

# Uni-Encoder Meets Multi-Encoders: Representation Before Fusion for Brain Tumor Segmentation with Missing Modalities

## Supplementary Material

This supplementary material details data preprocessing and inference settings, reports additional metrics on BraTS 2023 and BraTS 2024, presents extended ablation studies on crop and patch sizes, analyzes model complexity, and discusses limitations of UniME.

### A. Data Preprocessing

We follow standard BraTS preprocessing practices used in recent work [16, 24, 45, 68]. We load the four MRI modalities (FLAIR, T1ce, T1, and T2) and the associated segmentation labels from the original BraTS dataset. We do not resample the voxel grid or alter spatial coordinates. We compute a brain mask as the union of non-zero voxels across modalities and apply  $z$ -score normalization per modality and per case within this mask.

For comparability with previous work [36, 67, 68] that trains with  $128 \times 128 \times 128$  crop size, and because a few volumes have a dimension less than 128, we apply symmetric padding along the affected axis to 128 voxels. The impact of this padding is empirically negligible. For the BraTS 2023 dataset, only one case has a dimension of 127 voxels, so padding does not alter gross tumor morphology or bias evaluation.

BraTS 2024 corresponds to a post-operative glioma segmentation setting and introduces an additional annotation for the resection cavity, which is not present in the standard pre-operative BraTS datasets. To ensure compatibility with prior incomplete-modality segmentation methods and maintain consistent evaluation across BraTS 2023 and BraTS 2024, the resection cavity is merged into the enhancing tumor class during preprocessing. Then, the evaluation regions therefore follow the standard BraTS definition: the whole tumor (WT) includes all abnormal regions, the tumor core (TC) includes necrotic or non-enhancing tumor together with enhancing tumor, and the enhancing tumor (ET) region corresponds to the merged enhancing tumor category. Under this protocol, the enhancing tumor region in BraTS 2024 effectively represents the union of enhancing tumor and resection cavity.

### B. Inference Details

At inference, we use a 3D sliding-window strategy with overlap. For each model, the window size equals the training crop size, and the stride is half the window length, yielding 50% overlap. Per-window logits

are passed through softmax and accumulated as spatial probability maps. In overlapping regions, probabilities are averaged across overlapping windows, and the final voxel label is obtained via argmax.

To suppress false positives in the enhancing tumor region, we apply a standard post-processing rule [16, 24, 27, 45, 68]. If the predicted ET volume contains fewer than 500 voxels, we set the corresponding ET mask to zero. For completeness, Tab. 8 and Tab. 9 report ET segmentation performance on BraTS 2023 and BraTS 2024 before and after this post-processing.

### C. Extended Comparisons

In addition to the DSC comparisons in Tab.1 and Tab.2, we also evaluate boundary quality using the 95th percentile Hausdorff Distance (HD95), which provides a robust measure of boundary accuracy. Specifically, for two 3D segmentation surface point sets  $\partial A$  and  $\partial B$ , the directed 95th percentile Hausdorff distance is defined as:

$$h_{95}(A, B) = \mathcal{P}_{95\%} \left\{ \min_{b \in \partial B} \|a - b\|_2 : a \in \partial A \right\}, \quad (8)$$

where  $\mathcal{P}_{95\%}$  denotes the 95th percentile and  $\|\cdot\|_2$  the  $L_2$  norm. The HD95 between  $A$  and  $B$  is defined as:

$$\text{HD95} = \max(h_{95}(A, B), h_{95}(B, A)). \quad (9)$$

We compute HD95 for the WT, TC, and ET using the DisTorch implementation [42].

Tab. 10 and Tab. 12 report HD95 results on BraTS 2023 and BraTS 2024, respectively. Across both datasets, UniME consistently outperforms competing methods on WT and ET in HD95. Performance on TC is very close to the best-performing method. Overall, UniME accurately characterizes lesion boundaries, as reflected by HD95.

### D. Extended Ablation Studies

**Effect of Crop Size.** Transformers capture long-range dependencies, and previous work in general medical image segmentation indicates that increasing the input resolution can improve segmentation performance [9]. Therefore, we conduct ablations on crop size at 80, 96, and 128. Tab. 13 reports the detailed results.

