# From Noisy Labels to Intrinsic Structure: A Geometric-Structural Dual-Guided Framework for Noise-Robust Medical Image Segmentation

Tao Wang[a,b,1], Zhenxuan Zhang[b,1], Yuanbo Zhou[a], Xinlin Zhang[a], Yuanbin Chen[a], Tao Tan[f], Guang Yang[b,c,d,e,2], Tong Tong[a,2,*]

[a]College of physics and information engineering, Fuzhou University, Xueyuan Road No.2, Fuzhou, 350108, China
[b]Bioengineering Department and Imperial-X, Imperial College London, London W12 7SL, UK
[c]National Heart and Lung Institute, Imperial College London, London SW7 2AZ, UK
[d]Cardiovascular Research Centre, Royal Brompton Hospital, London SW3 6NP, UK
[e]School of Biomedical Engineering & Imaging Sciences, King's College London, London WC2R 2LS, UK
[f]Faculty of Applied Science, Macao Polytechnic University, Macao Special Administrative Region of China

## ARTICLE INFO

## ABSTRACT

The effectiveness of convolutional neural networks in medical image segmentation relies on large-scale, high-quality annotations, which are costly and time-consuming to obtain. Even expert-labeled datasets inevitably contain noise arising from subjectivity and coarse delineations, which disrupt feature learning and adversely impact model performance. To address these challenges, this study proposes a Geometric–Structural Dual-Guided Network (GSD-Net), which integrates geometric and structural cues to improve robustness against noisy annotations. It incorporates a Geometric Distance-Aware module that dynamically adjusts pixel-level weights using geometric features, thereby strengthening supervision in reliable regions while suppressing noise. A Structure-Guided Label Refinement module further refines labels with structural priors, and a Knowledge Transfer module enriches supervision and improves sensitivity to local details. To comprehensively assess its effectiveness, we evaluated GSD-Net on six publicly available datasets: four containing three types of simulated label noise, and two with multi-expert annotations that reflect real-world subjectivity and labeling inconsistencies. Experimental results demonstrate that GSD-Net achieves state-of-the-art performance under noisy annotations, achieving improvements of 2.52% on Kvasir, 22.76% on Shenzhen, 8.87% on BU_SUC, and 4.59% on BraTS2020 under $S_R$ simulated noise. The codes of this study are available at https://github.com/ortonwang/GSD-Net.

## 1. Introduction

Medical image segmentation is a foundational technique in image analysis, which plays a critical role in disease diagnosis and clinical decision-making. Convolutional neural networks (CNNs) have been widely adopted for their strong feature extraction capabilities, with U-Net Ronneberger et al. (2015) and its variants becoming the leading architectures in medical image segmentation Liu et al. (2024); Kuang et al. (2025); Çiçek et al. (2016). However, the performance of these methods relies heavily on the availability of high-quality, precisely annotated pixel-level labels Shen et al. (2023). In practice, challenges such as limited image quality, indistinct object boundaries, and low contrast between foreground and background

---

*Corresponding author: ttraveltong@gmail.com
[1]These authors contributed equally to this work.
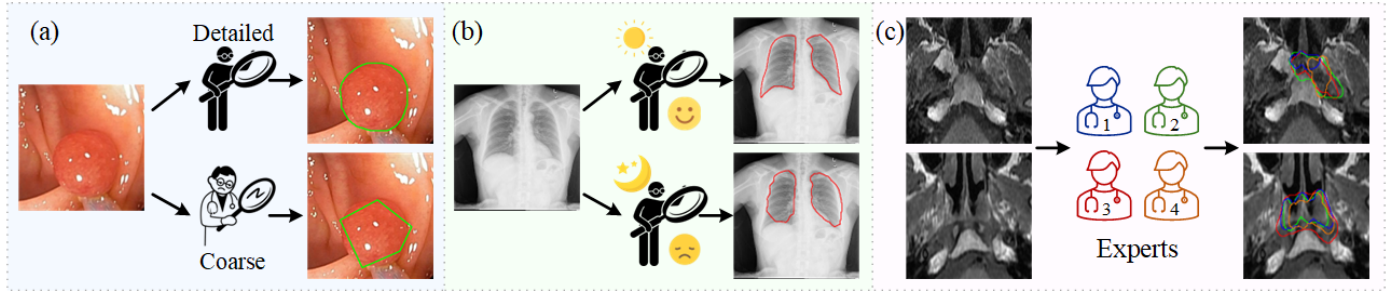[2]Co-last senior authors.

**Fig. 1.** Illustration of label noise sources: (a) coarse annotations, (b) intra-observer variability (inconsistencies by the same expert under different conditions), and (c) inter-observer variability (annotations from different experts shown in different colors).

Zhang et al. (2021) make precise delineation highly challenging. For instance, Computed Tomography (CT) images often suffer from motion artifact Ko et al. (2021), Magnetic Resonance Imaging (MRI) is prone to Rician noise and intensity inhomogeneity Mehrnia et al. (2024), while ultrasound commonly exhibits speckle noise and acoustic shadowing Singla et al. (2022). These modality-specific degradations further obscure anatomical boundaries and complicate accurate segmentation. As a result, reference annotations are typically provided by experienced clinicians through manual delineation.

However, this process is susceptible to label noise due to imprecise and biased annotations Ma et al. (2023). First, pixel-level labeling is costly and time-consuming. To save effort, annotators may produce coarse delineations Yu et al. (2020). As shown in Fig. 1 (a), a detailed annotation required 32 seconds, whereas a coarse one took only 5 seconds. Second, the inherent properties of medical images, including low contrast and indistinct boundaries, further increase the risk of boundary-localized errors. Third, annotation variability arises both within and between annotators. Intra-observer variability occurs when the same annotator produces inconsistent results due to timing, cognitive state, or emotional condition (Fig. 1 (b)). Inter-observer variability occurs when different experts provide inconsistent annotations due to subjectivity, fatigue, or inconsistent region definitions (Fig. 1 (c)) Wu et al. (2024); Schmidt et al. (2023). In large-scale projects, these variabilities accumulate and generate noisy labels which mislead the training process. Unlike classification noise, which is global and categorical, segmentation noise is spatially localized and structure-dependent. It often occurs near ambiguous boundaries or in anatomically complex regions Fang et al. (2023). Pixel-level noise further distorts spatial features and impairs the model's ability to capture object morphology and boundaries. As a result, both robustness and generalizability are diminished Yi et al. (2021). Therefore, mitigating label noise and improving robustness under noisy supervision remain critical challenges in medical image segmentation.

Extensive research has been conducted to address the performance degradation of deep learning models caused by noisy labels. Existing methods can be generally categorized into three categories: 1) Noise-robust loss functions, which modify the loss formulation to reduce the impact of mislabeled samples Barron (2019); Zhang and Sabuncu (2018). These approaches are simple to integrate and can effectively down-weight noisy data, but they rely heavily on loss statistics, which may misclassify hard examples as noise and fail to capture structural dependencies. 2) Sample selection, which identifies and prioritizes reliably annotated data Fang et al. (2023); Han et al. (2018). Such strategies mitigate the influence of noisy supervision by focusing on low-loss or high-confidence regions, yet they often discard ambiguous boundary pixels and underutilize supervision from unselected samples. and 3) Label correction, which improves noisy annotations using model predictions Yang et al. (2022). This paradigm can progressively refine noisy labels and recover useful supervision, but it is sensitive to early model errors and, without explicit structural constraints, may generate artifacts or discontinuous boundaries.

