

UNLOCKING CLINICAL POTENTIAL: BEYOND SINGLE-TO-TRI-PHASE CT WITH DYNAMIC FUSION FOR LIVER TUMOR SEGMENTATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Liver tumor segmentation is essential for treatment planning and disease monitoring. Most existing methods rely on single-phase computed tomography (CT), they often suffer from low contrast and incomplete lesion depiction. Contrast-Enhanced CT (CECT) offers multiple imaging phases: arterial (ART), portal venous (PV), and delayed (DL), which provide complementary anatomical and functional information. This study begins with a systematic quantitative evaluation of each enhanced phase using standard segmentation models to investigate their individual contributions and validate phase-specific clinical insights. Guided by this analysis, a Multi-phase Attention Deep Fusion Network (**MADF-Net**) is proposed to hierarchically integrate ART, PV, and DL features across the input, feature, and decision levels. Experiments on the clinically collected multi-phase liver lesion (MPLL) dataset (the largest and most clinically comprehensive multi-phase liver cancer CECT dataset) demonstrate that the proposed method achieves state-of-the-art segmentation performance. **MADF-Net** achieves a Dice score of **78.65%**, which is **9.39%** higher than the best single-stage baseline, by deeply fusing information from three phases, and consistently improves across all evaluation metrics. Our codes are available at https://anonymous.4open.science/r/ICLR26_unlocking_clinical_potential-EFE8/.

1 INTRODUCTION

Liver tumor segmentation is a critical task in quantitative medical image analysis, providing essential morphological and spatial information for surgical planning, radiotherapy, and post-treatment monitoring (Bilic et al., 2023). With the advent of deep learning, fully convolutional networks

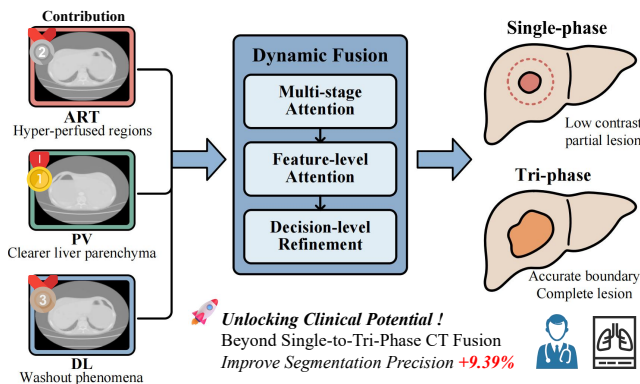


Figure 1: Motivation of our multi-phase CT fusion strategy. Different CT phases provide complementary information: ART highlights vessels, PV improves lesion–parenchyma contrast, and DL captures delayed enhancement. Their integration yields more accurate, robust liver tumor segmentation.

(FCNs), particularly U-Net and its variants (Ronneberger et al., 2015; Ren & Li, 2025; Du et al., 2022; Huang et al., 2017), have achieved notable success in automated segmentation tasks. These models extract features at either the 2D slice level or the 3D volumetric level, enabling robust representation learning for complex anatomical structures (Li et al., 2018).

However, the majority of existing methods rely solely on single-phase Computed Tomography (CT) images (Bilic et al., 2023; Wang et al., 2023; Hatamizadeh et al., 2022; Cao et al., 2022), often ignoring the phase-specific characteristics inher-

ent in clinical imaging protocols (Jun et al., 2023). Due to low tissue contrast and resolution limitations, single-phase methods frequently fail to achieve the precision required for clinical deployment

(Song et al., 2023; Zhang et al., 2024). Contrast-Enhanced CT (CECT), which captures dynamic changes in tissue attenuation following contrast agent administration, provides a valuable alternative by acquiring images at multiple time points—typically including the non-contrast (NC), arterial (ART), portal venous (PV), and delayed (DL) phases (Chi et al., 2013).

Among these, the NC phase offers a baseline anatomical information, but lacks enhancement patterns relevant to the tumor vasculature and lesion contrast, and is therefore generally not emphasized in liver tumor segmentation studies (Ni et al., 2024; Xu et al., 2021; Liu et al., 2024). The ART phase captures early vascular features, highlighting hyper-perfused regions and enhancing lesion boundary delineation (Kulkarni et al., 2021; Urban et al., 2000). The PV phase provides clearer liver parenchyma and structural completeness, facilitating more accurate segmentation (Kulkarni et al., 2021; Schneider et al., 2014). The DL phase captures delayed enhancement and washout phenomena, aiding in the identification of fibrotic or hypo-perfused tumors (Monzawa et al., 2007; Lim et al., 2002). The complementary nature of these phases offers a compelling opportunity for improved segmentation through multi-phase fusion (As shown in Figure 1). Therefore, how to effectively extract and fuse the features from different phases has attracted the attention of many researchers.

Existing multi-phase fusion strategies can be broadly categorized into three types (Zhang et al., 2021b): (1) *Input-level fusion* (Ouhmich et al., 2019), where multiple phases are concatenated as input and processed via a shared encoder; (2) *Feature-level fusion* (Zhang et al., 2021b; Zhu et al., 2022; Zhang et al., 2023; Hazirbas et al., 2016; Liu et al., 2023), which extracts features from each phase independently before combining them at intermediate layers; and (3) *Decision-level fusion* (Sun et al., 2017; Raju et al., 2020), where each phase is processed by a separate network and results are fused at the output level. While these approaches have shown potential, they often suffer from limitations such as insufficient modeling of nonlinear inter-phase relationships (Sun et al., 2017; Zhang et al., 2023), reduced reliability in ambiguous or low-contrast regions, and vulnerability to missing-phase scenarios common in clinical workflows (Xu et al., 2021; Zhu et al., 2022). Moreover, many existing methods treat all phases with equal importance during fusion, overlooking their distinct clinical value and the complementary information they offer (Xu et al., 2021; Zhong et al., 2024; Qiao et al., 2024). This results in suboptimal performance, especially in cases with blurred lesion boundaries or small lesions. **Therefore, how to effectively fuse multi-phase CT features while leveraging their individual strengths and mitigating their limitations remains an open challenge** (Jiang et al., 2020).

In this paper, we begin by systematically evaluating the segmentation performance of each enhanced CT phase using standard deep learning models. Our quantitative analysis reveals that the PV phase contributes most significantly to segmentation accuracy, consistent with its known clinical role. Guided by this observation, we propose a novel framework, Multi-phase Attention Deep Fusion Network (MADF-Net), to exploit the complementary advantages of the ART, PV, and DL phases through hierarchical fusion. MADF-Net introduces full-stage attention-based fusion across the input, feature, and decision levels, enabling deep inter-phase information interaction. Extensive experiments on the MPLL dataset demonstrate that our method achieves state-of-the-art segmentation performance, reaching a Dice score of 78.65% when using all three phases, representing a 9.39% improvement over the best single-phase baseline, confirming its robustness and generalizability.

Our contributions are as follows:

- ① We conduct a comprehensive quantitative analysis of liver tumor segmentation across different CT phases and demonstrate the predominant contribution of the PV phase, providing both empirical and clinical insights.
- ② We propose MADF-Net, a multi-phase attention-based fusion network that integrates ART, PV, and DL phase features at multiple stages, enhancing liver tumor segmentation performance through deep inter-phase feature interaction.
- ③ Extensive experiments on a newly collected multi-phase liver lesion (MPLL) benchmark (the largest and most clinically comprehensive multi-phase liver cancer CECT dataset) demonstrate that the proposed method achieves state-of-the-art liver tumor segmentation performance.

2 RELATED WORKS

Single-Phase Based Liver Tumor Segmentation. Deep learning has significantly advanced single-phase liver tumor segmentation in CT images. Ronneberger et al. (Ronneberger et al., 2015)

introduced U-Net, whose encoder-decoder structure with skip connections became foundational, effectively capturing both local details and global context for handling low contrast and fuzzy boundaries. H-DenseUNet (Li et al., 2018) enhanced feature reuse through hybrid dense connections and achieved state-of-the-art results on the LiTS2017 dataset. Variants such as UNet++ (Zhou et al., 2018) further improved efficiency and multi-scale accuracy, particularly for small tumor detection. However, CNN-based models often struggled with global context in complex tumor structures. To address this, TransUNet (Chen et al., 2021) combined CNNs for low-level feature extraction with Transformers for global dependency modeling. UNETR (Hatamizadeh et al., 2022) and UNETR++ (Shaker et al., 2024) integrated global context and local detail via a Transformer-based U-shaped architecture, achieving strong performance on 3D CT tasks. Nevertheless, single-phase methods still suffer from limited tissue contrast and resolution, resulting in information loss and reduced clinical applicability.