As shown in Tab. 13, a crop size of 96 yields the best overall performance, balancing context coverage and computational efficiency. Larger crops improve WT accuracy, indicating that UniME benefits from additional

Type	Flair T1 T1ce T2	○	○	○	●	○	○	●	○	●	●	●	●	○	●	Avg.	
		○	○	●	○	○	●	●	○	○	○	●	●	○	●		
ET	HeMIS	4.96	30.62	5.91	18.92	57.34	46.31	18.09	7.63	21.47	65.83	68.14	18.48	68.17	60.23	67.41	37.30
	U-HEVD	35.96	71.10	22.41	36.63	75.89	73.97	41.62	39.12	45.07	76.27	77.08	45.84	77.15	76.75	77.49	58.16
	RFNet	46.89	79.36	38.73	40.97	80.56	80.82	48.62	50.45	55.10	81.10	82.20	55.96	81.27	81.08	81.59	65.65
	mmFormer	46.19	78.64	38.70	41.46	80.75	80.84	47.01	48.96	51.54	81.35	81.82	53.64	81.28	81.23	81.48	64.99
	M <sup>3</sup> AE	58.68	82.61	54.02	57.17	83.93	83.13	60.73	60.35	63.09	84.26	84.63	63.52	84.56	83.65	84.62	72.60
	M <sup>2</sup> FTrans	53.01	82.12	44.80	48.49	83.06	83.11	54.22	55.29	58.48	83.28	83.56	58.88	83.33	83.28	83.40	69.22
	IM-Fuse	50.57	81.30	48.19	51.15	82.00	82.31	55.00	54.78	57.18	82.51	82.95	58.61	82.25	82.41	82.50	68.91
	M <sup>2</sup> SegMamba	54.06	81.32	50.29	54.44	82.95	83.09	59.13	57.92	60.28	82.95	83.48	62.13	83.07	83.40	83.44	70.80
	Ours	<b>59.16</b>	<b>84.15</b>	<b>54.24</b>	<b>58.69</b>	<b>85.55</b>	<b>84.29</b>	<b>62.16</b>	<b>61.07</b>	<b>63.83</b>	<b>84.71</b>	<b>84.89</b>	<b>64.90</b>	<b>85.40</b>	<b>85.41</b>	<b>85.51</b>	<b>73.60</b>
ET*	HeMIS	4.96	30.62	5.91	18.92	57.34	46.31	18.89	7.60	21.46	65.83	68.54	18.76	68.56	60.23	67.74	37.44
	U-HEVD	35.95	71.10	22.41	36.63	75.89	73.96	42.02	38.98	45.07	76.27	77.48	46.24	77.50	77.11	78.67	58.35
	RFNet	48.87	80.51	39.40	41.63	81.76	81.62	49.22	51.44	55.80	82.70	84.20	57.05	82.47	83.08	83.59	66.89
	mmFormer	46.59	79.84	39.05	41.98	82.35	82.79	47.09	49.67	52.50	82.90	83.35	54.39	82.88	83.23	83.08	66.11
	M <sup>3</sup> AE	58.68	83.95	54.00	57.09	84.79	85.12	60.90	60.25	63.03	84.90	85.80	63.50	85.43	85.23	85.74	73.23
	M <sup>2</sup> FTrans	54.58	84.04	45.79	49.81	85.03	85.51	55.34	56.60	60.18	84.87	85.56	60.38	84.90	85.64	85.78	70.93
	IM-Fuse	51.74	82.85	48.80	51.74	84.40	84.30	56.04	56.68	58.71	84.51	85.34	60.94	84.25	84.81	84.90	70.67
	M <sup>2</sup> SegMamba	54.84	81.72	51.09	55.24	84.55	83.49	59.48	58.69	60.59	84.53	84.67	62.37	84.27	85.40	85.03	71.73
	Ours	<b>61.12</b>	<b>86.95</b>	<b>55.35</b>	<b>60.00</b>	<b>87.54</b>	<b>87.40</b>	<b>63.15</b>	<b>62.11</b>	<b>65.11</b>	<b>87.59</b>	<b>87.94</b>	<b>65.93</b>	<b>87.80</b>	<b>87.87</b>	<b>87.98</b>	<b>75.59</b>

Table 8. Performance comparison of ET region segmentation (DSC %) on BraTS 2023. Available and missing modalities are denoted by ● and ○, respectively. ET represents the enhancing tumor region without post-processing, while ET\* indicates results after applying post-processing. ■, ■, and ■ indicate the first-, second-, and third-best performance, respectively. WT, TC, and ET represent whole tumor, tumor core, and enhancing tumor, respectively.