While these approaches have demonstrated strong performance in image classification, directly transferring them to segmentation is non-trivial. Pixel-level noise in segmentation is spatially structured and often concentrated near object boundaries, making it challenging for loss-based or sample-selection strategies to preserve reliable information. Moreover, label correction methods are prone to error accumulation and may yield anatomically inconsistent results, underscoring the need for segmentation-specific, structure-aware solutions. These limitations prevent them from providing reliable supervision under diverse and clinically realistic noise conditions. Despite recent advances, several critical issues remain: 1) Limited noisy-pixel detection. Many approaches rely solely on loss-based criteria (e.g., the small-loss strategy), which may discard correctly labeled yet ambiguous boundary pixels while retaining mislabeled regions with artificially low loss due to early overfitting Shi et al. (2024). 2) Insufficient anatomical modeling. Current methods typically rely on local smoothing or morphological constraints, failing to explicitly capture anatomical topology and regional coherence. This often results in fragmented boundaries, shape distortions, and spurious contours under severe noise. 3) Lack of cross-sample complementarity. Most strategies operate at the single-sample level and disregard structural information across images, thereby limiting their ability to recover reliable supervision. Additionally, many studies simulate label noise using simple morphological operations (e.g., erosion, dilation) Gonzalez-Jimenez et al. (2025); Li et al. (2021), which is overly simplistic and fails to capture the diverse, irregular annotations common in clinical practice. Taken together, these observations underscore the urgent need for noise-robust segmentation methods capable of delivering re-
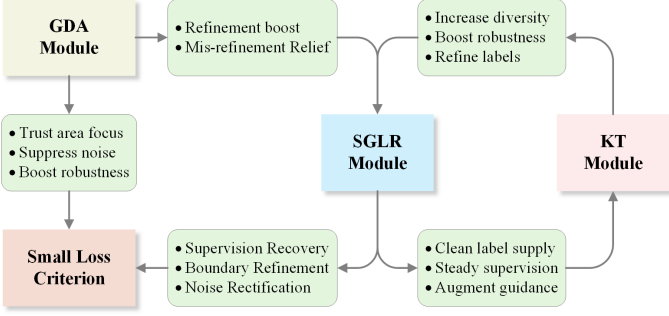
**Fig. 2.** Illustration of collaborative learning among modules. GDA: Geometric Distance-Aware; SGLR: Structure-Guided Label Refinement; KT: Knowledge Transfer.

liable performance under diverse and clinically realistic conditions.

To address the challenges of noisy annotations, we propose the Geometric-Structural Dual-Guided Network (GSD-Net), a unified framework that integrates multiple complementary modules in a forward-collaborative manner (Fig. 2). Specifically, GSD-Net first employs a dual-model co-regularization strategy and leverages the small-loss criterion to identify reliable regions, thereby providing a stable supervisory foundation for subsequent refinement. Building on this, the Geometric Distance-Aware (GDA) module reweights losses using geometric priors to further strengthen supervision in trustworthy regions and suppress noise-induced interference. Subsequently, the Structure-Guided Label Refinement (SGLR) module dynamically fuses the predictions of two networks under superpixel-based structural constraints to correct noisy labels and recover the supervisory signals skipped by the small-loss strategy. At the same time, the GDA module alleviates the negative effects of erroneous corrections, preventing the model from overfitting to erroneous pseudo-labels. Finally, the Knowledge Transfer (KT) module transfers local structural information between randomly paired images to enhance cross-sample diversity, thereby improving robustness and further facilitating more accurate label correction in the SGLR module. Through this progressive and synergistic design, GSD-Net establishes a positive interaction among reliability, label refinement, and cross-sample diversity, ultimately enabling robust segmentation under diverse and clinically realistic noise conditions. The main contributions of this study are summarized as follows:

- We propose GSD-Net, a unified framework for noise-robust medical image segmentation that holistically integrates reliability, structural guidance, and cross-sample diversity into a forward-collaborative design.

- This framework enhances robustness through a synergistic mechanism that consolidates reliable supervision with geometric priors, refines noisy regions via structural constraints, and enriches supervision by transferring diverse local patterns across images, with these processes working in concert within a unified framework.

- Extensive experiments on six public datasets, including simulated noise and real annotation noise caused by inter-

expert variability, demonstrate that our method achieves robust performance and effectively mitigates inter-expert inconsistencies.

## 2. Related Work

Numerous studies have demonstrated that label noise can severely degrade the generalization of deep neural networks, with high-parameter models such as CNNs and Transformers being especially susceptible to overfitting Zhang et al. (2020); Xu et al. (2022); Northcutt et al. (2021). To address this, researchers have developed various noise-robust learning approaches aimed at sustaining high segmentation performance despite noisy annotations. Existing work primarily focuses on the following aspects:

### 2.1. Noise-Robust Loss Functions

In noisy-label learning, loss function design is critical for ensuring model robustness and generalization. Standard cross-entropy (CE) loss assigns equal weight to all samples, which often causes overfitting to noisy labels Zhang et al. (2020); Northcutt et al. (2021). To address this, several noise-robust loss functions have been proposed. The Generalized Cross-Entropy (GCE) loss Zhang and Sabuncu (2018) uses a tunable parameter to interpolate between mean absolute error and CE loss, balancing robustness and convergence. The Symmetric Cross-Entropy loss Wang et al. (2019) combines CE with reverse cross-entropy to reduce degradation under both symmetric and asymmetric noise. The Dynamics-Aware Loss Li et al. (2023) adapts to the learning dynamics of deep networks by emphasizing easy examples in early training, progressively enhancing robustness, and incorporating a bootstrapping term. The multi-class unhinged loss and its smooth variants, SGCE and $\alpha$-MAE, stem from a decomposition theory of multi-class losses, enabling smooth transitions between unhinged loss and MAE for dynamic robustness control Paquin et al.. These approaches have shown strong performance in classification tasks, offering valuable insights for medical image segmentation. In segmentation, Gonzalez-Jimenez et al. Gonzalez-Jimenez et al. (2025) proposed T-Loss, a robust loss function derived from the Student-t distribution, which adaptively controls sensitivity to label noise and outliers via a learnable parameter. These losses have demonstrated effectiveness in classification tasks, where label noise is less spatially correlated. However, segmentation requires dense, pixel-level predictions and is more sensitive to structured and boundary-localized noise, limiting their effectiveness in this setting.

### 2.2. Reliable Sample Selection

Small-loss sample selection is a widely adopted strategy for noise-robust learning. During early training, models tend to fit correctly labeled "clean" samples first, while noisy samples typically yield higher loss values Arpit et al. (2017). Building on this idea, Han et al. Han et al. (2018) introduced the Co-Teaching framework, where two networks are jointly trained and exchange small-loss samples for mutual supervision. Wei et al. proposed JoCoR Wei et al. (2020), which adds

co-regularization to maintain prediction consistency alongside low-loss selection. Zhang et al. Zhang et al. (2020) extended Co-teaching to Tri-teaching, where three networks are trained simultaneously, with each pair collaboratively selecting reliable samples to guide the third. Fang et al. Fang et al. (2023) further proposed a collaborative learning framework where two models clean and distill reliable knowledge from each other using consistency-based regularization. Despite their effectiveness and structural simplicity, these methods rely on loss values as indicators for label correctness, may lead to the unintended removal of "high-loss clean samples" near ambiguous boundaries or complex structures. This not only restricts the model's ability to learn fine-grained details but also ignores potentially valuable image information present in noisy regions.

