# Robust Divergence Learning for Missing-Modality Segmentation

Runze Cheng[1,2], Zhongao Sun[1], Ye Zhang[1,3], and Chun Li[1,*]

[1]*MSU-BIT-SMBU Joint Research Center of Applied Mathematics, Shenzhen MSU-BIT University, Shenzhen, 518172, China.*
[2]*Institute of Control Theory and Control Engineering, School of Automation, Beijing Institute of Technology, 100081, Beijing, China.*
[3]*School of Mathematics and Statistics, Beijing Institute of Technology, 100081, Beijing, China.*
Email: 3120220916@bit.edu.cn; sunzhongao0224@gmail.com; ye.zhang@smbu.edu.cn.
*Corresponding author: Chun Li (E-mail: lichun2020@smbu.edu.cn).

*Abstract*—**Multimodal Magnetic Resonance Imaging (MRI) provides essential complementary information for analyzing brain tumor subregions. While methods using four common MRI modalities for automatic segmentation have shown success, they often face challenges with missing modalities due to image quality issues, inconsistent protocols, allergic reactions, or cost factors. Thus, developing a segmentation paradigm that handles missing modalities is clinically valuable. A novel single-modality parallel processing network framework based on Hölder divergence and mutual information is introduced. Each modality is independently input into a shared network backbone for parallel processing, preserving unique information. Additionally, a dynamic sharing framework is introduced that adjusts network parameters based on modality availability. A Hölder divergence and mutual information-based loss functions are used for evaluating discrepancies between predictions and labels. Extensive testing on the BraTS 2018 and BraTS 2020 datasets demonstrates that our method outperforms existing techniques in handling missing modalities and validates each component's effectiveness.**

*Index Terms*—**Missing modality learning, brain-tumor segmentation, divergence learning, knowledge distillation**

## I. INTRODUCTION

Brain tumors are aggressive diseases requiring early detection for effective treatment. MRI is widely used for evaluating brain tumors due to its superior soft tissue contrast and lack of radiation exposure. MRI segmentation is crucial for isolating healthy tissue from abnormal cells, offering various modalities for effective tumor detection, including T1-weighted, T1-weighted post-contrast- enhancement, T2-weighted, and FLAIR. Several existing methods [1]–[6] achieved high accuracy in tumor segmentation when all modalities are available. However, in real-world scenarios, one or more modalities may be missing due to patient movement, hardware issues, or other factors. This "missing modality" problem arises when one or more modalities (e.g., T1w, T2, T1c, and FLAIR) are missing during inference but available during training [7].

Several approaches have been developed to address this issue, which can be categorized into two types [8]: modeling each missing case individually or using a single model to handle all cases. For the former, knowledge distillation is commonly used to transfer competencies from a well-trained

teacher model to a student model designed for specific missing modalities. SMU-Net [9] employed a novel distillation strategy where a multimodal teacher network transfers knowledge to unimodal student networks at both the latent space and network output levels. ProtoKD [10] integrated prototype learning with knowledge distillation, effectively capturing the underlying data distribution. MMCFormer [11] leveraged transformers with auxiliary tokens to facilitate modality-specific representation transfer.

The latter category aims to address all missing-modal situations with a single model, typically involving separate modality encoders to project each modality into a shared latent space before feature fusion. RFNet [12] integrated characteristics from various sources using a region-cognizant component. Moreover, Ting and Liu [13] used modality-specific encoders, a shared decoder, and a strategy to complement incomplete data with complete data. Furthermore, Wang et al. [14] employed specific and shared encoders to handle missing modalities for both segmentation and classification tasks.

However, these approaches have some shortcomings. Using a specific model for each missing modality scenario is training-costly, for example, $2^N - 1$ models need to be trained when there are $N$ modalities [8]. Conversely, a single model for all cases often results in performance deficiencies with few modalities available [9] and high inference costs due to numerous parameters.

Inspired by Chang et al. [15] and high mutual information knowledge transfer learning [16], we process four different modalities individually to preserve unique information and enhance the model's ability to recognize diverse data features, which can handle all cases with signal model with shared backbone. Specifically, we propose a novel mutual information-based metric with Hölder divergence [17] that evaluate discrepancies between the predictions and labels. What is more, a dynamic sharing framework is introduced, which allows the model to adapt its parameters depending on the availability of different modalities.

The main contributions of this paper are summarized as follows: 1. Novel Network Architecture: We propose a new network architecture for parallel computing based on 3D U-Net. This framework combines unimodal parallel processing

TABLE I
MAIN NOTATIONS USED IN THIS WORK. THIS TABLE PROVIDES AN
OVERVIEW OF THE MAIN NOTATIONS USED THROUGHOUT THIS WORK,
OFFERING A CONCISE REFERENCE FOR UNDERSTANDING THE SYMBOLS
AND TERMINOLOGY EMPLOYED IN THE ALGORITHMS DISCUSSED.

| Notation | Definition |
|---|---|
| $x_i$ | the $i^{th}$ modality data of the sample. |
| $d_i$ | The deepest-level feature of the $i^{th}$ modality. |
| $d_f$ | The deepest-level full-modality feature. |
| $h_i$ | The generated single-modality representation of the $i^{th}$ modality. |
| $\widehat{Y}$ | Integrated output under missing modalities. |
| $Y$ | real sample. |
| $p(d_f \mid d_m)$ | The conditional distribution of the feature f given the missing modality information m. |
| $q(d_f \mid d_m)$ | The conditional distribution approximated using variational methods. |

and dynamic network module combinations to handle missing modalities during brain tumor segmentation training. 2. New Metric Introduction: the Hölder divergence and mutual information are introduced to evaluate discrepancies between model predictions and labels. By minimizing the distances, we achieve more accurate feature alignment. 3. Extensive Validation: Extensive experiments on the BraTS 2018 and BraTS 2020 medical image datasets [18] are conducted. The results demonstrate that our method achieves state-of-the-art performance, showcasing its efficiency and practicality. The main notations used in this work is shown in Table I.

