a strong representation learner”. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer. 2023, pp. 508–518.
- [55] Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. “Dataset condensation with gradient matching”. In: *arXiv preprint arXiv:2006.05929* (2020).
- [56] Hong-Yu Zhou et al. “A unified visual information preservation framework for self-supervised pre-training in medical image analysis”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023).
- [57] Tongxue Zhou, Su Ruan, and Haigen Hu. “A literature survey of MR-based brain tumor segmentation with missing modalities”. In: *Computerized Medical Imaging and Graphics* 104 (2023), p. 102167.
- [58] Yongchao Zhou, Ehsan Nezhadarya, and Jimmy Ba. “Dataset distillation using neural feature regression”. In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 9813–9827.
- [59] Zongwei Zhou et al. “Models genesis”. In: *Medical image analysis* 67 (2021), p. 101840.

APPENDIX AND THE USE OF SUPPLEMENTAL FILES

In this supplementary material, we provide additional information regarding the downstream fine-tuning stage of BrainMVP, along with supplementary results for the HD95 metric results of six downstream segmentation tasks. Appendix I provides a detailed explanation of the implementation details for downstream fine-tuning processes of BrainMVP. Fig. 5 showcases visual experimental results of BrainMVP and other methods on datasets BraTS-PED [24], ISLES22 [21], VSseg [38] and MRBrainS13 [33]. Tabel V and Table VI present experimental results for the HD95 metric on six segmentation tasks, comparing BrainMVP with other methods.

APPENDIX I
IMPLEMENTATION DETAILS

Segmentation. We have built a data augmentation pipeline for segmentation tasks based on the MONAI² framework. The input mpMRI images are first reoriented to the RAS coordinate system, then the image spacing is adjusted to a uniform $1.0mm \times 1.0mm \times 1.0mm$ (for the ISLES22 [21] dataset it's $1.5mm \times 1.5mm \times 1.5mm$) using bilinear interpolation. Subsequently, the pixel grayscale values of the input mpMRI images are normalized from the 5th to the 95th percentile to a range between 0 and 1 for each channel. After cropping the foreground area of the image, we randomly crop a fixed area of $96 \times 96 \times 96$. To avoid over-segmentation, we allow the sampling center to be in the background area. Then, we validate the input patch for random mirror flipping along three axes with a probability of 0.5. Similarly, we perform random intensity offset with 0.1 offset and 1.0 probability and random intensity scaling with a scale factor of 0.1, also with a probability of 1.0. For network training, we use the AdamW [32] optimizer. The initial learning rate is $3e-4$ and is equipped with cosine learning rate decay, weight decay is $1e-3$ for UNETR [18] based models, $1e-4$ for UniFormer [29] and Swin-UNETR [17] based models, and $1e-5$ for UNET3D [37] based models. We train the network with a batch size of 3 for 500 epochs and λ_{cons} is set to 0.1.

Classification. The data augmentation part is different from segmentation in that we resize the input image to a fixed size of $128 \times 128 \times 64$ after normalizing it to fit the training of the comparison method (this generally leads to a decrease in accuracy). Subsequently, we randomly crop a fixed region of $96 \times 96 \times 64$ and then perform the same random data augmentation as segmentation. In the inference stage, we crop an area of $96 \times 96 \times 64$ at the center of the input image. we set batch size to 64 based on gradient accumulation and train all networks for 200 epochs. The remaining hyper-parameters are the same as in segmentation.

TABLE V
EXPERIMENTAL RESULTS ON DATASETS BRATS-PED [24],
BRATS2023-MET [34] AND ISLES22 [21]. WE REPORT THE MEAN
HD95 (\downarrow) ON EACH DATASET.