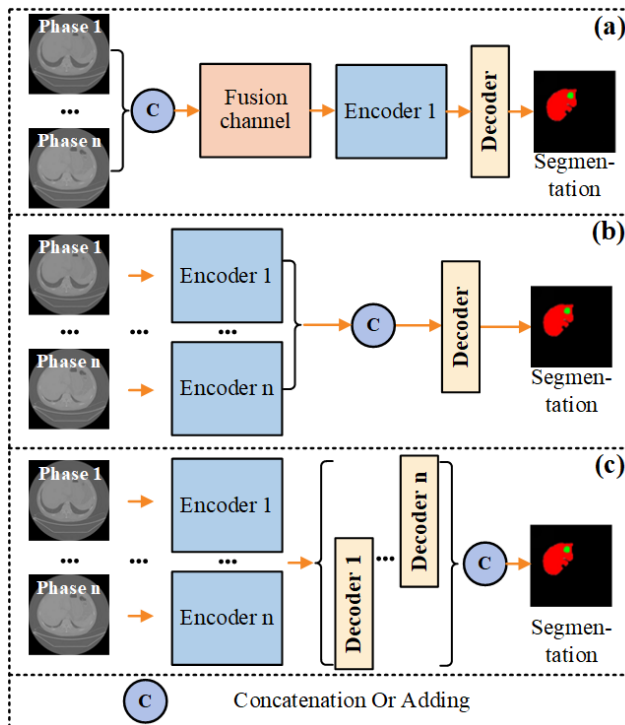


Figure 2: Multi-phase fusion method of enhanced CT. (a), (b), and (c) correspond to input-level, feature-level, and decision-level fusion architectures, respectively.

enable dense information exchange. Wu et al. (Wu et al., 2019) treated non-contrast and enhanced CT scans equally and applied feature-level fusion at selected U-Net layers. In decision-level fusion, features are independently extracted from each phase and fused at higher layers. Raju et al. (Raju et al., 2020) proposed an integrated joint and semi-supervised training strategy that leveraged limited plain and enhanced CT data to achieve robust cross-domain segmentation. Despite progress in single-stage fusion, these methods still face challenges such as information loss and limited ability to capture phase-specific characteristics. **② Multi-stage Fusion.** Feature-level and decision-level combinations currently dominate multi-stage fusion network designs (Ni et al., 2024; Liu et al., 2024; Zhu et al., 2022; Kuang et al., 2024). PA-ResSeg (Xu et al., 2021) introduced intra- and inter-phase attention mechanisms to capture both channel-wise dependencies and cross-phase interactions, embedding attentional modules at each encoder layer to fuse multi-scale information from ART and PV phases. Building on this, SA-Net (Zhang et al., 2021b) added a spatial aggregation module for encoding-stage interaction and an uncertainty correction module at the decision stage to refine fuzzy tumor boundaries. To address spatial misalignment in multi-phase CT, Zhang et al. (Zhang et al., 2023) incorporated differentiable deformation operations (Jaderberg et al., 2015) for enhanced feature alignment. Raju et al. (Raju et al., 2020) proposed a joint and semi-supervised training strategy that effectively leveraged limited non-contrast and enhanced CT data, though at the cost of increased training time. HRadNet (Liang et al., 2023) utilized a feature

Multi-Phase Based Liver Tumor Segmentation.

Recently, an increasing number of studies have investigated how to leverage multi-phase CT information to improve liver tumor segmentation performance. Multiphase fusion is typically performed at one of three stages: input-level, feature-level, or decision-level fusion (Zhang et al., 2021b), referred to as single-stage fusion in this paper. Alternatively, fusion can occur across multiple stages, which we define as multi-stage fusion. **① Single-stage Fusion.**

An early example of single-stage input-level fusion was proposed by Ouhmich et al. (Ouhmich et al., 2019), who concatenated PV and ART phase images as input to a U-Net, significantly improving tumor segmentation performance over single-phase training. Feature-level fusion is currently the most active research area. Zhou et al. (Zhou et al., 2019) introduced a dual-path 3D fully convolutional network with cross-phase skip connections to

pyramid and a metadata fusion layer to incorporate clinical features such as tumor size and patient age, improving generalizability. However, most existing multi-stage approaches adopt only two fusion stages and still suffer from potential information loss, limiting segmentation accuracy and robustness. To address this, we propose a three-stage fusion network designed to better preserve information throughout the extraction and fusion process.

3 PRELIMINARY

To fully exploit the complementary anatomical and pathological information provided by the three phases of CECT, we propose a novel three-stage fusion framework, named Multi-phase Attention Deep Fusion Network (**MADF-Net**). This section first introduces preliminary knowledge on fusion strategies, and then describes the proposed **MADF-Net**. As shown in Figure 2, three common fusion strategies are illustrated. A detailed pseudocode description of the overall procedure is provided in Appendix A.

Input-level Fusion. As shown in Figure 2 (a), this strategy concatenates images from different phases (Phase 1, Phase 2, ..., Phase n) along the channel dimension at the input stage to form a unified input tensor. To enhance flexibility, we introduce learnable phase-wise modulation weights $\{\alpha_i\}_{i=1}^n$ and a normalization operator $\mathcal{N}(\cdot)$:

$$\mathbf{I}_{\text{input}} = \mathcal{N}\left(\left\|_{i=1}^n (\alpha_i \cdot \Gamma(\mathbf{I}_i) + \beta_i \cdot \mathbf{1}_{H \times W \times C})\right.\right), \alpha_i = \frac{\exp(\theta_i)}{\sum_{j=1}^n \exp(\theta_j)}, \quad (1)$$

where $\Gamma(\cdot)$ denotes intensity standardization, $\mathbf{I}_i \in \mathbb{R}^{H \times W \times C}$ is the i -th phase image, θ_i are learnable logits and $\|$ denotes channel-wise concatenation. This formulation adaptively highlights more informative phases while suppressing noisy ones.

Feature-level Fusion. As shown in Figure 2 (b), this strategy integrates features from multiple phases by using attention-guided gating and nonlinear projections. Let $\mathbf{E}_i \in \mathbb{R}^{H' \times W' \times d}$ be the feature map extracted from the i -th phase. We compute phase attention maps \mathbf{A}_i from global descriptors \mathbf{g}_i via a softmax-normalized MLP, and then fuse features as:

$$\mathbf{E}_{\text{fusion}} = \Psi\left(\sum_{i=1}^n \left(\underbrace{\frac{\exp(\mathbf{W}_a \mathbf{g}_i + \mathbf{b}_a)}{\sum_{j=1}^n \exp(\mathbf{W}_a \mathbf{g}_j + \mathbf{b}_a)}}_{\mathbf{A}_i} \odot \underbrace{(\mathbf{W}_v * \mathbf{E}_i + \mathbf{b}_v)}_{\bar{\mathbf{E}}_i}\right)\right), \mathbf{g}_i = \text{GAP}(\mathbf{E}_i), \quad (2)$$

where $*$ denotes convolution, $\text{GAP}(\cdot)$ is global average pooling, and $\Psi(\cdot)$ is a residual refinement block. This design allows adaptive semantic fusion guided by global context cues.

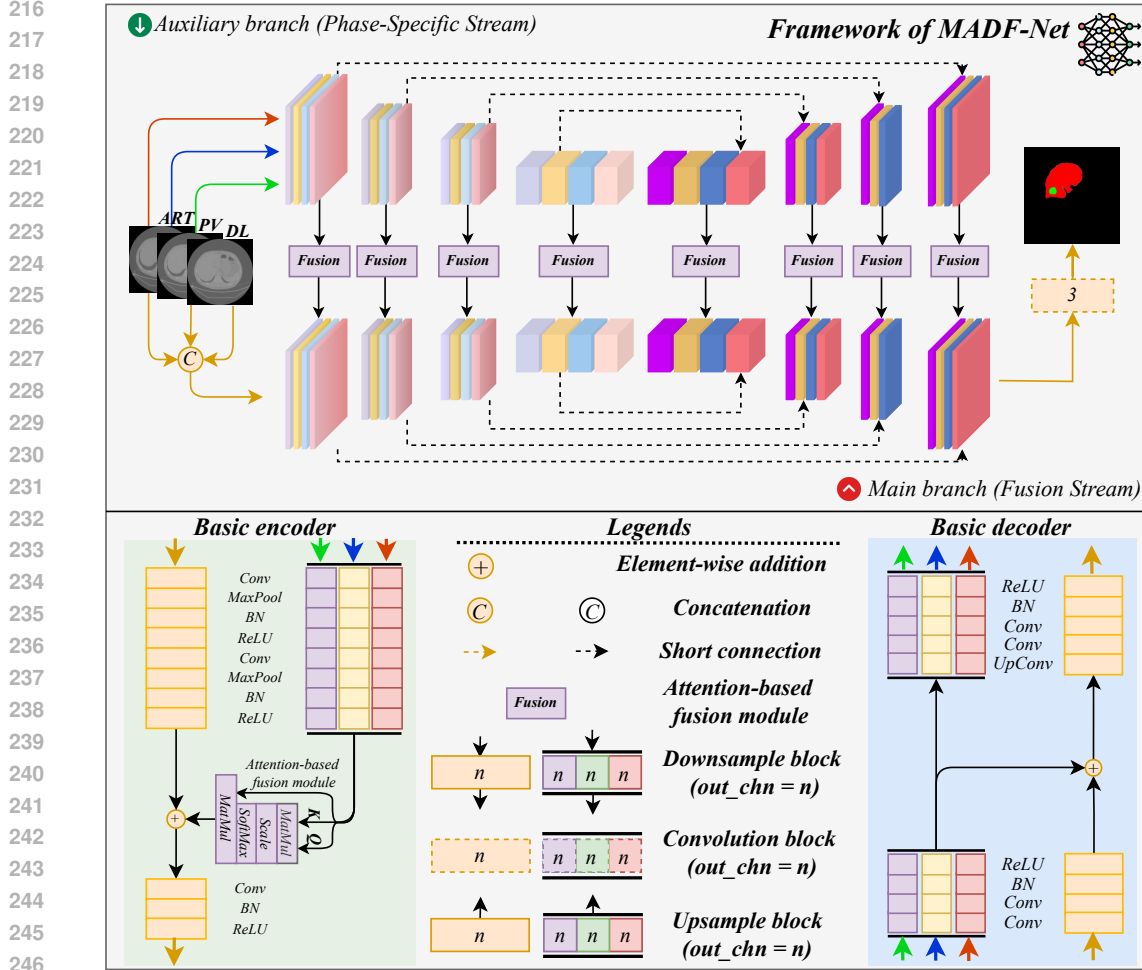
Decision-level Fusion. As shown in Figure 2 (c), this strategy constructs separate segmentation heads for each phase, and aggregates the resulting predictions $\{\mathbf{S}_i\}_{i=1}^n$ based on their confidence. We employ an uncertainty-aware soft weighting scheme with temperature scaling:

$$\mathbf{S}_{\text{final}} = \sum_{i=1}^n \left(\underbrace{\frac{\exp(-\tau \cdot \mathcal{H}(\mathbf{S}_i))}{\sum_{j=1}^n \exp(-\tau \cdot \mathcal{H}(\mathbf{S}_j))}}_{w_i} \cdot \underbrace{\sigma(\mathbf{S}_i)}_{\text{probability map}} \right), \mathcal{H}(\mathbf{S}_i) = -\frac{1}{|\Omega|} \sum_{p \in \Omega} \sum_c \mathbf{S}_i^{(p,c)} \log \mathbf{S}_i^{(p,c)}, \quad (3)$$

where $\mathcal{H}(\cdot)$ computes the spatial entropy over pixel set Ω , and τ controls weight sharpness. This formulation emphasizes confident predictions and suppresses noisy ones for decision-level fusion.