Type	Flair T1 T1ce T2	○	○	○	●	○	○	●	○	●	●	●	●	○	●	Avg.	
		○	○	●	○	○	●	●	○	○	○	○	○	○	○		
ET	HeMIS	15.76	26.47	8.68	1.12	41.52	34.96	15.65	23.69	26.04	40.02	40.82	24.03	46.25	41.10	44.30	28.69
	U-HEVD	39.93	54.21	47.12	34.97	63.32	62.44	49.10	51.77	46.12	66.85	69.19	51.66	69.06	67.13	70.42	56.22
	RFNet	56.49	71.94	55.26	51.93	74.42	74.59	60.02	59.88	60.31	75.91	76.90	61.92	76.90	75.59	77.52	67.30
	mmFormer	57.57	72.00	55.28	52.76	75.77	75.10	61.45	62.66	62.80	76.61	78.14	64.35	78.07	77.59	79.30	68.63
	M <sup>3</sup> AE	63.44	73.61	61.44	60.03	78.32	76.88	64.86	65.43	66.74	77.87	80.41	67.24	80.42	78.70	80.86	71.75
	M <sup>2</sup> FTrans	56.97	73.48	58.11	55.70	76.89	76.80	62.07	62.03	62.15	77.78	79.09	64.39	78.87	77.99	79.45	69.45
	IM-Fuse	60.71	74.14	59.02	55.12	77.16	76.97	64.75	64.24	65.17	78.61	80.14	67.72	80.35	77.89	80.44	70.83
	M <sup>2</sup> SegMamba	59.61	74.67	60.13	57.43	78.11	77.27	64.40	64.27	65.54	78.90	80.45	66.98	79.50	79.03	81.20	71.17
	Ours	<b>64.75</b>	<b>77.56</b>	<b>65.12</b>	<b>63.15</b>	<b>80.75</b>	<b>79.87</b>	<b>69.30</b>	<b>69.47</b>	<b>69.73</b>	<b>81.11</b>	<b>81.90</b>	<b>71.55</b>	<b>83.04</b>	<b>81.19</b>	<b>83.17</b>	<b>74.78</b>
ET*	HeMIS	16.08	26.47	8.62	2.15	41.52	34.96	15.59	24.06	26.38	40.02	41.51	24.71	46.57	41.10	44.98	28.98
	U-HEVD	39.93	54.21	47.12	34.97	63.24	62.44	49.81	51.74	46.06	67.22	69.93	52.31	69.06	67.12	70.79	56.40
	RFNet	57.23	73.53	57.04	53.70	75.06	76.25	61.37	60.60	61.18	77.14	78.59	63.73	77.67	76.20	78.75	68.54
	mmFormer	58.54	73.29	56.63	53.53	77.49	76.55	61.92	63.25	62.57	77.17	78.92	65.42	79.18	78.64	80.24	69.56
	M <sup>3</sup> AE	63.55	73.98	61.92	60.08	79.04	77.53	64.91	65.46	66.89	78.30	80.81	67.34	81.03	79.82	82.12	72.19
	M <sup>2</sup> FTrans	58.47	75.19	58.51	57.33	78.23	77.71	63.51	63.16	63.76	79.48	80.83	65.97	80.36	79.31	81.25	70.87
	IM-Fuse	61.36	76.02	61.01	55.87	77.99	78.33	64.59	64.87	65.06	79.28	80.43	67.08	80.01	78.94	80.77	71.44
	M <sup>2</sup> SegMamba	60.66	75.19	60.81	58.50	78.51	77.79	65.04	64.53	65.59	79.24	80.56	67.66	80.32	79.44	81.33	71.68
	Ours	<b>66.00</b>	<b>78.57</b>	<b>66.04</b>	<b>64.40</b>	<b>80.99</b>	<b>80.34</b>	<b>69.58</b>	<b>69.43</b>	<b>69.73</b>	<b>81.75</b>	<b>82.41</b>	<b>70.89</b>	<b>82.30</b>	<b>81.58</b>	<b>82.81</b>	<b>75.12</b>

Table 9. Performance comparison of ET region segmentation (DSC %) on BraTS 2024. Available and missing modalities are denoted by ● and ○, respectively. ET represents the enhancing tumor region without post-processing, while ET\* indicates results after applying post-processing. ■, ■, and ■ indicate the first-, second-, and third-best performance, respectively. WT, TC, and ET represent whole tumor, tumor core, and enhancing tumor, respectively.

global context. However, larger crops degrade performance on smaller, more intricate substructures such as TC and ET. Specifically, TC and ET peak at a crop size of 96, while larger crops reduce accuracy.