## II. METHODOLOGY

### A. Knowledge Distillation for Segmentation Using Hölder Divergence

Brain tumor segmentation, particularly glioma segmentation, involves distinguishing four categories: background, whole tumor, tumor core, and enhancing tumor. Missing modalities can degrade segmentation accuracy. The Hölder divergence is employed for its flexibility and robustness, making it suitable for complex models and non-symmetric data. It supports brain tumor segmentation under missing modalities, maintaining high accuracy in clinical settings.

The loss function using Hölder divergence is:

$$\frac{1}{D \times H \times W} \sum_{dhw} D_\alpha^H(\sigma(\mathbf{S}_{dhw}^p)|\sigma(\mathbf{S}_{dhw}^l)), \quad (1)$$

where $\mathbf{S}_{dhw}^p$ and $\mathbf{S}_{dhw}^l$ are predicted and label probabilities for pixel $(d, h, w)$, and $D_\alpha^H$ denotes Hölder divergence, the definition is shown in Definition 1:

**Definition 1.** (*Hölder Statistical Pseudo-Divergence, HPD [17]*) *HPD pertains to the conjugate exponents $\alpha$ and $\beta$, where $\alpha\beta > 0$. In the context of two densities, $p(x) \in L^\alpha(\Omega, \nu)$ and $q(x) \in L^\beta(\Omega, \nu)$, both of which belong to positive measures absolutely continuous with respect to $\nu$, HPD is defined as the logarithmic ratio gap, as follows: $D_\alpha^H(p(x) : q(x)) =$*

$$-\log\left(\frac{\int_\Omega p(x)q(x)\mathrm{d}x}{\left(\int_\Omega p(x)^\alpha \mathrm{d}x\right)^{\frac{1}{\alpha}}\left(\int_\Omega q(x)^\beta \mathrm{d}x\right)^{\frac{1}{\beta}}}\right), \text{ when } 0 < \alpha < 1 \text{ and}$$
$$\beta = \bar{\alpha} = \frac{\alpha}{\alpha-1} < 0 \text{ or } \alpha < 0 \text{ and } 0 < \beta < 1.$$

### B. High Mutual Information Knowledge Transfer Learning

In clinical practice, the challenge of missing modality segmentation often leads to incomplete information, limited model generalization, and data application constraints. To address these issues, a high mutual information knowledge transfer learning strategy between full and missing modalities is introduced. This strategy maximizes existing modality information, compensating for the loss caused by missing data, thereby enhancing model accuracy and stability with incomplete datasets.

Our approach involves extracting $K$ pairs of feature vectors $\left\{\left(d_f^{(k)}, d_m^{(k)}\right)\right\}_{k=1}^K$ from the encoder layers of both the full and missing modality paths. By calculating the entropy $H(d_f)$ and conditional entropy $H(d_f \mid d_m)$ for each pair, we derive the mutual information $MI(d_f; d_m) = H(d_f) - H(d_f \mid d_m)$, which shows how the full modality path reduces uncertainty given the missing modality information. To estimate these mutual information values accurately, a variational information maximization method [19] is employed.

We approximate the conditional distribution $p(d_f \mid d_m)$ with the variational distribution $q(d_f \mid d_m)$ to optimize the layer-wise mutual information. The optimization process is defined by the following loss function $\mathcal{L}_{\mathcal{MI}}$:

$$-\sum_{k=1}^K \gamma_k \mathbb{E}_{d_f^{(k)}, d_m^{(k)} \sim p\left(d_f^{(k)}, d_m^{(k)}\right)}\left[\log q\left(d_f^{(k)} \mid d_m^{(k)}\right)\right], \quad (2)$$

In our framework, the parameter $\gamma_k$ increases with the layer level $k$, reflecting the richer semantic information in higher network layers. This ensures effective knowledge transfer by assigning higher weights to these layers. The implementation of the variational distribution is given by:

$$\begin{aligned}
&-\log q(d_f \mid d_m)\\
&= \sum_{c=1}^C \sum_{h=1}^H \sum_{w=1}^W \left(\log \sigma_c + \frac{\left(d_f^{c,h,w} - \mu^{c,h,w}(d_m)\right)^2}{2\sigma_c^2}\right)\\
&+ \text{constant},
\end{aligned} \quad (3)$$

where $\mu(\cdot)$ and $\sigma$ represent the heteroscedastic mean and homoscedastic variance of the Gaussian distribution, respectively. $W$ and $H$ denote the width and height of the image, $C$ represents the number of channels, and $\text{constant}$ is a fixed term.

### C. Overall Framework

Let $X$ and $Y$ denote samples from a multimodal dataset, where $X = \{x_j\}_{j=1}^M$ contains $M$ samples with $N$ modalities per sample: $x_j = \{x_j^i\}_{i=1}^N$, with $x_j^i \in X$ representing the $i^{th}$ modality data of the $j^{th}$ sample. Corresponding label is $y_j$. To leverage features from different modalities, a parallel 3D U-Net-based network is designed. Each sample's data is first encoded into a common feature space by channel encoders $f_i$, unifying data representation across channels.
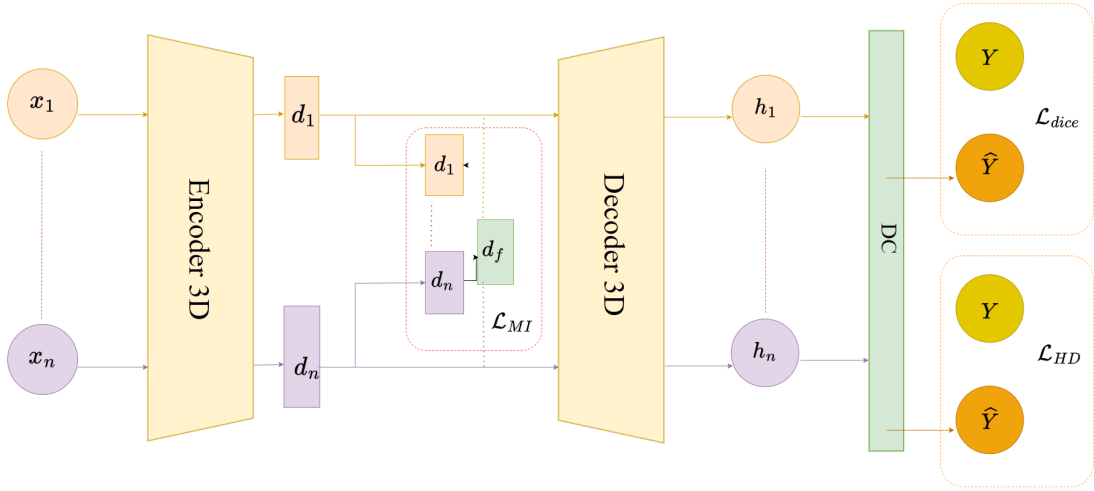
Fig. 1. The Framework of Robust Divergence Learning for Missing-Modality Segmentation. This figure illustrates the overall structure of the proposed robust divergence learning approach, specifically designed to address segmentation challenges in scenarios where certain modalities are missing.