### 2.3. Label Correction

Label correction and pseudo-label generation aim to improve supervision by refining suspected label errors during training. Unlike small-loss strategies that focus exclusively on assumed clean samples, these methods seek to exploit informative content from noisy regions, typically by leveraging high-confidence model predictions. For example, Xiao et al. proposed ProMix Xiao et al. (2022), which employs matched high-confidence selection to progressively expand the clean sample set. Qiu et al. Qiu et al. (2023) utilized a multimodal self-training framework to address label inconsistencies between Whole Slide Images and their patches. Liu et al. proposed ADELE Liu et al. (2022), which exploits early-learning dynamics to adaptively correct class-wise noise while enforcing multi-scale consistency. Shi et al. Shi and Wu (2021) introduced adaptive thresholding with prototype-guided correction for heavily corrupted subsets; and Jin et al. Jin et al. (2022) proposed pixel-level correction through noisy pixel estimation. Nevertheless, the predominant reliance on model-predicted confidence makes these methods vulnerable to biased predictions under severe noise, potentially leading to iterative reinforcement of incorrect labels. To address this limitation, our approach integrates confidence estimation with superpixel-based structural priors, thereby constraining correction within structurally coherent regions and reducing the risk of error propagation.

## 3. Methods

This study aims to address the challenge of inaccurate pixel-level annotations in medical image segmentation. To facilitate understanding, we first present an overview of the proposed framework, followed by detailed descriptions of its components. Key notations are summarized in Table 1 for clarity.

### 3.1. Overview

The overall architecture and pseudo-code of GSD-Net are presented in Fig. 3 and Algorithm 1. Given an input pair $(x_1, x_2)$, weak augmentations produce $(x'_1, x'_2)$ with predictions obtained as $p_1 = f_{\theta_1}(x_1)$, $p_2 = f_{\theta_2}(x_2)$, $p'_1 = f_{\theta_1}(x'_1)$, and

**Table 1. Descriptions of Key Notations**

| Notations | Descriptions |
|---|---|
| $x, nGT$ | Input image and its noisy annotation |
| $\mathcal{A}(\cdot)$ | Weak augmentation operation |
| $x'$ | Augmented image: $x' = \mathcal{A}(x)$ |
| $f_{\theta_i}(\cdot)$ | Model with parameters $\theta_i$ |
| $p_i$ | Prediction from model $i$: $p_i = f_{\theta_i}(x)$ |
| $\mathbb{D}$ | Pixel coordinate set of $x$ |
| $\mathbb{D}^{clean}$ | Selected small-loss subset |
| $\oplus$ | Element-wise addition |
| $\odot$ | Element-wise multiplication |
| $\mathcal{S}$ | Generated superpixel map: $\mathcal{S} = \text{SLIC}(x)$ |
| $y$ | Generated pseudo-label: $y = \text{SGLR}(x)$ |
| $\mathcal{W}$ | Weight map generated by the GDA module |
| $\mathcal{KL}$ | Symmetric Kullback−Leibler divergence |
| $x_{1\rightarrow2}$ | Knowledge transferred from $x_1$ to $x_2$ |
| $x_{2\rightarrow1}$ | Knowledge transferred from $x_2$ to $x_1$ |

$p'_2 = f_{\theta_2}(x'_2)$. A robust co-regularization mechanism is subsequently employed to optimize the model parameters. We first apply a small-loss strategy to identify reliably annotated regions in the simulated noisy ground truth ($nGT$). Then the GDA module suppresses potential noise and reinforces supervision in trustworthy areas. Next, a Structure-Guided Label Refinement module leverages superpixel-based spatial priors to dynamically weight predictions from two networks, thereby refine $nGT$. Finally, the Knowledge Transfer module enriches data diversity and improves the model's sensitivity to local details and structural variations across samples.

### 3.2. Preliminary research

JoCoR Wei et al. (2020) employs the co-regularization framework to clean up the training data. This method primarily incorporates the small-loss strategy with a contrastive loss and serves as a baseline in our study. The implementation begins by evaluating the loss between predictions and $nGT$, as formulated:

$$\mathcal{L}_{CE}(p_{[i,j]}, y_{[i,j]}) = -\sum_{c=1}^{C} y_{[i,j,c]} \log(p_{[i,j,c]}), \tag{1}$$

where $y_{[i,j,c]}$ indicates whether the pixel at position $[i, j]$ belongs to the $c$-th class, and $C$ denotes the total number of classes. Similarly, $p_{[i,j,c]}$ represents the probabilities of the pixel at position $(i, j)$ belonging to the $c$-th class. Therefore, the supervised loss is defined as:

$$\mathcal{L}_{sup}(x_{[i,j]}) = \mathcal{L}_{CE}\left\{f_{\theta_1}(x_{[i,j]}), nGT_{[i,j]}\right\} + \\ \mathcal{L}_{CE}\left\{f_{\theta_2}(x'_{[i,j]}), nGT_{[i,j]}\right\}, \tag{2}$$

According to the agreement maximization principle Sindhwani et al. (2005), models tend to agree more on correctly labeled regions than on mislabeled ones. JoCoR leverages a contrastive term for co-regularization, where agreement is quantified using
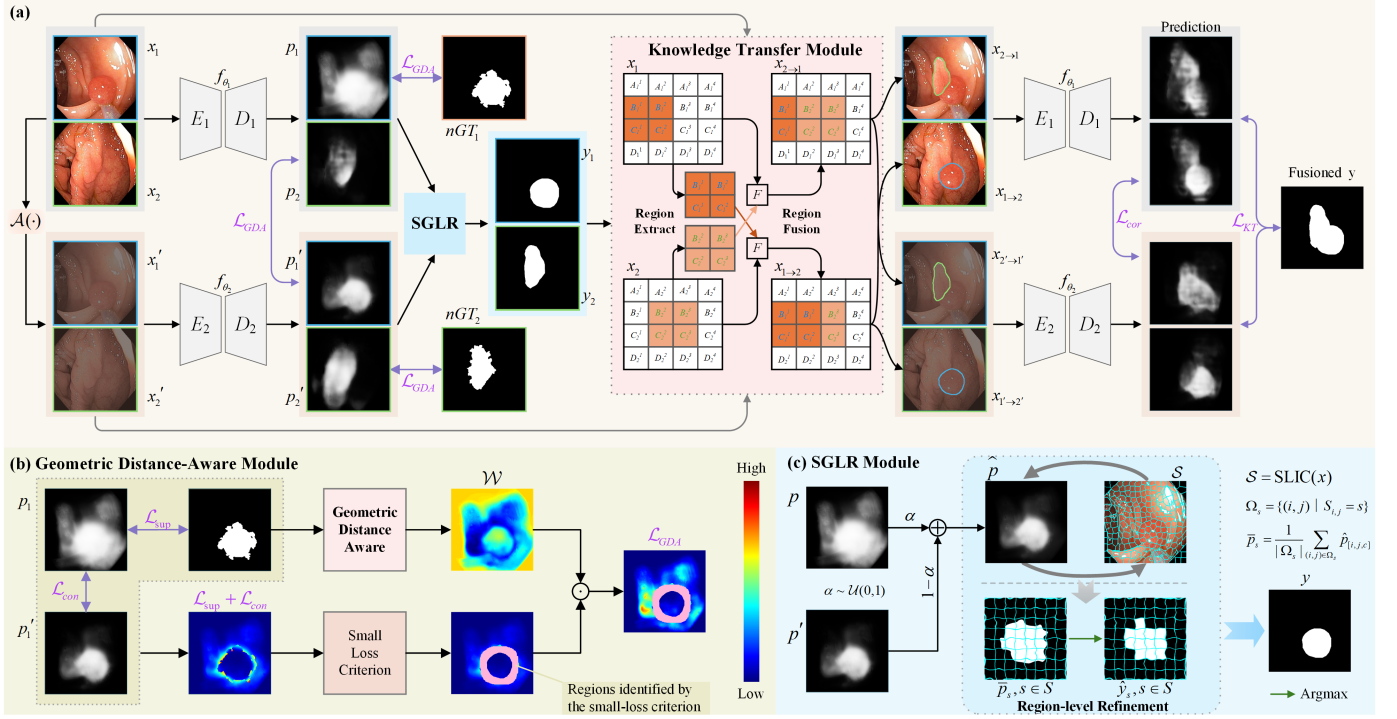
Fig. 3. Schematic diagram of the proposed GSD-Net framework: (a) overall workflow, (b) Geometric Distance-Aware module, and (c) Structure-Guided Label Refinement module.