**Effect of Patch Size.** Patch size governs the Uni-Encoder’s ability to model fine structures. Reducing the patch size significantly increases the sequence length, enabling the Uni-Encoder to model finer structures, but simultaneously leads to rapidly rising computational complexity. We conduct ablations on patch size at 8,

12, and 16. Tab. 14 reports the results.

A patch size of 8 achieves the best overall performance, indicating that finer tokenization better captures the fine-grained structure of the brain tumor. Increasing the patch size to 12 or 16 degrades accuracy across all regions, as larger patches lose critical spatial detail. Therefore, we select the smallest patch size feasible under our current memory budget. If the patch size were reduced to 4, the token sequence length would approach 14k tokens, significantly increasing computational over-

Type	Flair T1 T1ce T2	○	○	○	●	○	○	●	○	●	●	●	●	○	●	Avg.	
		○	○	●	○	○	●	●	○	○	○	○	○	○	○		
WT	HeMIS	81.36	75.92	73.24	71.03	72.54	69.19	63.02	65.25	68.50	65.00	56.92	53.89	59.12	59.66	49.03	65.58
	U-HEVD	48.40	39.78	52.47	42.74	28.02	34.52	22.85	27.84	27.30	25.11	16.90	17.82	18.42	19.56	13.73	29.03
	RFNet	31.14	27.21	24.27	22.79	23.29	22.82	14.66	21.27	20.20	18.23	15.00	16.58	18.51	21.90	16.72	20.97
	mmFormer	19.61	28.91	20.74	21.05	18.39	18.45	12.75	16.80	12.36	11.62	12.78	11.75	13.85	20.46	16.20	17.05
	M <sup>3</sup> AE	8.62	12.09	12.34	8.82	7.50	11.07	6.50	7.85	6.27	5.64	5.85	6.55	5.75	7.27	6.21	7.89
	M <sup>2</sup> FTrans	14.08	14.69	18.20	11.74	10.49	12.81	8.30	9.77	8.84	9.55	10.14	8.27	10.05	10.36	10.18	11.16
	LS3M	13.46	13.10	13.39	11.96	9.90	11.44	7.79	9.61	8.31	9.51	7.23	7.65	9.50	9.19	8.47	10.03
	IM-Fuse	21.11	13.63	14.30	9.64	10.36	12.20	8.22	12.03	11.34	6.53	7.48	9.40	8.77	9.97	8.90	10.93
	M <sup>2</sup> SegMamba Ours	13.25	12.33	13.59	7.63	8.93	9.53	5.90	8.21	6.89	6.40	6.04	6.14	6.76	8.79	6.36	8.45
		<b>7.40</b>	<b>9.57</b>	<b>9.28</b>	<b>5.75</b>	<b>6.26</b>	<b>8.28</b>	<b>5.30</b>	<b>6.14</b>	<b>5.37</b>	<b>4.83</b>	<b>5.34</b>	<b>5.15</b>	<b>4.87</b>	<b>5.88</b>	<b>5.13</b>	<b>6.30</b>
TC	HeMIS	98.31	89.32	86.38	88.90	87.78	82.19	78.90	81.77	85.02	81.12	70.33	71.74	74.80	76.24	59.74	80.84
	U-HEVD	72.26	57.41	68.95	62.95	48.52	52.11	36.71	37.04	42.19	45.15	28.82	26.35	35.46	26.82	20.87	44.11
	RFNet	56.41	46.73	43.31	39.22	46.36	42.76	24.39	42.11	38.92	31.91	28.13	32.06	36.98	45.90	35.02	39.35
	mmFormer	35.48	44.03	31.20	50.53	38.67	29.58	29.93	37.29	29.81	22.04	26.41	29.75	31.16	41.28	34.37	34.10
	M <sup>3</sup> AE	<b>10.63</b>	<b>6.38</b>	11.13	<b>10.08</b>	5.66	5.76	9.39	9.94	<b>9.55</b>	<b>4.98</b>	<b>4.85</b>	9.39	5.48	<b>5.52</b>	5.25	<b>7.60</b>
	M <sup>2</sup> FTrans	23.65	22.66	32.15	31.46	17.75	24.18	27.10	24.49	21.29	25.61	27.36	26.95	25.37	23.46	28.40	25.46
	LS3M	25.00	12.59	16.33	20.65	19.38	13.27	12.10	18.08	15.61	17.61	12.42	12.10	20.66	19.44	17.51	16.85
	IM-Fuse	49.18	17.18	19.30	21.37	20.78	12.73	16.42	29.53	29.66	12.84	13.07	22.70	19.92	18.67	17.67	21.40
	M <sup>2</sup> SegMamba Ours	31.94	11.06	26.43	16.71	20.44	10.74	11.89	18.70	17.29	13.44	10.82	12.32	15.93	16.42	10.91	16.34
		13.06	8.96	<b>10.18</b>	13.85	<b>5.29</b>	<b>5.51</b>	<b>8.97</b>	<b>9.88</b>	9.94	4.99	5.35	<b>7.04</b>	<b>4.48</b>	5.65	<b>4.36</b>	7.83
ET	HeMIS	98.93	89.82	86.73	89.58	89.28	82.50	79.26	82.64	85.13	81.07	68.63	70.17	74.03	77.09	57.81	80.84
	U-HEVD	66.63	51.51	66.72	56.39	39.35	44.72	30.40	30.17	35.30	36.84	22.95	20.39	27.13	20.17	15.16	37.59
	RFNet	45.19	37.43	31.76	30.66	38.64	33.09	17.13	30.42	30.33	24.30	20.99	22.42	28.81	37.50	27.79	30.43
	mmFormer	26.45	36.05	20.69	40.66	29.88	22.64	20.53	24.98	21.31	17.35	19.05	20.80	23.82	31.99	27.51	25.58
	M <sup>3</sup> AE	<b>9.31</b>	<b>4.72</b>	10.27	<b>8.95</b>	4.41	4.46	8.58	8.83	8.40	3.80	3.67	8.34	3.98	<b>4.35</b>	3.77	6.39
	M <sup>2</sup> FTrans	16.26	17.33	21.89	26.26	14.41	19.03	19.38	14.89	14.43	20.33	21.38	17.04	20.73	17.44	21.74	18.84
	LS3M	17.44	8.83	11.42	15.37	16.24	10.30	8.76	11.01	11.98	13.49	8.68	8.75	15.68	14.08	12.10	12.28
	IM-Fuse	39.33	13.81	13.30	17.38	16.90	10.43	12.19	21.97	23.77	10.74	10.24	17.11	15.48	15.31	14.76	16.85
	M <sup>2</sup> SegMamba Ours	26.63	8.46	19.73	14.27	16.27	8.40	10.58	15.36	14.90	11.17	9.20	11.78	12.93	13.08	9.53	13.49
		9.34	6.02	<b>7.86</b>	10.56	<b>4.24</b>	<b>4.38</b>	<b>7.40</b>	<b>7.57</b>	<b>7.21</b>	<b>3.51</b>	<b>3.52</b>	<b>6.17</b>	<b>3.73</b>	4.53	<b>3.58</b>	<b>5.97</b>