Subsequently, each modality's data is independently input into a parameterized shared backbone $T(\cdot; \theta)$ to generate unique single-modality representations $h_i$: $h_i = T(f_i(x_i); \theta)$.

The final output in our shared network architecture includes a Dynamic Combination Network Module (DC). For missing modalities, we exclude their representations and use a flexible fusion operator $M(\cdot)$ to integrate remaining single-modality representations $H \subseteq \{h_1, h_2, \ldots, h_n\}$: $\widehat{Y} = M(H)$.

The Dice loss function [20] is utilized to optimize consistency between the fused predicted image $\widehat{Y}$ and target labels $Y$, ensuring precise pixel-wise training:

$$\mathcal{L}_{Dice}(\widehat{Y}, Y) = 1 - \frac{2}{J} \sum_{j=1}^{J} \frac{\sum_{i=1}^{I} \widehat{Y}_{i,j} Y_{i,j}}{\sum_{i=1}^{I} \widehat{Y}_{i,j}^2 + \sum_{i=1}^{I} Y_{i,j}^2}, \quad (4)$$

where $I$ is the total number of voxels and $J$ is the number of classes, $\widehat{Y}_{i,j}$ is the predicted probability of voxel $i$ belonging to class $j$, and $Y_{i,j}$ is the one-hot encoded label.

Facing the challenges of medical image processing with missing modalities, inspired by Hinton's knowledge distillation [21], the high mutual information knowledge transfer loss $\mathcal{L}_{\mathcal{MI}}$ is introduced to enhance model accuracy. After that, the total loss can be obtained:

$$\mathcal{L}_{all} = \mathcal{L}_{Dice}(\widehat{Y}, Y) + \mathcal{L}_{\mathcal{MI}} + \mathcal{L}_{HD}(\widehat{Y}, Y). \quad (5)$$

## III. EXPERIMENTS AND ANALYSIS

### A. Datasets and Evaluation Metrics

To improve the model's logical coherence, result reliability, and algorithm robustness in brain tumor segmentation, this study utilizes the BraTS 2018 and BraTS 2020 datasets [18]. These datasets are widely recognized in the field of medical imaging for multi-classification and segmentation tasks. They are extensive collections of multi-modal MRI scans (T1,

T1Gd, T2, and FLAIR) from patients with high-grade and low-grade gliomas. Expertly annotated, these datasets mark tumor subregions like the enhancing tumor, peritumoral edema, and necrotic core. They are essential for advancing and validating automated brain tumor segmentation algorithms. To evaluate the effectiveness of our method, the Dice Similarity Coefficient (DSC) [20], $\text{Dice}(P, G) = \dfrac{2 \times |P \cap G|}{|P| + |G|}$ is used, which is a common performance metric in medical image analysis. The DSC measures the overlap between the model's output ($P$) and the ground truth ($G$). A higher Dice coefficient indicates better predictive performance.

### B. Training Details

In this study, a PyTorch-based framework [22] (version 2.3.0) is utilized for training all models on a server equipped with dual NVIDIA RTX A6000 GPUs. The standard 3D U-Net architecture [23] is adopted, featuring a single encoder-decoder parallel processing structure that incorporates residual blocks and group normalization techniques. During training, a batch size of 8 is set, and the Adam optimizer [24] is employed to update model parameters, starting with a learning rate of 0.0008 and a weight decay of 0.00001. Training is conducted for 600 epochs to ensure comprehensive learning and performance optimization. Post-training, thorough testing of the model is conducted under all possible channel dropout configurations.

*1) Compare Experimental Models:* The comparative experimental models employed in this study are RFNet [25], MMFormer [26], MA3E [27], MTI [28], and GGDM [29]. Each model makes unique contributions to the field of missing modality segmentation, as outlined below. The results for RFNet, MMFormer, MA3E, and MTI are sourced from their respective original research papers, all adhering to the same experimental configuration as RFNet. Additionally, the GGDM

TABLE II

QUANTITATIVE EVALUATION OF SEGMENTATION RESULTS (DSC ↑) ON BraTS 2018. THIS TABLE PRESENTS THE QUANTITATIVE RESULTS OF SEGMENTATION PERFORMANCE, MEASURED BY THE DICE SIMILARITY COEFFICIENT (DSC), ON THE BraTS 2018 DATASET. THE RESULTS PROVIDE A COMPARATIVE EVALUATION OF THE EFFECTIVENESS OF DIFFERENT SEGMENTATION METHODS, WHERE HIGHER DSC VALUES INDICATE BETTER SEGMENTATION ACCURACY.