the symmetric Kullback-Leibler ($\mathcal{KL}$) divergence, defined as:

$$\mathcal{KL}(P_{[i,j]} \| Q_{[i,j]}) = \sum_{c=1}^{C} P_{[i,j,c]} \log \left( \frac{P_{[i,j,c]}}{Q_{[i,j,c]}} \right), \quad (3)$$

Therefore, the agreement between the two networks is defined as:

$$\mathcal{L}_{con}(x_{[i,j]}) = \mathcal{KL}\left\{ f_{\theta_1}(x_{[i,j]}) \| f_{\theta_2}(x'_{[i,j]}) \right\} + \\ \mathcal{KL}\left\{ f_{\theta_2}(x'_{[i,j]}) \| f_{\theta_1}(x_{[i,j]}) \right\}, \quad (4)$$

Next, JoCoR selects "clean" regions $\mathbb{D}^{clean}$ using the "small-loss" criterion:

$$\mathbb{D}^{clean} = \arg \min_{\mathbb{D}': |\mathbb{D}'| \geq \mathcal{R}(e)|\mathbb{D}|} \sum_{(i,j) \in \mathbb{D}} \mathcal{L}_{sup}(x_{[i,j]}) \oplus \mathcal{L}_{con}(x_{[i,j]}), \quad (5)$$

The retention rate is defined as $\mathcal{R}(e) = 1 - \min(\frac{e}{10}\tau, \tau)$, where $e$ denotes the current epoch and $\tau$ is a constant. At the beginning, $\mathcal{R}(e)$ is close to 1 and it decreases toward $1 - \tau$ as training progresses, gradually reducing the selected regions to mitigate overfitting to $nGT$ Han et al. (2018). Therefore, the following loss for JoCoR is defined as:

$$\mathcal{L}_{JoCoR}(x_{[i,j]}) = \frac{\sum_{(i,j) \in \mathbb{D}^{clean}} \mathcal{L}_{sup}(x_{[i,j]}) \oplus \mathcal{L}_{con}(x_{[i,j]})}{|\mathbb{D}^{clean}|}, \quad (6)$$

### 3.3. Geometric Distance-Aware Module

Although the small-loss criterion is widely adopted for retaining clean regions, it may mistakenly regard high-confidence noisy regions as clean and overlook mislabeled pixels near decision boundaries. To reduce the impact of misidentification, we

propose the GDA Module, as illustrated in the Fig. 3 (a) and Fig. 4, which dynamically adjusts loss weights during training. Since annotation errors frequently occur near lesion boundaries Lee et al. (2020), we define a geometric-aware weight map $\mathcal{W}$ that assigns greater weights to pixels farther from the boundaries, which are less likely to be mislabeled. Specifically, we first extract the boundaries $\partial_{nGT}$ across all categories of $nGT$, formulated as:

$$\partial_{nGT} = \{(i,j) \in \mathbb{D} \mid \exists (u,v) \in \mathcal{N}(i,j), \, nGT(u,v) \neq nGT(i,j)\}, \quad (7)$$

where $\mathbb{D}$ denotes the set of all pixel coordinates in the image $x$. $\mathcal{N}(i,j)$ represents the connected neighborhood of pixel $(i,j)$. For each pixel, its distance to the boundary $\partial_{nGT}$ is mapped to a corresponding weight, as formulated:

$$\mathcal{W}_{raw}(x_{[i,j]}) = \min_{(m,n) \in \partial_{nGT}} \sqrt{(i-m)^2 + (j-n)^2}, \quad (8)$$

Here, $\mathcal{W}_{raw}$ assigns smaller weights to pixels near boundaries and larger weights to those farther away. To prevent overemphasis on distant pixels, which could lead to overfitting or gradient instability, we adopt a conservative weight-clipping strategy:

$$\mathcal{W}(x_{[i,j]}; e) = \max \left\langle \min \left\{ \mathcal{W}_{raw}(x_{[i,j]}), T \right\} - \frac{e}{E} \times T, \, 1 \right\rangle, \quad (9)$$

where $T$ represents the threshold for $\mathcal{W}_{raw}$, $e$ denotes the current training epoch, and $E$ denotes the maximum training epoch. The refined weight map $\mathcal{W}$ is then integrated into the
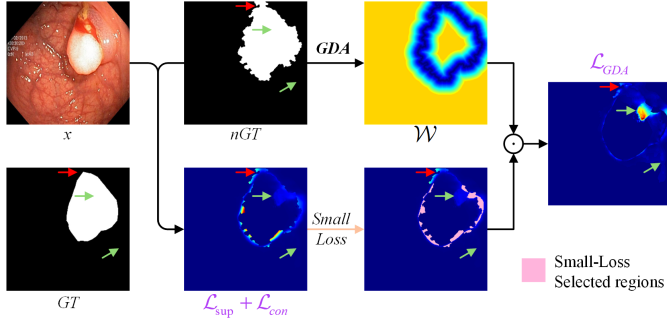
**Fig. 4.** The schematic diagram of the Geometric Distance-Aware Module.

supervision loss to form the Geometric Distance-Aware Loss $\mathcal{L}_{GDA}$:

$$\mathcal{L}_{GDA} = \frac{1}{|\mathbb{D}^{clean}|} \sum_{(i,j)\in\mathbb{D}^{clean}} \left\{ \mathcal{L}_{sup}(x_{[i,j]}) \oplus \mathcal{L}_{con}(x_{[i,j]}) \right\} \odot \mathcal{W}. \quad (10)$$

The module enhances supervision in reliably labeled regions (green arrows in Fig. 4) while attenuating the effect of noisy regions (red arrows in Fig. 4). This is particularly beneficial in the early training stages, where it reduces the influence of label noise and improves both stability and robustness.