Table 10. **Performance comparison (HD95 mm) on BraTS 2023.** Lower is better. Available and missing modalities are denoted by ● and ○, respectively.  ,  , and   indicate the first-, second-, and third-best performance, respectively. WT, TC, and ET represent whole tumor, tumor core, and enhancing tumor, respectively.

Setting	Configuration					Params (M)			FLOPs (GMacs)			Memory (GiB)				Average Region DSC (%)			Overall (%)
	Crop Size	Patch Size	Heads	L	$\alpha_{embed}$	Trans	CNN	Total	Trans	CNN	Total	Trans	CNN	Other	Total	WT	TC	ET	Mean
Scale	96	8	12	12	864	111.01	11.19	122.20	252.30	208.28	460.61	1.57	2.43	0.29	4.29	89.50	80.77	69.72	80.00
	96	8	12	16	864	146.92	11.19	158.11	335.37	208.28	543.69	2.08	2.43	0.29	4.80	90.38	<b>84.51</b>	<b>75.59</b>	<b>83.49</b>
	96	8	16	24	1056	325.75	11.85	337.61	714.21	209.42	923.68	4.01	2.44	0.29	6.74	<b>90.50</b>	84.09	75.08	83.22
Crop Size	80	8	12	16	864	146.29	11.19	157.48	174.06	120.53	294.61	1.43	1.43	0.17	3.03	90.17	84.22	74.54	82.97
	96	8	12	16	864	146.92	11.19	158.11	335.37	208.28	543.69	2.08	2.43	0.29	4.80	90.38	<b>84.51</b>	<b>75.59</b>	<b>83.49</b>
	128	8	12	16	864	148.97	11.19	160.16	1064.20	493.70	1557.99	4.18	5.71	0.69	10.58	<b>90.87</b>	84.25	74.64	83.25
Patch Size	96	8	12	16	864	146.92	11.19	158.11	335.37	208.28	543.69	2.08	2.43	0.29	4.80	<b>90.38</b>	<b>84.51</b>	<b>75.59</b>	<b>83.49</b>
	96	12	12	16	864	150.07	11.19	161.27	84.59	204.11	288.74	1.02	2.43	0.29	3.74	89.99	83.89	74.33	82.74
	96	16	12	16	864	158.00	11.19	169.19	36.01	203.09	239.14	0.78	2.43	0.29	3.50	89.71	83.21	73.74	82.22