| Task | Methods | Fl | T2 | T1c | T1 | T2,Fl | T1c,Fl | T1c,T2 | T1,Fl | T1,T2 | T1,T1c | $\sim$ T1 | $\sim$ T1c | $\sim$ T2 | $\sim$ Fl. | Full | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| WT | RFNet | 85.8 | 85.1 | 73.6 | 74.8 | 89.3 | 89.4 | 85.6 | 89.0 | 85.4 | 77.5 | 90.4 | 90.0 | 89.9 | 86.1 | 90.6 | 85.5 |
| | mmFormer | 86.1 | 81.2 | 72.2 | 67.5 | 87.6 | 87.3 | 83.0 | 87.1 | 82.2 | 74.4 | 88.1 | 87.8 | 87.3 | 82.7 | 89.6 | 82.9 |
| | M3AE | 88.7 | 84.8 | 75.8 | 74.4 | 89.9 | 89.7 | 86.3 | 89 | 86.7 | 77.2 | 90.2 | 89.9 | 88.9 | 85.7 | 90.1 | 85.8 |
| | MTI | 88.4 | 86.6 | 77.8 | 78.7 | 90.3 | 89.5 | 88.2 | 89.7 | 88.1 | 81.8 | 90.6 | 89.7 | 90.4 | 88.4 | 90.6 | 87.3 |
| | GGDM | 89.3 | 87.0 | 79.9 | 75.9 | 90.7 | 90.6 | 88.6 | 90.2 | 88.2 | 81.1 | **91.3** | **91.0** | 90.5 | 87.9 | 91.0 | 87.6 |
| | OUR | **89.8** | **88.2** | **80.5** | **78.8** | **90.8** | **90.7** | **89.3** | **90.4** | **88.7** | **82.0** | **91.3** | **91.0** | **90.9** | **89.2** | **91.3** | **88.2** |
| TC | RFNet | 62.6 | 66.9 | 80.3 | 65.2 | 71.8 | 81.6 | 82.4 | 72.2 | 71.1 | 81.3 | 82.6 | 74.0 | 82.3 | 82.9 | 82.9 | 76.0 |
| | mmFormer | 61.2 | 64.2 | 75.4 | 56.6 | 69.8 | 77.9 | 78.6 | 65.9 | 69.4 | 78.6 | 79.6 | 71.5 | 79.8 | 80.4 | 85.8 | 73.0 |
| | M3AE | 66.1 | 69.4 | 82.9 | 66.4 | 70.9 | 84.4 | 84.2 | 70.8 | 71.8 | 83.4 | 84.6 | 72.7 | 84.1 | 84.4 | 84.5 | 77.4 |
| | MTI | 66.7 | 68.8 | 81.5 | 65.6 | 71.8 | 84.8 | 84.8 | 72.0 | 72.3 | 83.5 | 85.8 | 74.1 | 85.2 | 85.8 | 85.9 | 77.9 |
| | GGDM | **77.3** | 76.3 | 85.3 | 58.1 | 78.5 | **87.0** | **87.6** | 76.3 | 76.8 | 85.6 | **87.1** | 78.3 | 86.5 | 86.2 | 85.8 | 80.8 |
| | OUR | 76.2 | **77.6** | **86.5** | **72.6** | **79.3** | 86.6 | 87.2 | **78.6** | **79.1** | **86.9** | **87.1** | **80.1** | **87.1** | **87.4** | **87.3** | **82.6** |
| ET | RFNet | 35.5 | 43.0 | 67.7 | 32.3 | 45.4 | 72.5 | 70.6 | 38.5 | 42.9 | 68.5 | 73.1 | 46.0 | 71.1 | 70.9 | 71.4 | 56.6 |
| | mmFormer | 39.3 | 43.1 | 72.6 | 32.5 | 47.5 | 75.1 | 74.5 | 43.0 | 45.0 | 74.0 | 75.7 | 47.7 | 75.5 | 74.8 | 77.6 | 59.9 |
| | M3AE | 35.6 | 47.6 | 73.7 | 37.1 | 45.4 | 75.0 | 75.3 | 41.2 | 48.7 | 74.7 | 73.8 | 44.8 | 74.0 | 75.4 | 75.5 | 59.9 |
| | MTI | 40.5 | 41.4 | 75.7 | 44.5 | 48.3 | 76.8 | 77.7 | 44.4 | 47.7 | 77.1 | 76.6 | 50.0 | 77.4 | 78.5 | 80.4 | 62.5 |
| | GGDM | 47.4 | **53.4** | 81.6 | 34.7 | 55.2 | 82.0 | 82.6 | 51.1 | 54.7 | 82.0 | 82.1 | 56.0 | 82.2 | 82.8 | 82.1 | 67.6 |
| | OUR | **48.6** | 52.9 | **82.4** | **48.7** | **55.8** | **83.0** | **83.2** | **53.8** | **56.9** | **82.8** | **83.7** | **58.4** | **83.2** | **83.5** | **84.1** | **69.4** |

method is tested, according to authors' code. These models are outlined below:

1. **RFNet (Ding et al., ICCV 2021) [25]**: RFNet is a region-aware fusion network designed for the segmentation of brain tumors in scenarios with incomplete multi-modal data. 2. **MMFormer (Zhang et al., MICCAI 2022) [26]**: MMFormer is a multimodal medical transformer developed to improve brain tumor segmentation in scenarios involving incomplete multimodal data. 3. **MA3E (Liu et al., AAAI 2023) [27]**: M3AE is a multimodal representation learning approach for brain tumor segmentation that effectively handles missing modalities. 4. **MTI (Ting and Liu, JBHI 2024) [28]**: MTI is a multimodal transformer designed to enhance brain tumor segmentation using incomplete MRI data. 5. **GGMD (Wang et al., AAAI 2024) [29]**: GGMD is a method designed to enhance robustness in brain tumor segmentation when handling missing modalities.