4 METHODOLOGY

To address the challenge of liver tumor segmentation using multiphase CT data, we propose a unified framework, **MADF-Net**, that integrates input-level, feature-level, and decision-level fusion. As shown in Figure 3, the network consists of two parallel branches (main and auxiliary) with symmetric encoder-decoder structures, enabling hierarchical feature aggregation across three CT phase.



247 Figure 3: Overview of MADF-Net. Three single-phase CT inputs (ART, PV, DL) are processed via parallel
248 main/auxiliary branches. Fusion occurs at three levels: (1) Input-level (concatenation of phases in main branch);
249 (2) Feature-level (cross-branch fusion in encoder/decoder blocks); (3) Decision-level (final result aggregation).

251 **Input-Level Fusion: Multi-Phase Data Initialization.** Given the three CT phases I_{ART} , I_{PV} ,
252 and I_{DL} , the main branch performs early-stage fusion by channel-wise concatenation to directly
253 expose the encoder to cross-phase correlations: $I_{fusion} = \text{Concat}(I_{ART}, I_{PV}, I_{DL}) \in \mathbb{R}^{H \times W \times 3C}$,
254 where H and W denote the spatial resolution and C denotes the number of channels per phase.
255 In parallel, the auxiliary branch independently forwards each phase through three isolated sub-
256 encoders: $I_{phase}^{(p)} = \{I_{ART}, I_{PV}, I_{DL}\}$, $p \in \{1, 2, 3\}$, preserving phase-specific characteristics
257 that might otherwise be suppressed by early fusion.

259 **Encoder Blocks: Hierarchical Feature Extraction with Cross-Branch Fusion.** Both branches
260 comprise four encoder blocks indexed by $l \in \{1, \dots, 4\}$, progressively downsample spatial resolu-
261 tion while expanding feature dimensionality (e.g., $C \rightarrow 4C \rightarrow 16C \rightarrow 64C$). Encoder block performs
262 following operations: **(I) Self-Attention-based Phase Reweighting.** For auxiliary features $X_{ART}^{(l)}$,
263 $X_{PV}^{(l)}$, and $X_{DL}^{(l)}$, a shared self-attention module computes attention weights across phases:
264

265
266
267

$$\alpha^{(l)} = \text{softmax}\left(\frac{Q^{(l)} \cdot (K^{(l)})^\top}{\sqrt{d_k}}\right), \quad Q^{(l)} = W_Q X^{(l)}, \quad K^{(l)} = W_K X^{(l)}, \quad (4)$$

268 where $W_Q, W_K \in \mathbb{R}^{d_k \times d_{in}}$ are learnable projections and $X^{(l)} = [X_{ART}^{(l)}, X_{PV}^{(l)}, X_{DL}^{(l)}]$. The
269 reweighted auxiliary feature is then: $X_{aux}^{(l)} = \sum \alpha_p^{(l)} \odot X_p^{(l)}$, $p \in \{ART, PV, DL\}$

(II) Feature-level Cross-Branch Fusion. After obtaining the reweighted auxiliary representation, we further integrate it with the main branch feature $\mathbf{X}_{main}^{(l)}$. To this end, we design a *gated residual summation* mechanism:

$$\mathbf{X}_{enc}^{(l)} = \sigma\left(\mathbf{W}_g * [\mathbf{X}_{main}^{(l)}, \mathbf{X}_{aux}^{(l)}]\right) \odot \mathbf{X}_{main}^{(l)} + \left(1 - \sigma\left(\mathbf{W}_g * [\mathbf{X}_{main}^{(l)}, \mathbf{X}_{aux}^{(l)}]\right)\right) \odot \mathbf{X}_{aux}^{(l)}, \quad (5)$$

where the gating factor is adaptively determined by the concatenated representations from both branches. This allows the model to dynamically balance their contributions, avoiding redundancy and gradient dilution caused by naive concatenation or summation. Specifically, when the gate approaches 1, the model emphasizes high-level semantics from the main branch, while values closer to 0 highlight fine-grained cues from the auxiliary branch.

Nevertheless, simply stacking the above fusion operation across multiple depths may hinder gradient propagation, thus limiting the representation capacity of deeper layers. To address this, we introduce a *residual preservation regularization* to facilitate cross-layer information flow. Concretely, at the l -th layer, the fused feature $\mathbf{X}_{enc}^{(l)}$ is enhanced with a gated residual connection from previous layer:

$$\hat{\mathbf{X}}_{enc}^{(l)} = \mathbf{X}_{enc}^{(l)} + \lambda \cdot \left(\alpha^{(l)} \odot \hat{\mathbf{X}}_{enc}^{(l-1)} + (1 - \alpha^{(l)}) \odot f\left(\hat{\mathbf{X}}_{enc}^{(l-1)}\right)\right), \quad l > 1, \quad (6)$$

where λ is a learnable global balancing coefficient, and $\alpha^{(l)}$ is a layer-wise gating vector that adaptively controls the trade-off between direct residual propagation and a transformed path. The function $f(\cdot)$ denotes a lightweight non-linear mapping (e.g., a convolutional projection or an MLP). This design ensures that shallow features can effectively penetrate deeper layers to improve gradient flow, while the nonlinear transformation path enriches cross-layer feature diversity.

Decoder Blocks: Multi-Scale Feature Reconstruction. The decoder consists of four blocks that mirror the encoder structure, progressively upsampling the fused representations $\mathbf{X}_{enc}^{(l)}$ back to the original resolution. Each decoder block not only restores spatial resolution but also selectively incorporates complementary information from shallow layers through gated skip connections. Specifically, at each stage l , we compute:

$$\mathbf{Y}^{(l-1)} = \text{ReLU}\left(\mathbf{G}^{(l)} \odot \text{UpConv}(\mathbf{X}_{enc}^{(l)}) + (1 - \mathbf{G}^{(l)}) \odot \mathcal{F}_{att}(\mathbf{X}_{enc}^{(l-1)}, \text{UpConv}(\mathbf{X}_{enc}^{(l)}))\right), \quad (7)$$

where $\mathbf{G}^{(l)}$ is a learned gating map, $\text{UpConv}(\cdot)$ denotes an upsampling convolution block with batch normalization and activation, and $\mathcal{F}_{att}(\cdot, \cdot)$ is an attention-based fusion module for shallow-deep interaction. This stage-wise design ensures that high-resolution details from earlier layers are progressively blended with the deep semantic context from later layers.

To further enhance multi-scale consistency and stabilize gradient flow, we augment the reconstruction with a residual-preserving multi-scale aggregation term:

$$\hat{\mathbf{Y}}^{(l-1)} = \mathbf{Y}^{(l-1)} + \mu \cdot \left(\sum_{k=1}^K \beta_k^{(l)} \odot \text{UpConv}^{(k)}(\mathbf{X}_{enc}^{(l-k)})\right), \quad (8)$$

where μ is a learnable global scaling factor, $\beta_k^{(l)}$ are adaptive weights normalized by a softmax constraint, and $\text{UpConv}^{(k)}(\cdot)$ denotes k -step hierarchical upsampling operators. This formulation explicitly aggregates contextual evidence from multiple encoder depths, enabling the decoder to reconstruct fine details while preserving long-range semantic dependencies.

In summary, the decoder leverages a combination of gated skip fusion, residual-preserving connections, and multi-scale aggregation to ensure both spatial fidelity and semantic consistency. Such a design alleviates the common issue of blurred boundaries in dense prediction tasks, while also enhancing robustness against vanishing gradients during backpropagation.

Decision-Level Fusion: Final Segmentation Output. Finally, the outputs of the main and auxiliary decoders are aggregated to produce the segmentation mask. Specifically,

$$\mathbf{O}_{final} = \sigma\left(\mathbf{W}_{out} * [\mathbf{O}_{ART}, \mathbf{O}_{PV}, \mathbf{O}_{DL}, \mathbf{O}_{fusion}] + \mathbf{b}_{out}\right), \quad (9)$$

where \mathbf{O}_{ART} , \mathbf{O}_{PV} , \mathbf{O}_{DL} are the three auxiliary outputs, \mathbf{O}_{fusion} is the main-branch output, and \mathbf{W}_{out} , \mathbf{b}_{out} denote the parameters of the final convolutional projection. This decision-level fusion enforces complementary exploitation of both phase-specific and cross-phase knowledge, yielding a precise tumor segmentation mask.