### 3.4. Structure-Guided Label Refinement Module

Although small-loss strategies and GDA modules suppress noisy supervision, they inevitably overlook certain regions whose correction could improve robustness. To address this, we propose the SGLR Module (Fig. 3 (c)), which integrates predictions from two networks to enhance pseudo-label accuracy and diversity. By enhancing pseudo-label diversity and preserving model independence, the SGLR module mitigates performance degradation from over-enforcing prediction consistency (Eq. (3)). Specifically, a random coefficient $\alpha \sim \mathcal{U}(0, 1)$ is used to linearly combine the softmax outputs of the two segmentation branches ($f_{\theta_1}, f_{\theta_2}$), formulated as:

$$\hat{p}_{[i,j,c]} = \left\langle \alpha \cdot \sigma \left\{ f_{\theta_1}(x_{[i,j]}) \right\} \right\rangle \oplus \left\langle (1-\alpha) \cdot \sigma \left\{ f_{\theta_2}(x'_{[i,j]}) \right\} \right\rangle, \quad (11)$$

where $\sigma\{\cdot\}$ denotes the softmax operation. To improve pseudo-label reliability, we incorporate structural priors extracted using the SLIC superpixel algorithm Achanta et al. (2012), i.e., $\mathcal{S} = $ SLIC$(x)$, configured with $n_{\text{segments}} = 1024$ and compactness $= 10$, common settings in medical image segmentation. This configuration balances boundary adherence and computational efficiency while preserving anatomically coherent superpixel shapes. For each superpixel region, let $\Omega_s = \{(i, j) \mid S_{i,j} = s\}$ denote its pixel set. The confidence of superpixel $s$ is computed as:

$$\overline{p}_s = \frac{1}{|\Omega_s|} \sum_{(i,j)\in\Omega_s} \hat{p}_{[i,j,c]}, \quad (12)$$

and assign pseudo-labels by selecting the most probable class:

$$\hat{y}_{[i,j]} = \arg\max_c \hat{p}_{[i,j,c]}, \quad \text{for } (i, j) \in \Omega_s. \quad (13)$$

---

**Algorithm 1** A Geometric-Structural Dual-Guided Framework for Noise-Robust Medical Image Segmentation (Training)

---
1: **Input:** Training set $\{(x_k, nGT_k)\}_{k=1}^n$
2: **Initialize:** Two networks $f_{\theta_1}(\cdot)$ and $f_{\theta_2}(\cdot)$
3: **for** epoch = 1 to MaxEpoch **do**
4:     $x' \leftarrow \mathcal{A}(x)$
5:     $p \leftarrow f_{\theta_1}(x), \quad p' \leftarrow f_{\theta_2}(x')$
6:     Compute $\mathcal{L}_{sup}$ (Eq. (2)) and $\mathcal{L}_{con}$ (Eq. (4))
7:     $\mathbb{D}^{clean} \leftarrow$ Small-loss Criterion ($\mathcal{L}_{sup} \oplus \mathcal{L}_{con}$)(Eq. (5))
8:     $\mathcal{W} \leftarrow$ GDA($nGT$, epoch) (Eq. (9))
9:     // Reweighting Loss $\mathcal{L}_{GDA}$:
10:     $\mathcal{L}_{GDA} \leftarrow \dfrac{\sum_{(i,j)\in\mathbb{D}^{clean}} \left( \mathcal{L}_{sup} + \mathcal{L}_{con} \right) \cdot \mathcal{W}}{|\mathbb{D}^{clean}|}$
11:     // Structure-Guided Label Refinement:
12:     $\mathcal{S} \leftarrow$ SLIC$(x)$
13:     $y \leftarrow$ SGLR($p, p', \mathcal{S}, \mathbb{D}^{clean}$) (Sec. 3.4)
14:     // Knowledge Transfer:
15:     Sample $x_1, x_2 \sim \{x_k\}_{k=1}^n$
16:     $x_{2\rightarrow1}, x_{1\rightarrow2}, y_{2\rightarrow1}, y_{1\rightarrow2} \leftarrow$ KT($x_1, x_2, y_1, y_2$) (Sec. 3.5)
17:     $p_{2\rightarrow1} \leftarrow f_{\theta_1}(x_{2\rightarrow1}), \quad p'_{2\rightarrow1} \leftarrow f_{\theta_1}(x'_{2\rightarrow1})$
18:     $p_{1\rightarrow2} \leftarrow f_{\theta_2}(x_{1\rightarrow2}), \quad p'_{1\rightarrow2} \leftarrow f_{\theta_2}(x'_{1\rightarrow2})$
19:     Compute $\mathcal{L}_{KT}$ (Eq. (16)) and $\mathcal{L}_{cor}$ (Eq. (17))
20:     // Total Loss Function:
21:     $\mathcal{L}_{total} \leftarrow \mathcal{L}_{GDA} + \mathcal{L}_{KT} + \mathcal{L}_{cor}$
22:     Update parameters $\theta_1, \theta_2$ using $\mathcal{L}_{total}$
23: **end for**
24: **Return:** Trained models $f_{\theta_1}(\cdot), f_{\theta_2}(\cdot)$

---

The final pseudo-labels used for training are defined as:

$$y = (\mathbb{D}^{clean} \odot nGT) \oplus (\overline{\mathbb{D}^{clean}} \odot \hat{y}), \quad (14)$$

where $\mathbb{D}^{clean}$ denotes the clean regions as defined by Eq. (5), and $\overline{\mathbb{D}^{clean}}$ represents their complement, i.e., $\overline{\mathbb{D}^{clean}} = 1 - \mathbb{D}^{clean}$.

### 3.5. Knowledge Transfer Module

To enhance diversity and sensitivity to local structures, we propose a Knowledge Transfer (KT) module. By doing so, the module strengthens robustness and supports noisy label correction in the SGLR module. This module extracts local regions from randomly paired images and embeds them into their counterparts, enabling semantic-level integration. This strategy provides consistent and diverse supervision across different spatial contexts, enabling the model to capture fine-grained details and structural variability, and thus to learn more robust and generalizable feature representations. Let $\mathcal{M}_1$ and $\mathcal{M}_2$ denote the masks of regions extracted from $x_1$ and $x_2$, respectively. Each mask is generated by randomly sampling a foreground or background region, with the sampling probability determined by the proportion of the target foreground area in the image. The fusion process is formulated as:

$$x_{i\rightarrow j} = x_j \odot (1 - \mathcal{M}_i) + x_i \odot \mathcal{M}_i, \quad i, j \in \{1, 2\} \quad i \neq j,$$
$$\mathcal{W}_{i\rightarrow j} = \mathcal{W}_j \odot (1 - \mathcal{M}_i) + \mathcal{W}_i \odot \mathcal{M}_i, \quad i, j \in \{1, 2\} \quad i \neq j, \quad (15)$$

the fused images are fed into the networks, with supervision provided by the corresponding fused pseudo-labels generated with the same fusion strategy as in Eq. (15). Accordingly, the supervision after the KT module can be expressed as:

$$\mathcal{L}_{KT} = \sum_{(x,y,\mathcal{W}) \in \mathcal{P}} \{ \mathcal{L}_{CE} [f_{\theta_1}(x), y] \odot \mathcal{W} + \mathcal{L}_{DC} [f_{\theta_1}(x), y] \}$$
$$+ \sum_{(x,y,\mathcal{W}) \in \mathcal{P}'} \{ \mathcal{L}_{CE} [f_{\theta_2}(x), y] \odot \mathcal{W} + \mathcal{L}_{DC} [f_{\theta_2}(x), y] \}, \tag{16}$$

where the pairwise transformation sets are defined as follows:

$$\mathcal{P} = [(x_{1 \to 2}, y_{1 \to 2}, \mathcal{W}_{1 \to 2}), (x_{2 \to 1}, y_{2 \to 1}, \mathcal{W}_{2 \to 1})],$$
$$\mathcal{P}' = [(x_{1' \to 2'}, y_{1' \to 2'}, \mathcal{W}_{1 \to 2}), (x_{2' \to 1'}, y_{2' \to 1'}, \mathcal{W}_{2 \to 1})].$$

We further incorporate a co-regularization mechanism to encourage consistency between networks, as formulated:

$$\mathcal{L}_{cor} = \mathcal{KL} [f_{\theta_1}(x_{2 \to 1}) \| f_{\theta_2}(x_{2' \to 1'})] \cdot \mathcal{W}_{2 \to 1} +$$
$$\mathcal{KL} [f_{\theta_1}(x_{1 \to 2}) \| f_{\theta_2}(x_{1' \to 2'})] \cdot \mathcal{W}_{1 \to 2}, \tag{17}$$

this term enforces agreement on transferred samples while accounting for geometric-aware weights.