Tables II–III showcase the performance of our research method on the BraTS 2018 and BraTS 2020 datasets, benchmarked against five state-of-the-art brain tumor segmentation techniques, and the symbol $\sim (\cdot)$ in Tables II–III denotes the amissing of a specific modality, with optimal performance results highlighted in black across different tumor types. The results highlight our method's superior performance across all three evaluated tumor regions—Whole Tumor (WT), Tumor Core (TC), and Enhancing Tumor (ET)—achieving the highest average Dice Similarity Coefficient (DSC). Specifically, Table II demonstrates improvements of 0.6% in WT, 1.8% in TC, and 1.8% in ET regions compare to existing state-of-the-art methods on the BraTS 2018 dataset. Similarly, Table III shows enhancements of 0.8% in WT, 0.7% in TC, and 0.9% in ET
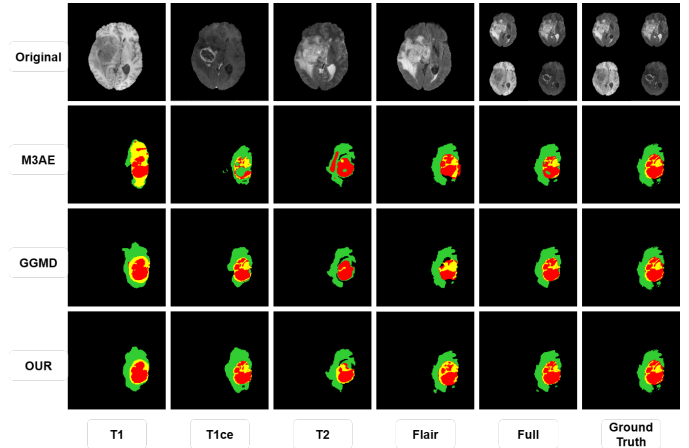


Fig. 2. This figure presents the segmentation results of three models on the BraTS 2018 dataset using different modality inputs. The second row shows the reproduced results of the M3AE, the third row shows the reproduced results of the GGMD, and the fourth row displays the results of ours. Each column represents different input settings: the first four columns show the results for single modality inputs (T1, T1ce, T2, and Flair, respectively), the fifth column displays the results using all four modalities as input simultaneously, and the last column shows the corresponding ground truth.

regions on the BraTS 2020 dataset.

Our method excels notably in scenarios with missing multimodal data, achieving substantial gains of 2.0% to 6.0% in Dice coefficients. For instance, in Table II, when only the T1 modality is available, our method outperforms other advanced algorithms by 0.1%, 6.2%, and 4.2% in WT, TC, and ET regions, respectively. Correspondingly, Table III indicates improvements of 2.4%, 2.9%, and 3.0% under similar conditions.

These findings underscore the robustness of our method in maintaining efficient segmentation performance despite significant deficits in multimodal data.Furthermore, Figure 2 compares our method with other advanced techniques like M3AE and GGMD under single-modal and full-modal input conditions. These segmentation results validate our conclusions from the evaluation metrics, demonstrating that our method surpasses other advanced technologies in handling multimodal data, especially when only a single modality is available. This performance advantage is particularly notable under single-modal input conditions.

*2) Exploration of the Superiority of Hölder Divergence:* To explore the superiority of Hölder divergence, this study conduct experimental comparisons using Hölder divergence and other $f$-divergences [30], including Total Variation [30], Squared Hellinger [30], Kullback-Leibler [30], Neyman $\chi^2$ [30], and Jensen-Shannon divergence [30]. As shown in the table IV, the average Dice coefficient of Hölder divergence reached 80.1% after adjusting the hyperparameter $\alpha$, which is 6.1% higher than the best-performing alternative methods. This significant performance advantage underscores the importance of Hölder divergence in improving model accuracy.

The experimental data consistently demonstrate that the application of Hölder divergence significantly enhances segmentation task performance, validating its effectiveness in the field of medical image processing. Additionally, our research reveals the critical role of Hölder divergence in enhancing knowledge distillation techniques to improve segmentation efficiency, providing valuable references and guidance for future research and development in related technologies. These findings not only deepen our understanding of the potential of Hölder divergence but also provide empirical evidence for optimizing deep learning models using this method.

### C. Exploring the Impact of Hölder Conjugate Exponents on Experimental Results

To further investigate the impact of Hölder conjugate exponents on experimental results, we explore various Hölder conjugate exponents. As shown in Table V, this study compares the performance under different Hölder hyperparameters ($\alpha$), KLD, and without the application of knowledge distillation. The experimental results indicate that when the Hölder divergence hyperparameter $\alpha = 1.1$, the performance improves by an average of 0.7% compared to the case without knowledge distillation and by 6.1% compared to KL divergence. This result underscores the crucial role of selecting an appropriate Hölder conjugate exponent ($\alpha$) in significantly enhancing model performance.

Our findings clearly demonstrate the critical role of Hölder divergence in enhancing knowledge distillation techniques to improve segmentation task efficiency. Throughout our experiments, we consistently observe that setting the Hölder conjugate exponent to $\alpha = 1.1$ markedly improves the model's segmentation performance, further validating the effectiveness of Hölder divergence.

### D. Ablation Study

In this study, we conduct a series of ablation experiments to demonstrate the effectiveness of mutual information knowledge transfer, denoted as $\mathcal{L}_{MI}$, between full and missing modalities. Additionally, we explore the impact of the Hölder divergence-based loss function, $\mathcal{L}_{HD}$, on model performance. Compared to traditional multimodal processing methods, our parallel network framework selectively activates only the data relevant to the available modalities. This approach effectively preserves the unique information of each modality, enhancing the model's ability to recognize diverse data features. Furthermore, the introduction of the Hölder divergence loss function improves the model's performance in handling multimodal data by precisely quantifying the mutual information between modalities, thereby promoting better feature alignment.