5 EXPERIMENTS

Dataset Curation (Multi-phase Dataset). The multi-phase liver lesion (MPLL) dataset, consists of 952, 601 2D slices with liver disease from the "Anonymous Authoritative Hospitals (Information will be made public after the paper is accepted)". The dataset includes patients aged between 9 and 72 years, and the number of axial slices per scan varying from 48 to 777. This is the largest and most valuable multi-phase CECT liver cancer dataset to date, all the images in MPLL dataset contain three enhanced phases (ART, PV, and DL). The registered images were annotated using ITK-SNAP software by two experienced attending radiologists, and subsequently reviewed by a third attending radiologist to ensure the accuracy and consistency of the annotations. All data have been anonymized and contain only image information. The MPLL dataset has received approval from the institutional ethics committee under certification number 2022-BE(H)-194. Figure 6 shows example images from the datasets. The training, validation, and test splits (7:1:2, following previous work (Jiang et al., 2023)), along with image dimensions and other details, are summarized in Table 1. A more detailed description is provided in Appendix B.

Table 1: Dataset characteristics.

Dataset	Attribute	Value
MPLL	Phase	ART, PV, DL
	Slice thickness	0.62mm–5.0mm
	Slice resolution	512 × 512
	Disease type	ABS, HCC, HEM, ICC, Lipoma

Evaluation Metrics. We employed the Dice Similarity Coefficient (DSC), Jaccard Similarity Coefficient (JSC), Average Symmetric Surface Distance (ASSD), and 95% Hausdorff Distance (HD_{95}) (Jiang et al., 2025) to evaluate the experimental results. In the experiments on single-phase (1P), two-phase (2P) and three-phase (3P) input, we additionally employed Volume Overlap Error (VOE) and Relative Volume Difference (RVD) as supplementary evaluation metrics.

Implementation Details. All models were trained for 100 epochs with a batch size of 8. The Stochastic Gradient Descent (SGD) optimizer was adopted with a learning rate of 0.01 and 4 parallel data loading workers. Data augmentation techniques include horizontal flipping (with probability 0.5) and vertical flipping (with probability 0.5), and no post-processing is used.

The proposed method was implemented on a Linux 5.4.0 system using PyTorch 1.13.1. All experiments were conducted on two NVIDIA GeForce RTX 3090 GPUs (24 GB × 2), providing sufficient computational resources for efficient model training and evaluation.

5.1 MAIN RESULT

Table 2: Quantitative comparison of segmentation performance across different phase combinations (relative to 3P, based on numerical differences).

Input Phases	DSC(%) \uparrow	JSC(%) \uparrow	HD_{95} \downarrow	ASSD \downarrow
1P (PV)	69.26 _{-9.39}	64.08 _{-10.73}	40.298 _{+13.492}	16.958 _{+6.222}
2P (ART+PV)	76.09 _{+2.56}	71.41 _{+3.40}	28.838 _{+2.032}	15.659 _{+4.923}
3P (ART+PV+DL)	78.65	74.81	26.806	10.736

proposed MADF-Net in leveraging complementary information across multiple imaging phases. Specifically, the best result of 1P input (Zheng et al., 2024b) using only the PV phase achieves a DSC of 69.26% and a JSC of 64.08%. When the ART phase is added to form the 2P input (ART+PV), both DSC and JSC show moderate improvements to 76.09% and 71.41%, respectively, while the HD_{95} drops significantly from 40.298 to 28.838, indicating better boundary localization. The 3P input (ART+PV+DL) achieves the best overall performance, with the highest DSC (78.65%) and JSC (74.81%), along with the lowest HD_{95} (26.80) and Average ASSD of 10.73. These results suggest that the additional information from the ART and DL phases enhances both global overlap and local boundary accuracy. Compared to the 1P setting, the 3P input yields substantial gains of 9.39% in DSC and 10.73% in JSC, underscoring the benefit of multi-phase fusion in capturing diverse tumor characteristics.

Obs. ②: Comparison of Single Phase Performance. We conducted a segmentation performance comparison using different single-phase input images on the MPLL dataset, evaluating several state-of-the-art models, including FANet (Tomar et al., 2022), GRENet (Wang et al., 2023), ASSNet (Zheng et al., 2024a), TransUNet (Chen et al., 2021), KiU-Net (Valanarasu et al., 2021),

Obs. ①: Phase Combination Analysis. We conducted three groups of experiments to explore the optimal input combination for 1P, 2P, and 3P settings. As listed in Table 2, the number of input phases increases, segmentation accuracy improves accordingly, demonstrating the effectiveness of the

Table 3: Quantitative comparison of different methods on the MPLL dataset. Best results are **bold**, and differences relative to PV are shown as colored arrows.

Methods	DSC(%) \uparrow			Jaccard(%) \uparrow			HD ₉₅ (mm) \downarrow			ASSD (mm) \downarrow		
	ART	PV	DL	ART	PV	DL	ART	PV	DL	ART	PV	DL
FANet	67.56 \downarrow _{1.52}	69.08	66.01 \downarrow _{3.07}	61.98 \downarrow _{1.11}	63.09	60.74 \downarrow _{2.35}	60.384 \uparrow _{17.529}	42.855	62.083 \uparrow _{19.228}	30.300 \uparrow _{10.107}	20.193	25.206 \uparrow _{5.013}
GRENet	66.59 \downarrow _{1.67}	68.26	66.21 \downarrow _{2.05}	60.14 \downarrow _{0.94}	61.08	59.91 \downarrow _{1.17}	60.651 \uparrow _{3.390}	57.261	65.412 \uparrow _{8.151}	28.069 \uparrow _{4.059}	24.010	30.982 \uparrow _{6.972}
ASSNet	67.91 \downarrow _{1.35}	69.26	66.49 \downarrow _{2.77}	60.01 \downarrow _{4.07}	64.08	60.07 \downarrow _{4.01}	45.260 \uparrow _{4.962}	40.298	59.946 \uparrow _{19.648}	22.074 \uparrow _{5.116}	16.958	25.031 \uparrow _{8.073}
TransUNet	67.51 \downarrow _{1.07}	68.48	67.15 \downarrow _{1.33}	58.65 \downarrow _{2.58}	61.23	60.53 \downarrow _{0.70}	55.960 \downarrow _{-0.339}	55.301	58.594 \uparrow _{3.293}	26.983 \downarrow _{0.346}	27.329	24.972 \downarrow _{2.357}
KiU-Net	65.92 \downarrow _{1.70}	67.62	65.62 \downarrow _{2.00}	59.64 \downarrow _{1.24}	60.88	59.42 \downarrow _{1.46}	62.902 \downarrow _{3.642}	59.260	65.956 \uparrow _{6.696}	30.821 \downarrow _{2.178}	28.643	27.983 \downarrow _{0.660}
AttUNet	66.01 \downarrow _{2.51}	68.52	65.85 \downarrow _{2.67}	59.82 \downarrow _{1.84}	61.66	59.61 \downarrow _{2.05}	60.913 \uparrow _{5.289}	55.624	66.583 \uparrow _{10.959}	31.098 \uparrow _{4.117}	26.981	27.973 \downarrow _{-3.008}

and AttUNet (Chen et al., 2023). The results are summarized in Table 3. Across all models, segmentation performance was consistently better on the PV phase compared to the ART and DL phases. For example, with FANet, the DSC and JSC on PV were 1.52% and 1.11% higher than on ART, and 3.07% and 2.35% higher than on DL, respectively. In terms of boundary metrics, the HD₉₅ and ASSD on PV were 17.529 and 10.107 lower than on ART, and 19.228 and 5.013 lower than on DL, indicating more accurate boundary localization. These results suggest that segmentation outputs on the PV phase more closely match the ground truth. In clinical contexts, this may be attributed to the PV phase offering more distinct grayscale contrast between tissues, as well as between lesions and normal structures (Ni et al., 2024; Lam et al., 2017; Al-Battal et al., 2024). This contrast enhancement is particularly beneficial in cases with ambiguous or highly heterogeneous tumor boundaries (Liu et al., 2024).

Obs. ③: 2P Fusion: ART and PV. Single-phase experimental results indicate that segmentation performance is high when using the PV and ART phases as inputs. Based on this observation, we conducted 2P (A+P) experiments on the MPLL dataset, comparing the proposed MADF-Net with several state-of-the-art multi-phase segmentation methods, including MAML (Zhang et al., 2021a), MW-UNet (Zhu et al., 2022), SA-Net (Zhang et al., 2021b), PA-ResSeg (Xu et al., 2021), and MCDA-Net (Kuang et al., 2024). As listed in Table 4, MADF-Net achieves superior performance, improving the DSC metric by 5.27%, 4.71%, 4.18%, 1.90%, and 0.01% over the five comparison methods, respectively. It also consistently outperforms all other methods in terms of HD₉₅, demonstrating its ability to exploit complementary information across imaging phases to enhance both segmentation accuracy and boundary localization.

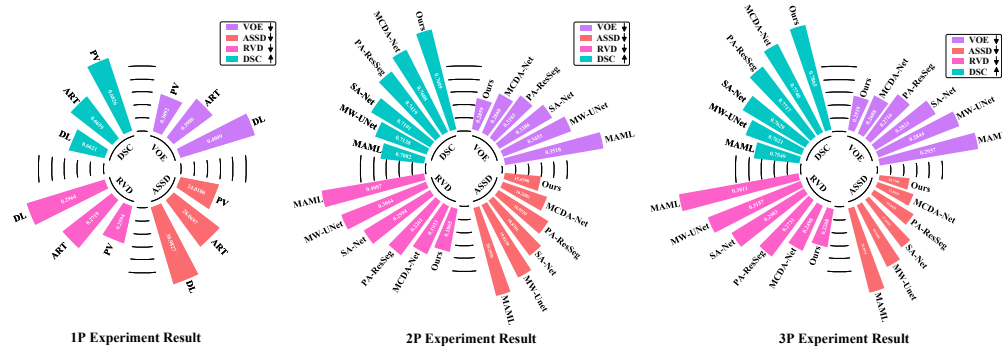


Figure 4: Comparison of the segmentation performance on 1P, 2P and 3P, evaluated using multiple quantitative metrics including DSC, VOE, RVD, and ASSD. Specifically, the black upward arrow (\uparrow) denotes that a higher metric value indicates superior performance, while the black downward arrow (\downarrow) signifies that a lower metric value reflects more favorable outcomes.