### 3.6. Total Loss Function

For end-to-end optimization, we integrate all module-specific losses into a unified objective, defined as:

$$\mathcal{L}_{total} = \mathcal{L}_{GDA} + \mathcal{L}_{KT} + \mathcal{L}_{cor}, \tag{18}$$

where $\mathcal{L}_{GDA}$, $\mathcal{L}_{KT}$, and $\mathcal{L}_{cor}$ are defined in Eq. (10), (16), and (17), respectively.

## 4. Experiments and Results

### 4.1. Datasets

We evaluate our method on six publicly available datasets against state-of-the-art approaches. We simulate label noise on Kvasir Jha et al. (2019), Shenzhen Candemir et al. (2013); Jaeger et al. (2013); Stirenko et al. (2018), BU_SUC Iqbal and Sharif (2024)[3], and BraTS2020 Menze et al. (2014) datasets, which all provide precise pixel-level annotations for controlled noise simulation and robustness evaluation. As widely recognized and representative benchmarks in medical image segmentation, they span diverse imaging modalities such as endoscopy, X-radiation (X-ray), ultrasound (US), and MRI, enabling a comprehensive assessment of the method's applicability and generalizability across different clinical contexts. To evaluate robustness under real-world noise, we further evaluate on the LIDC Armato III et al. (2011) and MMIS-2024 Luo et al. (2023); Wu et al. (2024); Bakas et al. (2021); Cepeda et al. (2023); Suter et al. (2022) datasets, where each image is annotated by multiple independent experts. These datasets capture label noise arising from inter-observer variability (Fig. 5). The details of each dataset are summarized in Table 2.

**Table 2. Summary of the datasets used in our experiments**

| Datasets | Modalities | Characteristics | Target | Train | Test | Size |
|---|---|---|---|---|---|---|
| Kvasir | Endoscopy | • Uneven illumination and glare<br>• Blurry tissue boundaries | Gastrointestinal polyp | 700 | 300 | $256^2$ |
| Shenzhen | X-Ray | • Foreground−background overlap<br>• Low soft-tissue contrast | Lungs area | 453 | 113 | $256^2$ |
| BU_SUC | US | • High noise and low contrast<br>• Heterogeneous echogenicity | Breast tumor | 568 | 243 | $256^2$ |
| BraTS2020 [1] | MRI | • Heterogeneous tissue<br>• Irregular tumor shapes | Glioma tumor | 296 | 73 | $128^3$ |
| LIDC [2] | CT | • Diverse nodule sizes and shapes<br>• Adhesion to surrounding vessels | Pulmonary nodules | 1287 | 322 | $128^2$ |
| MMIS-2024 [3] | MRI | • Complex anatomical structures<br>• Infiltration into adjacent tissues | Nasopharyngeal carcinoma | 100 | 20 | $128^2$ |

[1] Experiments use T1, T2, T1ce, and FLAIR sequences, with all tumor subtypes merged into a binary mask.
[2] Following Kohl et al. (2018); Wu et al. (2024), 1,609 patches, each annotated by four radiologists were selected.
[3] 100 volumes (2405 slices) were used for training and 20 volumes for testing; inference was performed slice-wise and subsequently reconstructed into 3D volumes for evaluation following Wu et al. (2024).

### 4.2. Implementation Details

The proposed framework was implemented using PyTorch and evaluated on an NVIDIA RTX 4090 GPU with 24 GB of memory. For 2D segmentation tasks, we used U-Net Ronneberger et al. (2015) as the backbone with a batch size of 16. For BraTS2020, we adopted 3D U-Net Çiçek et al. (2016) with a batch size of 4. All models were trained for 100 epochs using stochastic gradient descent (SGD) optimizer with a learning rate of $5 \times 10^{-3}$ and a weight decay of $1 \times 10^{-5}$. The maximum scaling factor $T$ in the GDA module was configured based on dataset resolution, with a value of 10 for 256×256 images and 5 for 128×128 images. In the small-loss strategy, the constant $\tau$ for the simulated datasets was set according to the actual noise rate estimated from the ground truth in simulated datasets. For the LIDC and MMIS-2024 datasets, $\tau$ is determined by the 15% of the foreground proportion, resulting in 0.003 for LIDC and 0.008 for MMIS-2024.

### 4.3. Noise simulation

Many studies use morphological operations (e.g., dilation or erosion) to generate simulated label noise, which we denote as $S_{DE}$. However, such noise is typically spatially uniform and structurally oversimplified, failing to capture the irregular, boundary-focused characteristics of real annotation errors. To better approximate realistic label noise, we employ a Markov-based boundary perturbation strategy Yao et al. (2023), which randomly distorts ground-truth contours. Building on this, we introduce two additional noise types beyond the commonly used $S_{DE}$: (1) foreground-reducing noise ($S_R$), simulating under-annotation; and (2) foreground-expanding noise ($S_E$), simulating over-annotation. These noise patterns capture irregular, boundary-focused errors often observed in real annotations and enable a more rigorous evaluation of model robustness under both under- and over-annotation scenarios. Examples of the generated noisy labels are illustrated in the upper panel of Fig. 5. For the simulated dataset, model weights from the final 10 epochs were used for evaluation, and the mean performance along with the standard deviation was reported.
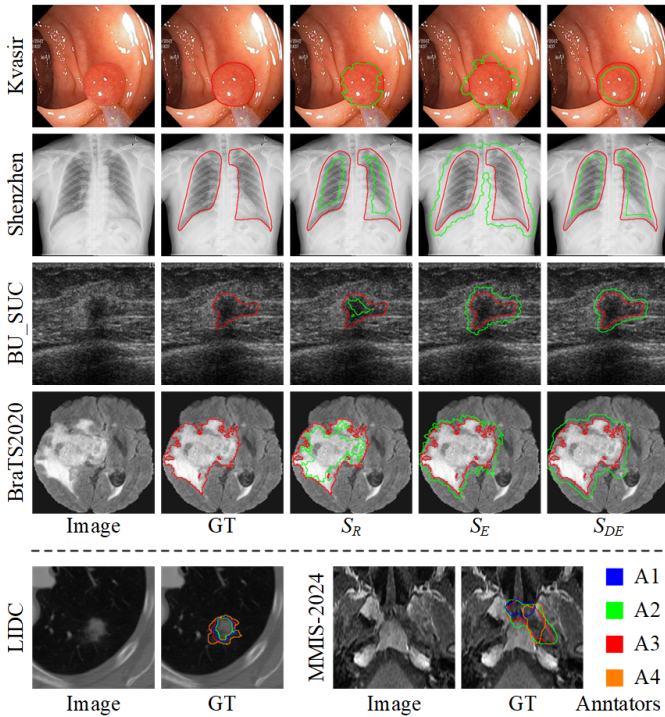
### 4.4. Comparison with State-of-the-Art Methods