Method	Modality	Network	BraTS-PED [24]				BraTS2023-MET [34]				ISLES22 [21]
			ET	TC	WT	AVG	ET	TC	WT	AVG	IS
<i>From Scratch</i>											
UNETR [18]	-	-	25.06	39.07	39.14	34.43	44.11	45.22	43.36	44.23	15.48
UNET3D [37]	-	-	22.48	34.02	33.07	29.86	45.68	46.85	39.93	44.15	4.43
UniFormer [29]	-	-	11.55	16.71	16.14	14.80	25.90	28.16	19.97	24.68	4.13
Swin-UNETR [17]	-	-	17.37	22.56	21.03	20.32	28.68	31.03	24.26	27.99	11.31
<i>With General SSL</i>											
MAE3D [19, 9]	Natural	UNETR	25.37	38.43	37.92	33.90	36.89	36.57	38.38	37.28	15.20
SimMIM [53]	Natural	UNETR	24.70	31.61	32.52	29.61	39.37	41.26	40.06	40.23	17.14
MoCoV3 [8]	Natural	UNETR	20.60	31.88	32.12	28.20	41.88	43.17	41.92	42.32	15.04
<i>With Medical SSL</i>											
MG [59]	CXR, CT	UNET3D	19.71	15.72	17.65	17.69	46.39	48.33	42.02	45.58	3.68
TransVW [16]	CT	UNET3D	18.36	25.42	24.67	22.82	47.85	48.06	39.41	45.11	7.93
GVSL [20]	CT	UNET3D	17.45	15.33	16.00	16.26	37.33	38.05	30.61	35.33	9.35
Swin-UNETR* [43]	MRI	Swin-UNETR	18.65	17.44	17.64	17.91	40.57	41.54	33.93	38.68	8.09
VoCo [51]	MRI	Swin-UNETR	18.98	17.21	17.16	17.78	38.52	39.79	34.73	37.68	12.22
DAE [45]	MRI	Swin-UNETR	19.33	21.41	21.71	20.82	37.63	37.37	38.74	37.91	12.50
M ³ AE [31]	MRI	UNET3D	13.48	11.91	10.88	12.09	22.40	23.87	18.96	21.74	4.58
M ³ AE [31]	MRI	UniFormer	16.19	15.95	19.78	17.31	25.89	28.37	24.35	26.21	2.64
BrainMVP	MRI	UNET3D	15.93	7.24	9.81	10.99	20.37	22.50	18.34	20.40	5.85
BrainMVP	MRI	UniFormer	13.93	7.88	14.56	12.12	22.60	25.88	19.83	22.77	2.69

CXR: Chest X-Ray; ET: enhancing tumor; TC: tumor core; WT: whole tumor; AVG: average; CF: Cerebrospinal Fluid; GM: Gray matter; WM: White matter; IS: Ischemic Stroke.

TABLE VI
EXPERIMENTAL RESULTS ON DATASETS MRBRAINS13 [33],
VSSEG [38] AND UPENN-GBM [6]. WE REPORT THE MEAN HD95 (\downarrow)
ON EACH DATASET.

Method	Modality	Network	MRBrainS13 [33]				VSseg [38]	UPENN-GBM [6]			
			CF	GM	WM	AVG	VS	ET	TC	WT	AVG
<i>From Scratch</i>											
UNETR [18]	-	-	4.16	3.46	5.04	4.22	24.54	16.97	24.80	31.00	24.26
UNET3D [37]	-	-	3.24	2.91	3.70	3.28	34.36	5.30	9.34	13.31	9.32
UniFormer [29]	-	-	2.38	2.43	4.04	2.95	5.68	4.46	6.97	11.32	7.58
Swin-UNETR [17]	-	-	3.38	2.65	4.00	3.34	14.12	1.86	7.22	9.15	6.08
<i>With General SSL</i>											
MAE3D [19, 9]	Natural	UNETR	3.69	2.62	3.59	3.30	24.17	15.41	20.10	35.71	23.74
SimMIM [53]	Natural	UNETR	3.84	2.67	3.55	3.35	26.82	17.23	20.71	32.11	23.35
MoCoV3 [8]	Natural	UNETR	3.84	2.99	4.74	3.86	21.35	17.08	19.83	34.35	23.75
<i>With Medical SSL</i>											
MG [59]	CXR, CT	UNET3D	3.47	9.43	12.67	8.52	14.87	2.27	4.29	12.67	6.41
TransVW [16]	CT	UNET3D	3.81	3.45	2.93	3.40	16.83	3.36	5.73	12.95	7.35
GVSL [20]	CT	UNET3D	3.73	3.44	3.28	3.48	11.58	2.23	3.71	9.17	5.03
Swin-UNETR* [43]	MRI	Swin-UNETR	3.33	2.26	2.33	2.64	20.73	2.44	4.07	9.79	5.43
VoCo [51]	MRI	Swin-UNETR	3.14	3.88	7.87	4.96	13.26	28.50	43.05	31.51	34.35
DAE [45]	MRI	Swin-UNETR	3.07	2.27	3.36	2.90	19.84	2.24	3.90	9.56	5.23
M ³ AE [31]	MRI	UNET3D	3.69	3.88	3.01	3.53	9.20	1.85	4.65	8.24	4.91
M ³ AE [31]	MRI	UniFormer	1.89	2.92	4.53	3.11	9.16	4.75	6.54	9.93	7.07
BrainMVP	MRI	UNET3D	3.71	4.92	3.84	4.14	16.41	2.35	4.60	9.13	5.36
BrainMVP	MRI	UniFormer	1.53	5.60	7.02	4.72	6.00	1.48	6.66	10.59	6.24

CXR: Chest X-Ray; ET: enhancing tumor; TC: tumor core; WT: whole tumor; AVG: average; CF: Cerebrospinal Fluid; GM: Gray matter; WM: White matter; VS: Vestibular schwannoma.