Table 4: Quantitative comparison of 2-phase (left) and 3-phase (right) inputs. ‘‘Phase’’ indicates the data modality used. The best are highlighted in **bold**.

Method	2-Phase (A+P)				3-Phase (A+P+D)			
	DSC(%) \uparrow	JSC(%) \uparrow	HD ₉₅ (mm) \downarrow	ASSD (mm) \downarrow	DSC(%) \uparrow	JSC(%) \uparrow	HD ₉₅ (mm) \downarrow	ASSD (mm) \downarrow
MAML	70.82 \downarrow _{5.27}	64.82 \downarrow _{6.59}	53.223 \uparrow _{24.38}	20.295 \uparrow _{4.64}	75.49 \downarrow _{3.16}	70.63 \downarrow _{4.18}	41.588 \uparrow _{14.78}	20.097 \uparrow _{9.36}
MW-UNet	71.38 \downarrow _{4.71}	65.47 \downarrow _{5.94}	37.519 \uparrow _{8.68}	19.013 \uparrow _{3.35}	76.21 \downarrow _{2.44}	71.56 \downarrow _{3.25}	35.221 \uparrow _{8.42}	18.640 \uparrow _{7.90}
SA-Net	71.91 \downarrow _{4.18}	66.14 \downarrow _{5.27}	36.946 \uparrow _{8.11}	18.870 \uparrow _{3.21}	76.29 \downarrow _{2.36}	71.67 \downarrow _{3.14}	34.126 \uparrow _{7.32}	17.444 \uparrow _{6.71}
PA-ResSeg	74.19 \downarrow _{1.90}	68.97 \downarrow _{2.44}	31.287 \uparrow _{2.45}	16.952 \uparrow _{1.29}	77.17 \downarrow _{1.48}	72.84 \downarrow _{1.97}	33.608 \uparrow _{6.80}	15.443 \uparrow _{4.71}
MCDA-Net	76.08 \downarrow _{0.01}	71.40 \downarrow _{0.01}	30.436 \uparrow _{1.60}	16.268 \uparrow _{0.61}	77.40 \downarrow _{1.25}	73.12 \downarrow _{1.69}	28.600 \uparrow _{1.79}	12.565 \uparrow _{1.83}
MADF-Net	76.09	71.41	28.8382	15.6590	78.65	74.81	26.8068	10.7366

Obs. ④: 3P Fusion: ART, PV, and DL. As listed in Table 4, the proposed MADF-Net achieved the state-of-the-art performance in terms of DSC (78.65%), JSC (74.81%), HD₉₅ (26.806), and ASSD (10.736) compared to the other five methods on the 3P fusion strategy. This indicates that our MADF-Net more accurately localizes the spatial positions and delineates the geometric shapes of the target regions. The quantitative comparison of the performance across 1P, 2P, and 3P fusion strategies is further illustrated in Figure 4, and the visualization of the segmentation results is shown in Figure 5. Compared to other methods, the proposed MADF-Net achieves the closest performance to the ground truth in terms of tumor contour localization and small-object boundary segmentation.

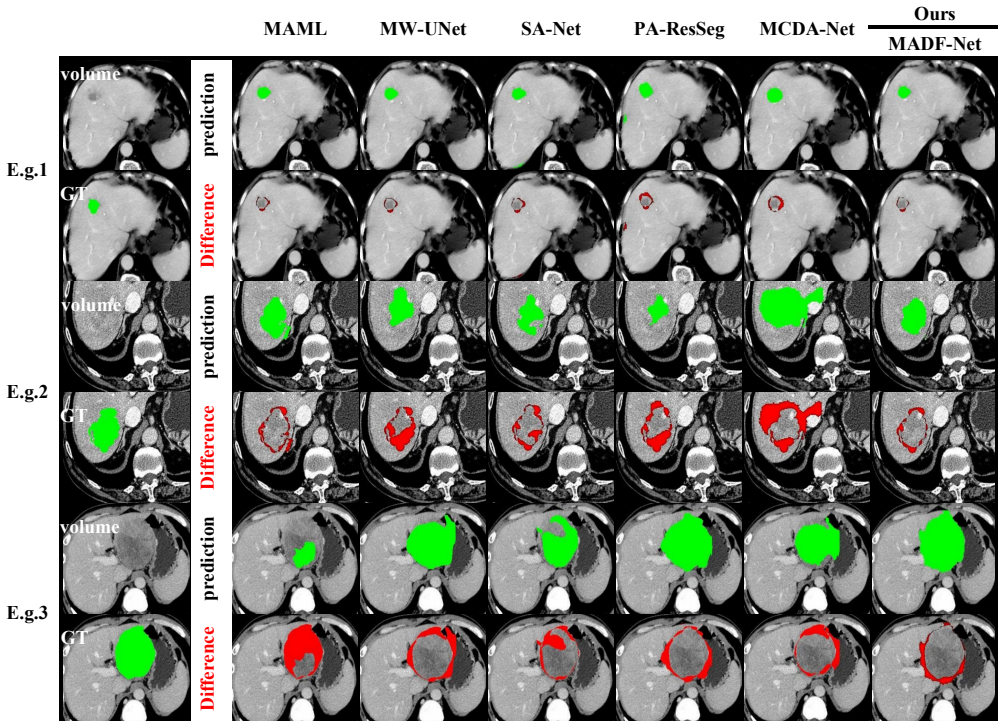


Figure 5: Result comparison of different three-phase networks. For better visualization, we performed appropriate cropping. The green region in the ground truth (GT) row represents the tumor, the green region in the prediction row indicates the predicted tumor area, and the red region in the difference row denotes the difference between the two.

5.2 EFFICIENCY ANALYSIS

Table 5 compares MADF-Net with baseline models. Our model requires 99.552×10^9 GFLOPs and 40.482 M parameters, achieving a favorable balance between computational cost and model size compared to SA-Net (152.965×10^9 GFLOPs, 170.852 M) and PA-ResSeg (64.660×10^9 GFLOPs, 67.732 M), while remaining competitive with lighter models such as MAML and MW-UNet. Further details on the experimental setup and efficiency comparisons are provided in Appendix C.

6 CONCLUSION AND FUTURE WORK

This paper presented MADF-Net, a novel multi-phase attention-based fusion network for liver tumor segmentation in contrast-enhanced CT images. Our approach is guided by a systematic quantitative evaluation of individual phases, confirming the predominant contribution of the PV phase and its alignment with clinical understanding. MADF-Net performs full-stage fusion across the input, feature, and decision levels to fully exploit the complementary information from ART, PV, and DL phases. The experiment on MPLL datasets demonstrate that our method achieves state-of-the-art performance and generalizes well across datasets. **Future Work:** We will (i) design phase-specific subnetworks and investigate phase-aware pretraining, and (ii) extend the paradigm to multi-modal/multi-omics fusion by integrating CT with digital pathology and radiomics-derived omics for comprehensive patient-level modeling.

REFERENCES

- 486
487
488 Abdullah F Al-Battal, Soan Duong, Van Ha Tang, Quang Duc Tran, Steven QH Truong, Chien Phan,
489 Truong Q Nguyen, and Cheolhong An. Multi-target and multi-stage liver lesion segmentation and
490 detection in multi-phase computed tomography scans. *arXiv preprint arXiv:2404.11152*, 2024.
- 491
492 Patrick Bilic, Patrick Christ, Hongwei Bran Li, Eugene Vorontsov, Avi Ben-Cohen, Georgios
493 Kaissis, Adi Szeskin, Colin Jacobs, Gabriel Efrain Humpire Mamani, Gabriel Chartrand, et al.
494 The liver tumor segmentation benchmark (lits). *Medical Image Analysis*, 84:102680, 2023.
- 495
496 Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xiaopeng Zhang, Qi Tian, and Manning Wang.
497 Swin-unet: Unet-like pure transformer for medical image segmentation. In *European Conference
on Computer Vision*, pp. 205–218. Springer, 2022.
- 498
499 Bingzhi Chen, Yishu Liu, Zheng Zhang, Guangming Lu, and Adams Wai Kin Kong. Transattunet:
500 Multi-level attention-guided u-net with transformer for medical image segmentation. *IEEE Trans-
actions on Emerging Topics in Computational Intelligence*, 8(1):55–68, 2023.
- 501
502 Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan Adeli, Yan Wang, Le Lu, Alan L Yuille,
503 and Yuyin Zhou. Transunet: Transformers make strong encoders for medical image segmentation.
504 *arXiv preprint arXiv:2102.04306*, 2021.
- 505
506 Yanling Chi, Jiayin Zhou, Sudhakar K Venkatesh, Qi Tian, and Jimin Liu. Content-based image
507 retrieval of multiphase ct images for focal liver lesion characterization. *Medical Physics*, 40(10):
508 103502, 2013.
- 509
510 Hao Du, Jiazheng Wang, Min Liu, Yaonan Wang, and Erik Meijering. Swinpa-net: Swin
511 transformer-based multiscale feature pyramid aggregation network for medical image segmen-
512 tation. *IEEE Transactions on Neural Networks and Learning Systems*, 35(4):5355–5366, 2022.
- 513
514 Ali Hatamizadeh, Yucheng Tang, Vishwesh Nath, Dong Yang, Andriy Myronenko, Bennett Land-
515 man, Holger R Roth, and Daguang Xu. Unetr: Transformers for 3d medical image segmentation.
516 In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp.
517 574–584, 2022.
- 518
519 Caner Hazirbas, Lingni Ma, Csaba Domokos, and Daniel Cremers. Fusetnet: Incorporating depth
520 into semantic segmentation via fusion-based cnn architecture. In *Asian Conference on Computer
521 Vision*, pp. 213–228. Springer, 2016.
- 522
523 Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected
524 convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern
525 recognition*, pp. 4700–4708, 2017.
- 526
527 Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. *Advances
528 in Neural Information Processing Systems*, 28, 2015.
- 529
530 Chunhui Jiang, Yi Wang, Qingni Yuan, Pengju Qu, and Heng Li. A 3d medical image segmentation
531 network based on gated attention blocks and dual-scale cross-attention mechanism. *Scientific
532 Reports*, 15(1):6159, 2025.
- 533
534 Linfeng Jiang, Jiajie Ou, Ruihua Liu, Yangyang Zou, Ting Xie, Hanguang Xiao, and Ting Bai.
535 Rmau-net: Residual multi-scale attention u-net for liver and tumor segmentation in ct images.
536 *Computers in Biology and Medicine*, 158:106838, 2023.
- 537
538 Xixi Jiang, Qingqing Luo, Zhiwei Wang, Tao Mei, Yu Wen, Xin Li, Kwang-Ting Cheng, and Xin
539 Yang. Multi-phase and multi-level selective feature fusion for automated pancreas segmen-
540 tation from ct images. In *International Conference on Medical Image Computing and Computer-
Assisted Intervention*, pp. 460–469. Springer, 2020.
- 541
542 Eunji Jun, Seungwoo Jeong, Da-Woon Heo, and Heung-Il Suk. Medical transformer: Universal
543 encoder for 3-d brain mri analysis. *IEEE Transactions on Neural Networks and Learning Systems*,
2023.