To evaluate the effectiveness of our method, we conducted comparisons with a wide range of advanced approaches. These include standard loss functions such as Cross-Entropy (CE) and noise-robust alternatives, such as Generalized Cross Entropy

**Table 3. Performance comparison on Kvasir, Shenzhen and BU_SUC dataset under three simulated label-noise settings.**

| Method | Kvasir dataset | | | Shenzhen dataset | | | BU_SUC dataset | | |
|---|---|---|---|---|---|---|---|---|---|
| | $S_R$ | $S_E$ | $S_{DE}$ | $S_R$ | $S_E$ | $S_{DE}$ | $S_R$ | $S_E$ | $S_{DE}$ |
| CE Loss | 66.86±4.46 | 70.33±3.25 | 63.11±2.54 | 60.36±1.68 | 73.82±1.58 | 76.08±3.01 | 58.14±2.75 | 73.01±1.52 | 79.54±2.99 |
| GCE Loss | 66.40±2.63 | 68.04±1.67 | 60.67±2.59 | 60.86±1.01 | 74.83±0.53 | 76.04±3.40 | 57.91±2.23 | 72.99±0.83 | 78.78±1.07 |
| RCE Loss | 73.51±1.58 | 73.68±1.57 | 66.17±1.97 | 59.64±2.78 | 74.67±0.67 | 78.82±2.67 | 59.37±2.12 | 73.54±0.59 | 79.74±1.74 |
| RMD | 68.34±2.18 | 71.57±1.09 | 66.90±1.75 | 60.20±0.26 | 74.84±0.09 | 86.88±0.32 | 48.20±1.62 | 76.31±0.48 | 79.10±1.26 |
| ADELE | 60.97±14.78 | 67.10±11.42 | 60.62±13.04 | 60.26±3.06 | 72.74±2.29 | 81.62±8.71 | 57.12±9.01 | 71.86±7.65 | 80.63±10.25 |
| CDR | 67.87±3.51 | 70.58±1.65 | 63.51±1.67 | 64.84±5.39 | 75.69±1.79 | 79.79±2.29 | 59.95±2.63 | 74.53±0.92 | 80.71±2.61 |
| Co-Teaching | 74.26±1.71 | 75.57±1.12 | 74.03±1.48 | 63.10±0.81 | 75.99±0.29 | 91.63±0.15 | 57.09±0.94 | 76.92±0.46 | 85.54±0.91 |
| IDMPS | 77.52±1.21 | 74.16±0.81 | 69.47±0.81 | 61.93±1.19 | 74.69±0.51 | 78.81±1.29 | 60.22±1.96 | 73.92±0.62 | 82.83±0.72 |
| JoCoR | 67.98±3.23 | 71.43±1.37 | 65.62±1.94 | 67.46±0.37 | 73.44±0.07 | 89.56±0.04 | 56.42±0.77 | 73.54±0.28 | 81.31±0.16 |
| SP-Guide | 69.24±2.09 | 62.97±1.65 | 61.74±1.56 | 70.55±0.63 | 74.92±0.24 | 81.87±0.42 | 71.40±0.74 | 71.07±0.47 | 81.32±0.54 |
| Ours | **80.04±0.42** | **79.39±0.33** | **79.97±0.53** | **93.31±0.09** | **89.25±0.25** | **94.68±0.07** | **80.27±0.68** | **84.59±0.41** | **89.30±0.25** |



**Fig. 5. Visualization of simulated label noise ($S_R$: foreground-reducing, $S_E$: foreground-expanding, $S_{DE}$: simulated via dilation or erosion) and inter-expert variability. In the upper panel, red contours represent ground truth boundaries, and green contours indicate simulated noisy labels.**

(GCE) Zhang and Sabuncu (2018) and Reverse Cross Entropy (RCE) Wang et al. (2019). We also consider advanced frameworks such as IDMPS Zhao et al. (2024), ADELE Liu et al. (2022), and CDR Xia et al. (2020). The comparison also covered co-training and small-loss-based approaches such as Co-Teaching Han et al. (2018), JoCoR Wei et al. (2020), and RMD Fang et al. (2023). Additionally, we compared with SP-Guide Li et al. (2021), which leverages superpixels as structural priors for noisy label refinement. For all experiments, the Dice score was used as the evaluation metric.

### 4.4.1. Results on the Kvasir, Shenzhen, and BU_SUC datasets

Table 3 summarizes the quantitative results on three datasets under simulated noise conditions ($S_R$, $S_E$, and $S_{DE}$), with qualitative comparisons provided in Fig. 6. Across all settings, our

method consistently outperforms existing approaches, demonstrating robustness across different imaging modalities and adaptability to diverse noise patterns.

On the Kvasir dataset, the improvement is particularly pronounced under the $S_R$ noise setting, where compared methods struggle with foreground-reducing noise and exhibit marked Dice score degradation. Our method effectively mitigates this issue and achieves the highest performance, even under the widely used $S_{DE}$ setting, reaching a Dice score of 79.97%, which surpasses the best competing method at 74.03%.

On the Shenzhen dataset, which contains larger foreground regions, simulated noise corrupts a greater portion of each image, thereby increasing task difficulty. Our method achieves Dice scores of 93.31% and 89.25% under $S_R$ and $S_E$ noise, respectively, outperforming competing methods that suffer from severe under- and over-segmentation. Even in the commonly adopted $S_{DE}$ setting, where strong baselines such as Co-Teaching reach 91.63%, our method still achieves the best result with a Dice score of 94.68%.

On the BU_SUC dataset, characterized by high speckle noise and small, low-contrast lesions, our method substantially outperforms competing methods, achieving 80.27% Dice under the challenging $S_R$ setting compared with 71.40% by the best method. While performance improves across all methods under $S_{DE}$ noise, our framework offers superior segmentation by retaining delicate lesion details.

Overall, these results highlight that our framework not only achieves state-of-the-art performance across different datasets but also exhibits strong robustness against diverse and severe noise corruptions.

### 4.4.2. Results on the BraTS2020 dataset

We further evaluated our method on the 3D BraTS2020 brain tumor segmentation dataset. As reported in Table 4, it consistently outperformed existing approaches across all three simulated noise settings. Specifically, it achieved Dice scores of 81.53%, 82.68%, and 83.84%, respectively, which surpass the best-performing baselines by a notable margin. These results highlight the robustness of our framework in handling complex tumor structures and mitigating annotation noise in 3D brain MRI segmentation task.

**Fig. 6.** Qualitative comparison on the Kvasir (upper), Shenzhen (middle) and BU_SUC (lower) datasets under three simulated label-noise settings. Red contours denote ground-truth boundaries, while green contours represent predicted boundaries.