First, we evaluate the utility of different components within our network architecture, including the Dice loss function $\mathcal{L}_{Dice}$, mutual information knowledge transfer $\mathcal{L}_{MI}$ between full and missing modalities, and Hölder divergence-based knowledge distillation $\mathcal{L}_{HD}$. The results, as shown in Table VI, indicate that the mutual information knowledge transfer $\mathcal{L}_{MI}$ and Hölder divergence-based knowledge distillation $\mathcal{L}_{HD}$ effectively improve model performance in various missing modality scenarios compared to the traditional segmentation loss $\mathcal{L}_{dice}$ alone. Specifically, when three modalities are missing, these strategies improve performance by 12.2% and 12.7%, respectively; with two missing modalities, they improve performance by 7.8% and 8.2%; with one missing modality, they improve performance by 6.8% and 7.2%; and with all modalities present, they improve performance by 6.5% and 6.9%. On average, the performance improvement for different modality inputs are 8.7% and 9.1%. Moreover, the combination of the parallel network architecture, mutual information knowledge transfer between full and missing modalities, and Hölder divergence-based knowledge distillation achieve the best results, further validating the effectiveness and superiority of our approach.

### E. Conclusion

In this work, we present the quantitative evaluation results of our proposed method for addressing missing modality segmentation, a common challenge in clinical practice. Our approach utilizes a 3D U-Net combined with a parallel network architecture, integrating mutual information knowledge transfer and knowledge distillation based on Hölder divergence. This method enhances brain tumor segmentation capabilities despite missing modalities by efficiently transferring knowledge and optimizing model generalization. Key components of our framework include: 1. Parallel U-Net network with single modality input to handle missing modalities. 2. Mutual information knowledge transfer to enhance model processing capabilities. 3. Optimization of knowledge transfer efficiency through Hölder divergence during knowledge distillation. Ablation experiments highlight the critical role of each component and reveal the impact of the Hölder conjugate exponent on model performance.

TABLE III
QUANTITATIVE EVALUATION OF SEGMENTATION RESULTS (DSC ↑) ON BraTS 2020. THIS TABLE PROVIDES A QUANTITATIVE ASSESSMENT OF THE SEGMENTATION PERFORMANCE ON THE BRATS 2020 DATASET, MEASURED USING THE DICE SIMILARITY COEFFICIENT (DSC). AN UPWARD ARROW (↑) INDICATES THAT HIGHER DSC VALUES CORRESPOND TO BETTER SEGMENTATION ACCURACY, ALLOWING A CLEAR COMPARISON OF THE EFFECTIVENESS OF DIFFERENT MODELS OR APPROACHES ON THIS DATASET.

| Task | Methods | Fl | T2 | T1c | T1 | T2,Fl | T1c,Fl | T1c,T2 | T1,Fl | T1,T2 | T1,T1c | $\sim$T1 | $\sim$T1c | $\sim$T2 | $\sim$Fl. | Full | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| WT | RFNet | 87.3 | 86.1 | 76.8 | 77.2 | 89.9 | 89.9 | 87.7 | 89.7 | 87.7 | 81.1 | 90.7 | 90.6 | 90.7 | 88.3 | 91.1 | 87.0 |
| | mmFormer | 86.5 | 85.5 | 78.0 | 76.2 | 89.4 | 89.3 | 87.5 | 88.7 | 86.9 | 80.7 | 90.4 | 89.8 | 89.7 | 87.6 | 90.5 | 86.4 |
| | M3AE | 86.5 | 86.1 | 73.9 | 76.7 | 89.3 | 89.5 | 87.4 | 89.4 | 87.2 | 78.1 | 90.2 | 90.4 | 90.0 | 88.6 | 90.6 | 86.3 |
| | MTI | 89.1 | 86.5 | 77.4 | 78.1 | 90.5 | 90.0 | 88.4 | 89.9 | 88.0 | 81.2 | 90.6 | 90.3 | 90.7 | 88.7 | 90.6 | 87.3 |
| | GGDM | 91.0 | 88.3 | 80.6 | 77.4 | 92.1 | 91.9 | 89.8 | 91.6 | 89.3 | 82.7 | 92.3 | 92.1 | 91.6 | 89.2 | 92.0 | 88.8 |
| | OUR | **91.9** | **89.2** | **81.5** | **80.5** | **92.5** | **92.4** | **90.0** | **92.1** | **89.8** | **83.2** | **92.8** | **92.6** | **92.3** | **89.9** | **92.7** | **89.6** |
| TC | RFNet | 69.2 | 71 | 81.5 | 66.0 | 74.1 | 84.7 | 83.5 | 73.1 | 73.1 | 83.4 | 85 | 75.2 | 85.1 | 83.5 | 85.2 | 78.2 |
| | mmFormer | 64.6 | 63.3 | 81.5 | 63.2 | 70.3 | 83.7 | 82.6 | 71.7 | 67.7 | 82.8 | 83.9 | 72.4 | 84.4 | 79.0 | 84.6 | 75.7 |
| | M3AE | 68.0 | 70.3 | 81.4 | 66.0 | 75.0 | 82.0 | 83.0 | 73.8 | 72.5 | 82.4 | 83.1 | 75.1 | 82.4 | 84.1 | 84.4 | 77.6 |
| | MTI | 69.3 | 71.5 | 83.4 | 66.8 | 75.5 | 85.6 | 86.4 | 73.9 | 73.3 | 85.2 | 86.4 | 75.9 | 86.5 | 86.5 | 87.4 | 79.6 |
| | GGDM | 77.0 | 79.1 | **87.6** | 70.5 | 81.6 | 88.0 | 88.5 | 80.3 | **81.1** | 88.0 | 88.1 | 82.4 | 87.9 | 88.4 | 88.0 | 83.8 |
| | OUR | **78.6** | **79.8** | 87.5 | **73.4** | 82.5 | **88.9** | **88.6** | **81.2** | 80.8 | **88.1** | **88.8** | 82.5 | **88.9** | **88.6** | **88.8** | **84.5** |
| ET | RFNet | 38.2 | 46.3 | 74.9 | 37.3 | 49.3 | 76.7 | 75.9 | 41.0 | 45.7 | 78.0 | 77.1 | 49.9 | 76.8 | 77.0 | 78.0 | 61.5 |
| | mmFormer | 36.6 | 49.0 | 78.3 | 37.6 | 49.0 | 79.4 | 77.2 | 42.9 | 49.1 | 81.7 | 78.7 | 50.0 | 80.6 | 68.3 | 79.9 | 62.6 |
| | M3AE | 40.5 | 46.0 | 72.4 | 39.9 | 47.3 | 74.7 | 76.8 | 43.2 | 46.6 | 75.4 | 77.1 | 48.2 | 75.9 | 77.4 | 78.0 | 61.3 |
| | MTI | 43.6 | 45.6 | 78.9 | 41.3 | 48.7 | 81.8 | 81.7 | 48.2 | 50.0 | 79.2 | 81.0 | 52.5 | 81.8 | 78.5 | 81.6 | 65.0 |
| | GGDM | 49.8 | 52.5 | 84.2 | 39.7 | 56.5 | 84.6 | 84.5 | 54.6 | **55.3** | 84.2 | 84.2 | **58.6** | 84.3 | **84.3** | 84.1 | 69.4 |
| | OUR | **51.9** | **54.6** | **84.6** | **44.3** | **57.7** | **84.9** | **84.8** | **55.4** | 55.1 | **84.3** | **84.9** | **58.6** | **84.4** | **84.3** | **84.3** | **70.3** |