- 540 Haopeng Kuang, Xue Yang, Hongjun Li, Jingwei Wei, and Lihua Zhang. Adaptive multiphase liver
541 tumor segmentation with multiscale supervision. *IEEE Signal Processing Letters*, 31:426–430,
542 2024.
- 543 Naveen Kulkarni, Alice Fung, Avinash R Kambadakone, and Benjamin M Yeh. Ct techniques,
544 protocols, advancements and future directions in liver diseases. *Magnetic Resonance Imaging
545 Clinics of North America*, 29(3):305, 2021.
- 546 A Lam, D Fernando, CC Sirlin, M Nayyar, SC Goodwin, DK Imagawa, and C Lall. Value of
547 the portal venous phase in evaluation of treated hepatocellular carcinoma following transcatheter
548 arterial chemoembolisation. *Clinical Radiology*, 72(11):994–e9, 2017.
- 549 Xiaomeng Li, Hao Chen, Xiaojuan Qi, Qi Dou, Chi-Wing Fu, and Pheng-Ann Heng. H-denseunet:
550 hybrid densely connected unet for liver and tumor segmentation from ct volumes. *IEEE Transac-
551 tions on Medical Imaging*, 37(12):2663–2674, 2018.
- 552 Yin hao Liang, Wenjie Tang, Ting Wang, Wing WY Ng, Siyi Chen, Kuiming Jiang, Xinhua Wei,
553 Xinqing Jiang, and Yuan Guo. Hradnet: A hierarchical radiomics-based network for multicenter
554 breast cancer molecular subtypes prediction. *IEEE Transactions on Medical Imaging*, 43(3):
555 1225–1236, 2023.
- 556 Jae Hoon Lim, Dongil Choi, Seung Hoon Kim, Soon Jin Lee, Won Jae Lee, Hyo Keun Lim, and
557 Seonwoo Kim. Detection of hepatocellular carcinoma: value of adding delayed phase imaging to
558 dual-phase helical ct. *American Journal of Roentgenology*, 179(1):67–73, 2002.
- 559 Lihao Liu, Angelica I Aviles-Rivero, and Carola-Bibiane Schönlieb. Contrastive registration for
560 unsupervised medical image segmentation. *IEEE Transactions on Neural Networks and Learning
561 Systems*, 2023.
- 562 Zhenbing Liu, Junfeng Hou, Xipeng Pan, Ruojie Zhang, and Zhenwei Shi. Pa-net: A phase attention
563 network fusing venous and arterial phase features of ct images for liver tumor segmentation.
564 *Computer Methods and Programs in Biomedicine*, 244:107997, 2024.
- 565 Shuichi Monzawa, Tomoaki Ichikawa, Hiroto Nakajima, Yuki Kitanaka, Kosaku Omata, and Tsu-
566 tomu Araki. Dynamic ct for detecting small hepatocellular carcinoma: usefulness of delayed
567 phase imaging. *American Journal of Roentgenology*, 188(1):147–153, 2007.
- 568 Yangfan Ni, Geng Chen, Zhan Feng, Heng Cui, Dimitris Metaxas, Shaoting Zhang, and Wentao
569 Zhu. Da-tran: Multiphase liver tumor segmentation with a domain-adaptive transformer network.
570 *Pattern Recognition*, 149:110233, 2024.
- 571 Farid Ouhmich, Vincent Agnus, Vincent Noblet, Fabrice Heitz, and Patrick Pessaux. Liver tissue
572 segmentation in multiphase ct scans using cascaded convolutional neural networks. *International
573 Journal of Computer Assisted Radiology and Surgery*, 14:1275–1284, 2019.
- 574 Shaohua Qiao, Mengfan Xue, Yan Zuo, Jiannan Zheng, Haodong Jiang, Xiangai Zeng, and
575 Dongliang Peng. Four-phase ct lesion recognition based on multi-phase information fusion frame-
576 work and spatiotemporal prediction module. *BioMedical Engineering OnLine*, 23(1):103, 2024.
- 577 Ashwin Raju, Chi-Tung Cheng, Yuankai Huo, Jinzheng Cai, Junzhou Huang, Jing Xiao, Le Lu,
578 ChienHung Liao, and Adam P Harrison. Co-heterogeneous and adaptive segmentation from
579 multi-source and multi-phase ct imaging data: A study on pathological liver and lesion segmen-
580 tation. In *European Conference on Computer Vision*, pp. 448–465. Springer, 2020.
- 581 Sucheng Ren and Xiaomeng Li. Hresformer: Hybrid residual transformer for volumetric medical
582 image segmentation. *IEEE Transactions on Neural Networks and Learning Systems*, 2025.
- 583 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed-
584 ical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–
585 MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceed-
586 ings, part III 18*, pp. 234–241. Springer, 2015.

- 594 J Gabriel Schneider, Zhen J Wang, Wilbur Wang, Judy Yee, Yanjun Fu, and Benjamin M Yeh.
595 Patient-tailored scan delay for multiphase liver ct: improved scan quality and lesion conspicuity
596 with a novel timing bolus method. *American Journal of Roentgenology*, 202(2):318–323, 2014.
597
- 598 Abdelrahman Shaker, Muhammad Maaz, Hanoona Rasheed, Salman Khan, Ming-Hsuan Yang, and
599 Fahad Shahbaz Khan. Unetr++: delving into efficient and accurate 3d medical image segmenta-
600 tion. *IEEE Transactions on Medical Imaging*, 43(9):3377–3390, 2024.
- 601 Youyi Song, Jeremy Yuen-Chun Teoh, Kup-Sze Choi, and Jing Qin. Dynamic loss weighting for
602 multiorgan segmentation in medical images. *IEEE Transactions on Neural Networks and Learn-
603 ing Systems*, 35(8):10651–10662, 2023.
- 604 Changjian Sun, Shuxu Guo, Huimao Zhang, Jing Li, Meimei Chen, Shuzhi Ma, Lanyi Jin, Xiaoming
605 Liu, Xueyan Li, and Xiaohua Qian. Automatic segmentation of liver tumors from multiphase
606 contrast-enhanced ct images based on fcns. *Artificial Intelligence in Medicine*, 83:58–66, 2017.
607
- 608 Nikhil Kumar Tomar, Debesh Jha, Michael A Riegler, Håvard D Johansen, Dag Johansen, Jens
609 Rittscher, Pål Halvorsen, and Sharib Ali. Fanet: A feedback attention network for improved
610 biomedical image segmentation. *IEEE Transactions on Neural Networks and Learning Systems*,
611 34(11):9375–9388, 2022.
- 612 Bruce A Urban, Patricia A McGhie, and Elliot K Fishman. Helical ct: diagnostic pitfalls of arterial
613 phase imaging of the upper abdomen. *American Journal of Roentgenology*, 174(2):455–461,
614 2000.
- 615 Jeya Maria Jose Valanarasu, Vishwanath A Sindagi, Ilker Hacihaliloglu, and Vishal M Patel. Ki-
616 net: Overcomplete convolutional architectures for biomedical image and volumetric segmenta-
617 tion. *IEEE Transactions on Medical Imaging*, 41(4):965–976, 2021.
- 618 Jinting Wang, Yujiao Tang, Yang Xiao, Joey Tianyi Zhou, Zhiwen Fang, and Feng Yang. Grenet:
619 gradually recurrent network with curriculum learning for 2-d medical image segmentation. *IEEE
620 Transactions on Neural Networks and Learning Systems*, 35(7):10018–10032, 2023.
- 621 Yichao Wu, Qiang Zhou, Haoji Hu, Guanghua Rong, Yongwu Li, and Shiyan Wang. Hepatic lesion
622 segmentation by combining plain and contrast-enhanced ct images with modality weighted u-net.
623 In *2019 IEEE International Conference on Image Processing (ICIP)*, pp. 255–259. IEEE, 2019.
624
- 625 Yingying Xu, Ming Cai, Lanfen Lin, Yue Zhang, Hongjie Hu, Zhiyi Peng, Qiaowei Zhang, Qingqing
626 Chen, Xiongwei Mao, Yutaro Iwamoto, et al. Pa-resseg: A phase attention residual network for
627 liver tumor segmentation from multiphase ct images. *Medical Physics*, 48(7):3752–3766, 2021.
628
- 629 Fan Zhang, Huiying Liu, Qing Cai, Chun-Mei Feng, Binglu Wang, Shanshan Wang, Junyu Dong,
630 and David Zhang. Federated cross-incremental self-supervised learning for medical image seg-
631 mentation. *IEEE Transactions on Neural Networks and Learning Systems*, 2024.
632
- 633 Yao Zhang, Jiawei Yang, Jiang Tian, Zhongchao Shi, Cheng Zhong, Yang Zhang, and Zhiqiang
634 He. Modality-aware mutual learning for multi-modal medical image segmentation. In *Internat-
635 ional Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 589–
636 599. Springer, 2021a.
- 637 Yue Zhang, Chengtao Peng, Liying Peng, Huimin Huang, Ruofeng Tong, Lanfen Lin, Jingsong
638 Li, Yen-Wei Chen, Qingqing Chen, Hongjie Hu, et al. Multi-phase liver tumor segmentation
639 with spatial aggregation and uncertain region inpainting. In *Medical Image Computing and Com-
640 puter Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France,
641 September 27–October 1, 2021, Proceedings, Part I 24*, pp. 68–77. Springer, 2021b.
- 642 Yue Zhang, Chengtao Peng, Ruofeng Tong, Lanfen Lin, Yen-Wei Chen, Qingqing Chen, Hongjie
643 Hu, and S Kevin Zhou. Multi-modal tumor segmentation with deformable aggregation and un-
644 certain region inpainting. *IEEE Transactions on Medical Imaging*, 42(10):3091–3103, 2023.
645
- 646 Fuchen Zheng, Xinyi Chen, Xuhang Chen, Haolun Li, Xiaojiao Guo, Guoheng Huang, Chi-Man
647 Pun, and Shoujun Zhou. Assnet: Adaptive semantic segmentation network for microtumors and
multi-organ segmentation. *arXiv preprint arXiv:2409.07779*, 2024a.