Table 11. **Model complexity with different configurations.**  ,  , and   indicate the first, second, and third best results, respectively. **Params** refers to the total number of parameters, measured in millions (M). **FLOPs** indicates floating-point operations, measured in giga multiply-accumulate operations (GMacs). **Memory** usage, measured in gigabytes (GiB), represents the estimated memory consumed by model parameters, buffers, and activations during a single forward pass. Average Region DSC and Overall DSC are reported based on evaluations on the BraTS 2023 dataset.

head. Since changing crop size and patch size substantially increases GPU memory usage and computation, we provide a detailed complexity analysis below.

## E. Model Complexity

As UniME models cross-modal information from high-resolution inputs with a single ViT encoder, one might expect longer sequences and the attendant  $\mathcal{O}(n^2)$  memory cost. In practice, the token sequence length remains modest. The Uni-Encoder tokenizes the channel-stacked volume into a sequence whose length is determined by

the crop and patch sizes. For patch size 8, crop size 96, and 4 register tokens, the sequence length is 1732. By contrast, prior designs [38, 67, 68] process features from four modality-specific CNNs concatenated after  $4\times$  downsampling, yielding  $8\times 8\times 8$  grids per modality; the fused sequence length is 2048. Consequently, UniME’s effective input sequence length is comparable to existing approaches.

Tab. 15 compares UniME’s efficiency against some existing methods [16, 38, 67, 68]. We reproduce the reported hyperparameters and training settings from the