TABLE IV
COMPARISON OF THE SUPERIORITY OF HÖLDER DIVERGENCE WITH DIFFERENT $f$-DIVERGENCES ON THE BRATS 2018 DATASET. THIS TABLE ILLUSTRATES THE COMPARATIVE PERFORMANCE OF HÖLDER DIVERGENCE AGAINST DIFFERENT TYPES OF $f$-DIVERGENCES ON THE BRATS 2018 DATASET.

| Methods | Dice | | | |
|---|---|---|---|---|
| $f$-divergence | WT | TC | ET | Avg. |
| **Total Variation** [30] | 67.2 | 1.9 | 0.9 | 23.3 |
| **Squared Hellinger** [30] | 85.3 | 75.4 | 60.1 | 73.6 |
| **Kullback-Leibler** [30] | 84.5 | 76.2 | 61.4 | 74.0 |
| **Neyman** $\chi^2$ [30] | 83.4 | 75.1 | 59.9 | 72.8 |
| **Jensen-Shannon** [30] | 84.6 | 76.5 | 59.8 | 73.6 |
| **Hölder** [17] | **88.2** | **82.6** | **69.4** | **80.1** |

TABLE V
EXPLORING THE IMPACT OF HÖLDER CONJUGATE EXPONENTS ON EXPERIMENTAL RESULTS BASED ON THE BRATS 2018 DATASET. THIS TABLE ILLUSTRATES HOW VARYING THE HÖLDER CONJUGATE EXPONENTS AFFECTS THE EXPERIMENTAL OUTCOMES DERIVED FROM THE BRATS 2018 DATASET. IT HIGHLIGHTS THE RELATIONSHIPS BETWEEN DIFFERENT CONJUGATE EXPONENT VALUES AND THEIR INFLUENCE ON THE PERFORMANCE METRICS WITHIN THE CONTEXT OF THE EXPERIMENTS.

| Methods | | Dice | | | |
|---|---|---|---|---|---|
| divergence | $\alpha$ | WT | TC | ET | Avg. |
| - | - | 87.8 | 82.9 | 67.4 | 79.4 |
| **KL** | - | 84.5 | 76.2 | 61.4 | 74.0 |
| **Hölder** | 1.05 | 87.8 | 82.6 | 69.2 | 79.9 |
| **Hölder** | 1.08 | 88.2 | 82.6 | 68.7 | 79.9 |
| **Hölder** | **1.10** | **88.2** | 82.6 | **69.4** | **80.1** |
| **Hölder** | 1.15 | 88.1 | 82.8 | 67.9 | 79.6 |
| **Hölder** | 1.20 | 87.9 | **83.1** | 68.1 | 79.7 |

TABLE VI
QUANTITATIVE EVALUATION RESULTS OF THE ABLATION STUDY ON THE BRATS 2018 DATASET. THIS TABLE PRESENTS THE IMPACT OF DIFFERENT MODEL COMPONENTS ON PERFORMANCE, HIGHLIGHTING THE EFFECTIVENESS OF EACH COMPONENT IN CONTRIBUTING TO OVERALL ACCURACY AND ROBUSTNESS.

| Methods | | | Number of Missing Modalities | | | | Avg. |
|---|---|---|---|---|---|---|---|
| $\mathcal{L}_{dice}$ | $\mathcal{L}_{MI}$ | $\mathcal{L}_{HD}$ | 3 | 2 | 1 | 0 | |
| ✓ | ✗ | ✗ | 60.2 | 71.9 | 77.1 | 80.5 | 70.7 |
| ✓ | ✓ | ✗ | 72.4 | 79.7 | 83.9 | 87.0 | 79.4 |
| ✓ | ✗ | ✓ | 72.9 | 80.1 | 84.3 | 87.4 | 79.8 |
| ✓ | ✓ | ✓ | **73.6** | **80.3** | **84.4** | **87.6** | **80.1** |

Despite its effectiveness, the method has some limitations: 1. Training costs due to the large number of parameters and extensive tuning requirements, and 2. Sensitivity to hyperparameter selection, necessitating extensive experimentation and validation.

REFERENCES

[1] F. Isensee, P. F. Jaeger, S. A. Kohl, J. Petersen, and K. H. Maier-Hein, "nnU-Net: a self-configuring method for deep learning-based biomedical image segmentation," *Nature Methods*, vol. 18, no. 2, pp. 203–211, 2021.

[2] D. Zhang, G. Huang, Q. Zhang, J. Han, J. Han, Y. Wang, and Y. Yu, "Exploring task structure for brain tumor segmentation from multi-modality MR images," *IEEE Transactions on Image Processing*, vol. 29, pp. 9032–9043, 2020.