648 Fuchen Zheng, Xinyi Chen, Xuhang Chen, Haolun Li, Xiaojiao Guo, Weihuang Liu, Chi-Man Pun,
649 and Shoujun Zhou. Affsegnet: Adaptive feature fusion segmentation network for microtumors
650 and multi-organ segmentation. *arXiv preprint arXiv:2409.07779*, 2024b.
651

652 W Zhong, F Liang, R Yang, and X Zhen. Prediction of microvascular invasion in hepatocellular
653 carcinoma based on multi-phase dynamic enhanced ct radiomics feature and multi-classifier hier-
654 archical fusion model. *Nan Fang yi ke da xue xue bao= Journal of Southern Medical University*,
655 44(2):260–269, 2024.

656 Yuyin Zhou, Yingwei Li, Zhishuai Zhang, Yan Wang, Angtian Wang, Elliot K Fishman, Alan L
657 Yuille, and Seyoun Park. Hyper-pairing network for multi-phase pancreatic ductal adenocarci-
658 noma segmentation. In *Medical Image Computing and Computer Assisted Intervention–MICCAI*
659 *2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part*
660 *II 22*, pp. 155–163. Springer, 2019.

661 Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++:
662 A nested u-net architecture for medical image segmentation. In *Deep Learning in Medical Image*
663 *Analysis and Multimodal Learning for Clinical Decision Support: 4th International Workshop,*
664 *DLMIA 2018, and 8th international workshop, ML-CDS 2018, held in conjunction with MICCAI*
665 *2018, Granada, Spain, September 20, 2018, proceedings 4*, pp. 3–11. Springer, 2018.
666

667 Xiner Zhu, Yichao Wu, Haoji Hu, Xianwei Zhuang, Jincan Yao, Di Ou, Wei Li, Mei Song, Na Feng,
668 and Dong Xu. Medical lesion segmentation by combining multimodal images with modality
669 weighted unet. *Medical Physics*, 49(6):3692–3704, 2022.
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A ALGORITHM.

The overall workflow of the **MADF-Net** is summarized in algorithm 1.

Algorithm 1: High-level pseudocode of the proposed **MADF-Net**.

```

702 A ALGORITHM.
703
704
705 The overall workflow of the MADF-Net is summarized in algorithm 1.
706
707 Algorithm 1: High-level pseudocode of the proposed MADF-Net.
708 Input: Phase images:  $I_A, I_B, I_C$ ; Clinical prior structure as a list of relations  $R$ ; Reference end
709 mask  $T$  (only for training)
710 Output: Segmentation output  $Y_{out}$ 
711 1 Initialize network parameters  $P$ ; set mode flag isTrain
712 2 foreach batch in training/eval do
713     /* basic preprocessing and branch init */
714     3  $I'_A \leftarrow \text{PrepData}(I_A)$ 
715     4  $I'_B \leftarrow \text{PrepData}(I_B)$ 
716     5  $I'_C \leftarrow \text{PrepData}(I_C)$ 
717     6  $I_m \leftarrow \text{CombineInit}(I'_A, I'_B, I'_C)$  // early merged input for main branch
718     7  $\text{BranchInputs} \leftarrow [I'_A, I'_B, I'_C, I_m]$ 
719     /* encode each branch to produce multi-level maps */
720     8 foreach entry  $J$  in  $\text{BranchInputs}$  do
721     9 |  $\text{FeatureMaps}[J] \leftarrow \text{EncodeBlock}(J)$  // returns list of maps at depths 1..D
722     end
723     /* clinical-aware message propagation (hierarchical depths) */
724     11 for depth  $d = 1$  to  $D$  do
725     12 | for node  $n$  in  $\text{RelationOrder}(R)$  do
726     13 | |  $\text{Parents} \leftarrow \text{GetParents}(n, R)$ 
727     14 | |  $\text{MsgIn} \leftarrow \text{PassMessages}([ \text{FeatureMaps}[\text{parent}][d] \text{ for parent in Parents } ])$ 
728     15 | |  $\text{UpdatedMap}[n][d] \leftarrow \text{UpdateUnit}(\text{FeatureMaps}[n][d], \text{MsgIn})$ 
729     16 | | // residual-style update
730     end
731     end
732     /* per-phase local refinement and prepare temporal stack */
733     18 for phase  $p$  in  $\{A, B, C\}$  do
734     19 |  $\text{LowFeat} \leftarrow \text{ExtractLow}(\text{UpdatedMap}[p])$ 
735     20 |  $\text{Refined}[p] \leftarrow \text{LocalRefine}(\text{LowFeat})$ 
736     21 |  $\text{StageOut}[p] \leftarrow \text{ProjectForTemporal}(\text{Refined}[p])$ 
737     end
738      $\text{TemporalStack} \leftarrow \text{Stack}(\text{StageOut}[A], \text{StageOut}[B], \text{StageOut}[C])$ 
739     /* per-pixel temporal attention */
740     24  $\text{TemporalEnhanced} \leftarrow \text{TemporalPerPixel}(\text{TemporalStack})$ 
741     /* neighbor-aware cross-temporal fusion per phase */
742     25 for phase  $p$  in  $\{A, B, C\}$  do
743     26 |  $\text{Attended}[p] \leftarrow \text{NeighborInteract}(\text{UpdatedMap}[p], \text{TemporalEnhanced})$ 
744     27 | // neighbor queries + relative-pos bias
745     28 |  $\text{Mix}[p] \leftarrow \text{BlendLinear}(\text{Attended}[p], \text{UpdatedMap}[p])$  // weighted linear mixing
746     29 |  $\text{FinalFeat}[p] \leftarrow \text{ChannelBoost}(\text{Mix}[p])$  // channel gating / enhancement
747     end
748     /* merge multi-phase features and decode */
749     30  $\text{MergedFeat} \leftarrow \text{FinalMerge}(\text{FinalFeat}[A], \text{FinalFeat}[B], \text{FinalFeat}[C],$ 
750     31  $\text{UpdatedMap}[I_m])$   $Y_{pred} \leftarrow \text{DecodeBlock}(\text{MergedFeat})$  // decoder with gated skips
751     /* auxiliary outputs aggregation (if enabled) */
752     32  $\text{AuxList} \leftarrow \text{GetAuxOutputs}()$  // possibly per-phase decoder heads
753     33  $Y_{out} \leftarrow \text{OutputMerge}(Y_{pred}, \text{AuxList})$ 
754     /* loss and update (training only) */
755     34 if isTrain then
756     35 |  $\text{LossVal} \leftarrow \text{CalcLoss}(Y_{out}, T)$ 
757     36 |  $\text{Backpropagate}(\text{LossVal}, P)$ 
758     end
759     return  $Y_{out}$ 

```

Given the three single-phase CT images, the network first performs input-level preprocessing and branch initialization, generating both phase-specific and early-fused representations. These features are then processed through hierarchical encoder blocks, where clinical-prior message propagation is applied to capture inter-phase dependencies. Next, local refinement modules enhance low-level cues, followed by per-pixel temporal attention to model cross-phase temporal correlations. Neighbor-aware cross-temporal fusion and channel enhancement further integrate complementary information before multi-phase features are merged and decoded. Finally, decision-level aggregation combines auxiliary and main-branch outputs to produce the final segmentation mask. During training, the predicted mask is supervised by the ground-truth labels via a composite loss function.