Type	Flair T1 T1ce T2	○	○	○	●	○	○	●	○	●	●	●	●	○	●	Avg.	
		○	○	●	○	○	●	●	○	○	○	○	○	○	○		
WT	HeMIS	63.79	66.71	63.09	55.90	63.31	63.99	57.63	60.44	56.19	59.15	53.67	51.97	52.88	58.26	49.79	58.45
	U-HEVD	71.87	71.24	47.93	46.59	60.43	50.65	33.42	44.79	41.75	47.63	39.89	31.48	41.51	44.79	35.42	47.29
	RFNet	26.03	23.50	21.29	22.30	26.05	20.82	16.19	21.27	17.19	17.41	16.53	17.84	18.76	21.76	19.13	20.41
	mmFormer	14.43	17.47	12.12	12.86	14.21	12.84	9.55	10.49	11.34	12.87	11.28	9.51	12.79	11.29	11.40	12.30
	M <sup>3</sup> AE	10.78	13.18	13.24	8.83	10.00	12.03	9.12	10.34	8.39	7.95	8.35	8.23	7.51	9.73	8.26	9.73
	M <sup>2</sup> FTrans	17.12	14.41	12.75	11.87	11.42	12.62	10.48	11.02	11.47	9.93	9.64	9.83	9.87	10.67	9.59	11.51
	LS3M	12.91	16.19	13.40	11.68	13.01	12.66	9.69	10.68	9.57	10.34	8.92	9.18	9.71	11.30	9.38	11.24
	IM-Fuse	15.95	14.85	12.68	11.13	11.07	11.99	8.71	10.01	8.60	8.47	8.69	8.29	7.68	10.02	8.64	10.45
	M <sup>2</sup> SegMamba Ours	11.62	13.03	12.50	10.52	9.25	11.02	8.68	10.62	9.18	8.79	8.29	9.21	8.51	8.97	7.66	9.86
		<b>8.04</b>	<b>9.53</b>	<b>9.47</b>	<b>7.66</b>	<b>7.51</b>	<b>8.90</b>	<b>5.73</b>	<b>7.08</b>	<b>5.82</b>	<b>6.43</b>	<b>6.12</b>	<b>5.27</b>	<b>5.07</b>	<b>7.71</b>	<b>5.29</b>	<b>7.04</b>
TC	HeMIS	79.77	87.53	85.83	78.90	79.94	86.65	81.50	76.02	71.32	78.52	74.79	70.11	66.15	77.27	65.36	77.31
	U-HEVD	97.01	99.31	78.18	77.27	86.67	80.35	57.57	70.33	69.32	74.48	64.60	54.94	66.59	71.10	59.78	73.83
	RFNet	55.27	36.69	29.48	39.12	50.36	30.15	30.32	39.22	37.16	32.01	29.33	34.15	38.03	39.28	36.49	37.14
	mmFormer	21.84	19.20	14.66	22.61	23.16	15.49	15.03	13.79	19.25	19.07	16.91	16.02	22.51	17.36	19.41	18.42
	M <sup>3</sup> AE	12.75	<b>9.27</b>	<b>14.18</b>	<b>14.19</b>	<b>6.90</b>	<b>9.34</b>	<b>13.63</b>	<b>12.40</b>	<b>12.22</b>	<b>7.96</b>	<b>6.79</b>	<b>11.70</b>	<b>7.04</b>	<b>6.70</b>	<b>6.49</b>	<b>10.11</b>
	M <sup>2</sup> FTrans	25.75	14.33	14.80	22.10	12.23	13.50	17.60	14.07	18.63	12.10	11.06	15.11	11.52	9.68	10.75	14.88
	LS3M	16.85	19.64	18.12	19.56	16.71	14.60	15.56	14.59	13.10	12.51	10.72	12.98	10.98	13.58	11.65	14.74
	IM-Fuse	33.49	19.47	18.26	22.81	15.23	12.41	14.63	17.08	14.59	9.39	8.99	13.27	10.34	12.34	8.73	15.40
	M <sup>2</sup> SegMamba Ours	16.80	12.89	14.24	17.89	10.07	<b>9.11</b>	<b>13.58</b>	<b>13.12</b>	<b>13.60</b>	<b>9.68</b>	<b>9.54</b>	<b>13.35</b>	<b>9.98</b>	<b>8.39</b>	<b>9.00</b>	<b>12.08</b>
		<b>12.53</b>	15.01	14.46	18.13	8.15	<b>9.41</b>	<b>11.90</b>	<b>11.71</b>	<b>11.99</b>	9.95	7.04	<b>10.70</b>	<b>6.57</b>	7.70	6.54	<b>10.79</b>
ET	HeMIS	80.05	87.34	85.75	77.74	79.23	86.29	81.95	75.56	70.77	77.76	73.31	69.64	64.81	76.60	63.78	76.70
	U-HEVD	96.37	98.56	77.26	76.82	84.14	79.10	53.57	68.30	68.65	71.29	61.53	52.10	63.89	68.43	56.61	71.78
	RFNet	49.65	32.06	25.88	36.18	45.81	25.92	26.04	34.95	32.60	28.11	25.23	30.17	32.83	35.74	31.77	32.86
	mmFormer	18.30	16.50	12.20	16.93	19.97	12.54	12.52	11.90	16.27	15.76	13.52	13.44	19.67	14.11	16.28	15.33
	M <sup>3</sup> AE	11.66	<b>8.39</b>	13.19	<b>12.90</b>	<b>6.97</b>	<b>8.22</b>	13.27	11.74	11.36	7.66	6.42	11.26	6.74	6.74	5.95	9.50
	M <sup>2</sup> FTrans	17.54	9.82	12.87	14.70	9.16	10.69	11.19	10.78	12.70	7.74	7.48	9.95	7.97	7.76	7.26	10.51
	LS3M	14.09	14.63	15.52	17.13	13.08	10.36	12.99	12.92	11.62	10.13	8.80	11.82	8.89	11.14	9.85	12.20
	IM-Fuse	27.45	15.20	14.13	19.07	12.34	10.20	13.60	14.01	12.55	<b>7.40</b>	7.32	11.58	8.72	10.61	7.42	12.77
	M <sup>2</sup> SegMamba Ours	13.42	10.55	13.26	15.72	<b>8.19</b>	<b>8.09</b>	12.66	11.63	12.36	7.53	8.00	11.87	7.92	6.90	7.04	10.34
		<b>10.73</b>	12.20	<b>11.24</b>	13.77	<b>6.34</b>	<b>7.25</b>	<b>10.04</b>	<b>10.67</b>	<b>10.60</b>	8.36	<b>5.74</b>	<b>9.84</b>	<b>5.46</b>	<b>6.05</b>	<b>5.83</b>	<b>8.94</b>

Table 12. **Performance comparison (HD95 mm) on BraTS 2024.** Lower is better. Available and missing modalities are denoted by ● and ○, respectively. ■, ■, and ■ indicate the first-, second-, and third-best performance, respectively. WT, TC, and ET represent whole tumor, tumor core, and enhancing tumor, respectively.