²<https://monai.io/>

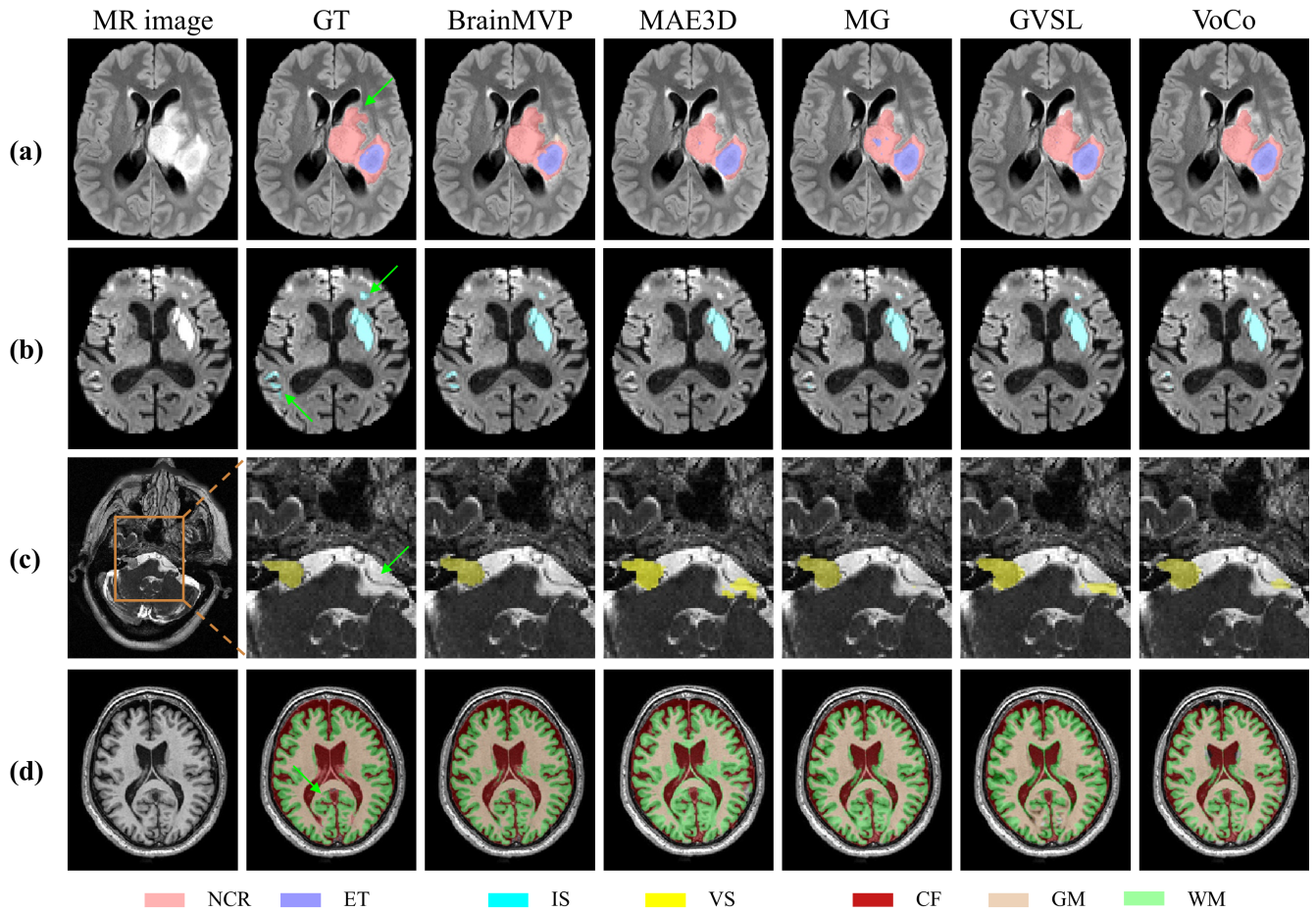


Fig. 5. Visualization results of segmentation tasks. (a) BraTS-PED [24]: pediatric tumor subregion segmentation. NCR: necrotic tumor core; ET: enhancing tumor. (b) ISLES22 [21]: Ischemic Stroke lesion (IS) segmentation. (c) VSseg [38]: Vestibular schwannoma (VS) segmentation. (d) MRBrainS13 [33]: brain structure segmentation. CF: Cerebrospinal Fluid; GM: Gray matter; WM: White matter. GT: ground truth. The green arrows highlight the regions where BrainMVP demonstrates superior performance over other methods.