B THE MPLL DATASET.

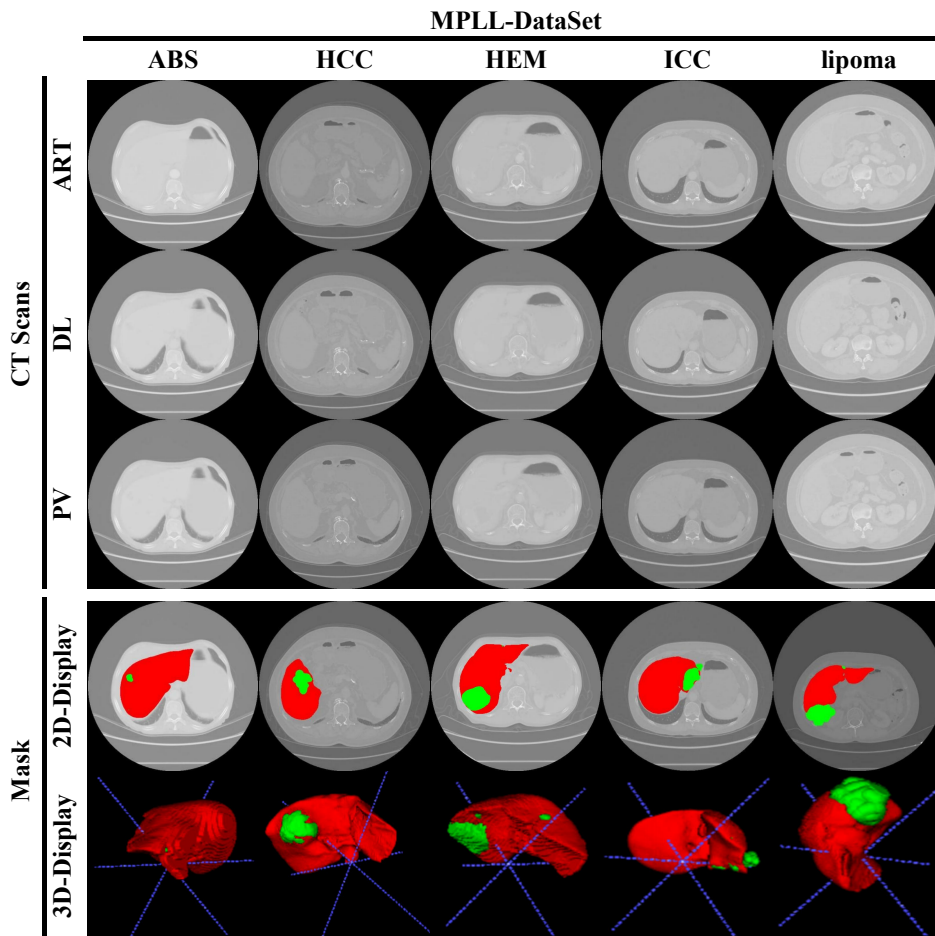


Figure 6: Example images from the MPLL dataset (red indicates liver regions, green indicates tumor regions).

Ⓛ: Patient Cohort and Imaging Protocol. The Multi-Phase Liver Lesion (MPLL) dataset was collected at the “Anonymous Authoritative Hospitals (information will be made public after the paper is accepted)”. The dataset comprises **952,601** 2D slices, making it one of the largest publicly reported multi-phase CT resources for liver tumor segmentation research, comprising 141 patients diagnosed with a wide spectrum of hepatic diseases. Imaging was performed between 2018 and 2022, covering both pediatric and adult populations (ages 9–72 years). All cases underwent standardized multi-phase contrast-enhanced CT examinations that included arterial, portal venous, and delayed phases, thereby capturing complementary hemodynamic information. Each scan was acquired at an in-plane resolution of 512×512 pixels, while slice thickness ranged from 0.62 mm to 5.0 mm. Due to differences in anatomical coverage, the number of slices varied considerably across

810 patients (48–777 slices per study). Data were de-identified before release, with ethical approval
811 obtained in advance.

812
813 **②: Clinical Diversity and Pathology Spectrum.** MPLL was intentionally designed to reflect
814 real-world clinical heterogeneity. It contains patients diagnosed with common malignant tumors
815 such as hepatocellular carcinoma (HCC) and intrahepatic cholangiocarcinoma, alongside a range
816 of benign lesions including cysts, hemangiomas, and abscesses. This diverse pathology coverage
817 ensures that the dataset does not disproportionately represent a single disease entity, but instead
818 provides a representative benchmark for developing algorithms that are robust across varying lesion
819 types, morphologies, and enhancement characteristics.

820
821 **③: Dataset Organization and Splitting Strategy.** To support reproducible research, the dataset
822 was partitioned into training, validation, and testing cohorts following a 7:1:2 split protocol consistent
823 with contemporary studies (Jiang et al., 2023). Importantly, the test set was fixed to 30 cases
824 and completely withheld during model design and training, thereby guaranteeing unbiased evaluation.
825 This design facilitates fair performance comparison across different methods and helps prevent
826 information leakage during algorithm development.

827
828 **④: Preprocessing and Annotation Pipeline.** One critical challenge of multi-phase imaging is the
829 misalignment across arterial, portal venous, and delayed acquisitions caused by respiration, patient
830 movement, or cardiac activity. To mitigate this, a B-spline deformable registration strategy was
831 applied using the portal venous phase as reference. This procedure significantly reduces inter-phase
832 variability and enables spatially consistent feature fusion. Ground-truth lesion masks were annotated
833 in ITK-SNAP by two board-certified radiologists, followed by an adjudication step by a senior
834 radiologist. This three-stage process was designed to maximize accuracy, reduce annotation bias,
835 and enhance inter-observer agreement.

836
837 **⑤: Comparative Advantages over Existing Datasets.** Unlike widely used liver CT datasets such
838 as LiTS2017 (Bilic et al., 2023) and Medical Segmentation Decathlon (Task 3: Liver), which primarily
839 focus on single-phase CT, MPLL offers multi-phase contrast-enhanced imaging across arterial,
840 portal venous, and delayed phases. This temporal richness provides unique opportunities for investigating
841 cross-phase fusion strategies, which are critical for accurate lesion delineation but are underexplored
842 in existing benchmarks. Moreover, MPLL is substantially larger in terms of slice count
843 (over 950k slices), contains a broader age range including pediatric cases, and offers a more diverse
844 pathology spectrum that includes both malignant and benign liver lesions. The dataset therefore not
845 only complements but also surpasses existing resources in its ability to support the development of
846 clinically relevant and generalizable liver lesion segmentation methods.

847
848 **⑥: Dataset Significance.** In summary, MPLL represents a large-scale, carefully curated, and clinically
849 diverse benchmark for multi-phase liver lesion segmentation. Its strengths lie in the combination
850 of temporal imaging information, broad pathology spectrum, rigorous preprocessing, and high-quality
851 expert annotations. Together, these characteristics make MPLL an invaluable resource
852 for advancing multi-phase fusion strategies in medical image analysis. Representative examples
853 highlighting inter-phase contrast variations and lesion depiction are illustrated in Figure 2.

854
855 Table 5: Efficiency Comparison of MADF-Net and Baseline Models (GFLOPs and Parameters). Performance is
856 evaluated on the MPLL dataset under the 3-phase experiment setting. The bold indicates the best.

Model	Gflops ($\times 10^9$)	Parameters (M)	Performance (%)
MAML (Zhang et al., 2021a)	23.802	4.216	75.49
MW-UNet (Zhu et al., 2022)	53.419	2.773	76.21
SA-Net (Zhang et al., 2021b)	152.965	170.852	76.29
PA-ResSeg (Xu et al., 2021)	64.660	67.732	77.17
MCDA-Net (Kuang et al., 2024)	89.480	48.717	77.40
Ours	99.552	40.482	78.65

C DETAILED EFFICIENCY ANALYSIS.

As shown in Table 5, **MADF-Net** demonstrates a favorable balance between computational complexity, model capacity, and segmentation accuracy. While lightweight models such as MW-UNet achieve a small parameter size (2.773 M), they still require a non-trivial computational cost (53.419 GFLOPs) and their performance (76.21%) remains noticeably lower than ours. Similarly, MAML achieves the lowest GFLOPs (23.802) but suffers from a relatively limited accuracy (75.49%), which constrains its clinical applicability. On the other hand, heavier architectures like SA-Net and PA-ResSeg demand extremely large computational budgets (up to 152.965 GFLOPs and 170.852 M parameters), yet the corresponding accuracy (76.29% and 77.17%, respectively) provides only marginal improvement over lightweight baselines.

MADF-Net maintains a moderate parameter count of 40.482 M and a competitive computational demand of 99.552 GFLOPs, while delivering the **highest segmentation accuracy** (78.65%) among all compared methods. This clearly illustrates the *efficiency-accuracy trade-off*: although **MADF-Net** is not the most lightweight in terms of FLOPs or parameters, it achieves the best performance, outperforming both lightweight and heavyweight counterparts. This balance highlights **MADF-Net**'s practicality for real-world clinical deployment, where both computational feasibility and reliable accuracy are crucial.

D LARGE LANGUAGE MODELS USAGE STATEMENT

LLMs were used only for language polishing in this work. The manuscript was drafted entirely by the authors, and LLMs were employed solely to refine grammar and clarity of English expression. All scientific ideas, methods, and results are original contributions of the human authors, with LLM assistance limited to post-writing editing akin to traditional proofreading.