Crop Size	Average Region DSC (%)			Overall
	WT	TC	ET	Mean
80	90.17	84.22	74.54	82.97
96	90.38	<b>84.51</b>	<b>75.59</b>	<b>83.49</b>
128	<b>90.87</b>	84.25	74.64	83.25

Table 13. **Ablation on the crop size.** ■, ■, and ■ indicate the first, second, and third best results, respectively.

Patch Size	Average Region DSC (%)			Overall
	WT	TC	ET	Mean
8	<b>90.38</b>	<b>84.51</b>	<b>75.59</b>	<b>83.49</b>
12	89.99	83.89	74.33	82.74
16	89.71	83.21	73.74	82.22

Table 14. **Ablation on the patch size.** ■, ■, and ■ indicate the first, second, and third best results, respectively.

released code. All experiments were conducted on a single NVIDIA RTX Pro 6000 GPU with 96 GiB of memory. To clearly quantify differences in training efficiency, we report the time required to complete 250 iterations using their recommended batch size and crop size.

Despite its larger parameter count, UniME attains the highest throughput. This significant efficiency gain is

primarily due to algorithmic optimizations, including the use of FlashAttention [11], mixed-precision training, and PyTorch graph compilation. These optimizations mitigate the computational overhead introduced by the Transformer layers, ensuring that increased model complexity does not compromise practical training efficiency. Moreover, UniME exhibits affordable GPU memory consumption during training (23.80 GiB), underscoring efficient use of computational resources.

## F. Limitations

Although UniME attains state-of-the-art results on BraTS 2023 and BraTS 2024, several limitations remain. **Two-stage Objective Difference.** Stage 1 uses pixel-wise reconstruction loss, whereas Stage 2 optimizes semantic segmentation. This difference may induce a suboptimal inductive bias. Although we mitigate this via (i) using a lightweight decoder in Stage 1 to steer the Uni-Encoder toward semantic features, and (ii) re-introducing a CNN decoder in Stage 2 to emphasize semantic expression, the objective gap may still limit further gains. Future work will explore better representation learning strategies for medical image scenarios, including missing-modality cases.

Methods	Configuration			Complexity				Average Region DSC (%)			Overall
	Batch Size	Crop	Optimizer	Params (M)	Memory (GiB)	Elapsed Time (s)	Throughput (samples/s)	WT	TC	ET	Mean
RFNet	2	80	Adam	8.99	19.08	88.36	5.66	87.60	79.20	66.89	77.90
mmFormer	4	128	Adam	36.65	57.64	293.80	3.40	87.48	78.49	66.11	77.36
IM-Fuse	2	128	RAdam	64.39	56.37	354.19	1.41	88.74	80.59	70.67	80.00
M <sup>2</sup> SegMamba	1	128	Adam	70.48	51.72	234.12	1.07	88.99	82.11	71.73	80.94
Ours	4	96	AdamW	158.11	23.80	79.52	12.58	<b>90.38</b>	<b>84.51</b>	<b>75.59</b>	<b>83.49</b>

Table 15. **Efficiency comparison across methods.** ■, ■, and ■ indicate the first-, second-, and third-best performance, respectively. Params refers to the total number of parameters, measured in millions (M). Memory usage, measured in gigabytes (GiB), represents the peak GPU memory consumed during training. Elapsed Time indicates the time for 250 iterations during training, measured in seconds (s). Throughput measures the number of samples processed per second. Average Region DSC and Overall DSC are reported based on evaluations on the BraTS 2023 dataset.

**Dataset Limitations.** We currently evaluate only on BraTS 2023 and BraTS 2024 datasets. Ideally, generalization should be validated on more datasets. However, suitable alternatives are limited. More specifically, many multimodal datasets exhibit challenges such as unaligned modalities in ISLES [26], incomplete modality coverage in CHAOS [29] and AMOS [28], and predominantly single-modality data in MSD [2] (with brain cases originating from BraTS).

In alignment with previous research [16, 17, 24, 36, 38, 45, 67–69], we have thus adopted the BraTS series as our primary evaluation benchmark. Although we recognize the importance of diverse datasets for comprehensive validation, our current assessments leverage the most suitable and widely accepted resources available. Future research efforts will seek additional datasets tailored to our specific task requirements and further investigate two-stage heterogeneous architectures in broader vision tasks.