[3] Y. Tang, D. Yang, W. Li, H. R. Roth, B. Landman, D. Xu, V. Nath, and A. Hatamizadeh, "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. 20 730–20 740.

[4] S. Wang, C. Li, R. Wang, Z. Liu, M. Wang, H. Tan, Y. Wu, X. Liu, H. Sun, R. Yang *et al.*, "Annotation-efficient deep learning for automatic medical image segmentation," *Nature Communications*, vol. 12, no. 1, p. 5915, 2021.

[5] Q. Chen, J. Zhang, R. Meng, L. Zhou, Z. Li, Q. Feng, and D. Shen, "Modality-Specific Information Disentanglement From Multi-Parametric MRI for Breast Tumor Segmentation and Computer-Aided Diagnosis," *IEEE Transactions on Medical Imaging*, vol. 43, no. 5, pp. 1958–1971, 2024.

[6] Y. Xie, J. Zhang, Y. Xia, and C. Shen, "Learning From Partially Labeled Data for Multi-Organ and Tumor Segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 12, pp. 14 905–14 919, 2023.

[7] D. Shah, A. Barve, B. Vala, and J. Gandhi, "A Survey on Brain Tumor Segmentation with Missing MRI Modalities," in *International Conference on Information Technology*. Springer, 2023, pp. 299–308.

[8] H. Liu, D. Wei, D. Lu, J. Sun, L. Wang, and Y. Zheng, "M3AE: Multimodal representation learning for brain tumor segmentation with missing modalities," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 2, 2023, pp. 1657–1665.

[9] R. Azad, N. Khosravi, and D. Merhof, "SMU-Net: Style matching U-Net for brain tumor segmentation with missing modalities," in *International Conference on Medical Imaging with Deep Learning, MIDL 2022, 6-8 July 2022, Zurich, Switzerland*, ser. Proceedings of Machine Learning Research, vol. 172. PMLR, 2022, pp. 48–62.

[10] S. Wang, Z. Yan, D. Zhang, H. Wei, Z. Li, and R. Li, "Prototype Knowledge Distillation for Medical Segmentation with Missing Modality," in *IEEE International Conference on Acoustics, Speech and Signal Processing ICASSP 2023, Rhodes Island, Greece, June 4-10, 2023*. IEEE, 2023, pp. 1–5.

[11] S. Karimijafarbigloo, R. Azad, A. Kazerouni, S. Ebadollahi, and D. Merhof, "Mmcformer: Missing modality compensation transformer for brain tumor segmentation," in *Medical Imaging with Deep Learning*. PMLR, 2024, pp. 1144–1162.

[12] Y. Ding, X. Yu, and Y. 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.

[13] H. Ting and M. Liu, "Multimodal Transformer of Incomplete MRI Data for Brain Tumor Segmentation," *IEEE Journal of Biomedical and Health Informatics*, vol. 28, no. 1, pp. 89–99, 2024.

[14] H. Wang, Y. Chen, C. Ma, J. Avery, L. Hull, and G. Carneiro, "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. 15 878–15 887.

[15] W.-G. Chang, T. You, S. Seo, S. Kwak, and B. Han, "Domain-specific batch normalization for unsupervised domain adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7354–7362.

[16] S. Ahn, S. X. Hu, A. Damianou, N. D. Lawrence, and Z. Dai, "Variational information distillation for knowledge transfer," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9163–9171.

[17] F. Nielsen, K. Sun, and S. Marchand-Maillet, "On Hölder projective divergences," *Entropy*, vol. 19, no. 3, p. 122, 2017.

[18] B. H. Menze, A. Jakab, S. Bauer, J. Kalpathy-Cramer, K. Farahani, J. Kirby, Y. Burren, N. Porz, J. Slotboom, R. Wiest *et al.*, "The multimodal brain tumor image segmentation benchmark (BRATS)," *IEEE Transactions on Medical Imaging*, vol. 34, no. 10, pp. 1993–2024, 2014.

[19] D. Barber and F. Agakov, "The IM algorithm: a variational approach to information maximization," *Advances in Neural Information Processing Systems*, vol. 16, no. 320, p. 201, 2004.

[20] F. Milletari, N. Navab, and S.-A. Ahmadi, "V-net: Fully convolutional neural networks for volumetric medical image segmentation," in *2016 Fourth International Conference on 3D Vision (3DV)*. IEEE, 2016, pp. 565–571.

[21] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.

[22] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[23] J. Chen and A. L. Martel, "Head and neck tumor segmentation with 3D UNet and survival prediction with multiple instance neural network," in *3D Head and Neck Tumor Segmentation in PET/CT Challenge*. Springer, 2022, pp. 221–229.

[24] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015.

[25] Y. Ding, X. Yu, and Y. 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.

[26] Y. Zhang, N. He, J. Yang, Y. Li, D. Wei, Y. Huang, Y. Zhang, Z. He, and Y. Zheng, "mmformer: Multimodal medical transformer for incomplete multimodal learning of brain tumor segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2022, pp. 107–117.

[27] H. Liu, D. Wei, D. Lu, J. Sun, L. Wang, and Y. Zheng, "M3AE: Multimodal representation learning for brain tumor segmentation with missing modalities," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 2, 2023, pp. 1657–1665.

[28] H. Ting and M. Liu, "Multimodal Transformer of Incomplete MRI Data for Brain Tumor Segmentation," *IEEE Journal of Biomedical and Health Informatics*, vol. 28, no. 1, pp. 89–99, 2024.

[29] H. Wang, S. Luo, G. Hu, and J. Zhang, "Gradient-Guided Modality Decoupling for Missing-Modality Robustness," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 14, 2024, pp. 15 483–15 491.

[30] F. Nielsen and K. Okamura, "On the f-divergences between densities of a multivariate location or scale family," *Statistics and Computing*, vol. 34, no. 1, p. 60, 2024.