**Table 4.** Performance comparison of GSD-Net and existing methods on the BraTS2020 dataset under three simulated label-noise settings.

| Method | $S_R$ | $S_E$ | $S_{DE}$ |
|---|---|---|---|
| CE Loss | 69.16±1.02 | 75.61±0.30 | 67.90±0.65 |
| ADELE | 67.60±5.07 | 80.82±1.45 | 71.24±2.40 |
| CDR | 70.00±3.96 | 79.04±1.14 | 73.06±4.21 |
| Co-Teaching | 76.94±1.72 | 79.40±1.55 | 80.81±1.08 |
| IDMPS | 74.01±4.20 | 78.99±2.27 | 73.21±2.05 |
| JoCoR | 67.97±0.44 | 76.36±0.21 | 76.85±0.16 |
| Ours | **81.53±1.44** | **82.68±0.65** | **83.84±0.67** |

### 4.4.3. Results on the LIDC and MMIS-2024 dataset

Table 5 summarizes the quantitative performance of our proposed GSD-Net compared to several existing methods on the LIDC and MMIS-2024 datasets. (1) U-Net models trained on a single expert's annotations perform well when evaluated on that expert's annotations but show notable performance drops on annotations from other experts, revealing subjectivity and inconsistency in labeling. In contrast, GSD-Net exhibits moderate degradation under cross-expert evaluation, indicating stronger robustness to inter-expert variability. (2) To better reflect real-world multi-expert scenarios, we adopt a sampling strategy that selects one expert annotation per image during training. This introduces natural label variability and enables evaluation under realistic annotation noise. Under this setting, GSD-Net achieves an average Dice score of 88.25% on the LIDC and 79.24% on the MMIS-2024 dataset, outperforming all compared approaches under the same training conditions. (3) We further compare with D-Persona Wu et al. (2024), a representative multi-rater method that uses all four expert annotations per image to generate personalized predictions. When trained on 50% of images, which results in double the annotation cost,

**Table 5. Performance comparison of GSD-Net and existing methods on the LIDC and MMIS-2024 dataset. The A1, A2, A3, and A4 refer to four independent annotators, "Sample" refers to randomly assigning a fixed expert annotation to each image.**

| Method | Ann | Image Num | Ann Cost | LIDC dataset | | | | | MMIS-2024 dataset | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | $Dice_{A1}$ | $Dice_{A2}$ | $Dice_{A3}$ | $Dice_{A4}$ | $Dice_{mean}$ | $Dice_{A1}$ | $Dice_{A2}$ | $Dice_{A3}$ | $Dice_{A4}$ | $Dice_{mean}$ |
| CE Loss | A1 | 100% | 100% | **89.41** | 86.75 | 87.23 | 86.77 | **87.54** | **86.54** | 73.02 | 72.55 | 76.47 | **77.15** |
| | A2 | 100% | 100% | 86.41 | **89.19** | 86.27 | 85.61 | 86.87 | **78.09** | 76.55 | 73.50 | 74.44 | 75.65 |
| | A3 | 100% | 100% | 87.35 | 86.39 | **89.28** | 86.59 | 87.40 | 75.83 | 75.58 | **78.46** | 74.10 | 75.99 |
| | A4 | 100% | 100% | 86.85 | 85.85 | 86.52 | **88.74** | 86.99 | **78.33** | 72.51 | 71.95 | 76.06 | 74.71 |
| **Ours** | A1 | 100% | 100% | **89.47** | 87.69 | 88.37 | 88.03 | **88.39** | **87.36** | 74.08 | 73.11 | 77.04 | **77.90** |
| | A2 | 100% | 100% | 88.04 | **89.45** | 87.85 | 87.35 | 88.17 | 78.52 | **78.68** | 74.82 | 75.57 | 76.90 |
| | A3 | 100% | 100% | 88.19 | 87.53 | **89.47** | 87.69 | 88.22 | 76.91 | 76.98 | **80.04** | 75.34 | 77.32 |
| | A4 | 100% | 100% | 87.80 | 86.96 | 87.68 | **89.36** | 87.95 | **80.42** | 74.97 | 74.01 | 78.43 | 76.96 |
| D-person | All | 25% | 100% | 86.73 | 88.00 | 88.42 | 86.40 | 87.39 | 81.41 | 75.20 | 74.57 | 76.30 | 76.87 |
| | All | 50% | **200%** | 87.36 | **88.61** | **89.17** | 87.33 | 88.12 | 82.07 | 77.22 | 77.27 | 77.72 | 78.57 |
| CE Loss | Sample | 25% | 25% | 84.21 | 83.36 | 83.95 | 83.17 | 83.67 | 78.13 | 73.95 | 72.41 | 73.58 | 74.52 |
| | Sample | 50% | 50% | 86.01 | 85.24 | 86.09 | 85.31 | 85.66 | 80.03 | 75.65 | 74.09 | 75.97 | 76.44 |
| | Sample | 100% | 100% | 87.41 | 87.19 | 87.09 | 86.65 | 87.09 | 81.45 | 76.36 | 76.29 | 76.86 | 77.74 |
| **Ours** | Sample | 25% | 25% | 86.63 | 86.12 | 86.69 | 86.19 | 86.40 | 80.88 | 75.68 | 74.68 | 76.24 | 76.87 |
| | Sample | 50% | 50% | 87.72 | 87.08 | 87.71 | 87.62 | 87.53 | 82.10 | 77.10 | 76.33 | 78.06 | 78.40 |
| | Sample | 100% | 100% | **88.53** | 87.98 | 88.33 | **88.16** | **88.25** | 82.37 | 78.62 | 77.64 | 78.35 | 79.24 |
| CE Loss | Sample | 100% | 100% | 87.41 | 87.19 | 87.09 | 86.65 | 87.09 | 81.45 | 76.36 | 76.29 | 76.86 | 77.74 |
| JoCoR | Sample | 100% | 100% | 87.58 | 87.02 | 87.26 | 86.75 | 87.15 | 82.27 | 77.54 | 76.70 | 77.60 | 78.53 |
| ADELE | Sample | 100% | 100% | 87.63 | 87.12 | 87.27 | 86.68 | 87.17 | **83.10** | 77.67 | 76.53 | 77.49 | 78.70 |
| CDR | Sample | 100% | 100% | 87.35 | 87.14 | 87.00 | 86.45 | 86.99 | 81.70 | 76.13 | 75.27 | 76.68 | 77.45 |
| Co-Teaching | Sample | 100% | 100% | 87.65 | 87.25 | 87.33 | 87.08 | 87.33 | 82.41 | 76.93 | 77.25 | 77.68 | 78.57 |
| IDAMP | Sample | 100% | 100% | 88.23 | 87.82 | 87.68 | 87.60 | 87.83 | 81.91 | 76.91 | 76.32 | 77.12 | 78.57 |
| JoCoR | Sample | 100% | 100% | 87.58 | 87.02 | 87.26 | 86.75 | 87.15 | 82.27 | 77.54 | 76.70 | 77.60 | 78.53 |
| SP-Guide | Sample | 100% | 100% | 80.92 | 80.35 | 79.91 | 80.40 | 80.40 | 60.20 | 55.08 | 55.24 | 57.41 | 57.48 |
| **Ours** | Sample | 100% | 100% | **88.53** | **87.98** | **88.33** | **88.16** | **88.25** | 82.37 | **78.62** | **77.64** | **78.35** | **79.24** |

D-Persona achieves an average Dice score of 88.12% on the LIDC dataset and 78.57% on MMIS-2024 dataset, both lower than those achieved by GSD-Net.

In addition to the quantitative analysis, the qualitative segmentation results in Fig. 7 further highlight the strengths of GSD-Net, showing that it produces segmentations closer to the true semantic regions, with improved accuracy and more precise delineation of lesion boundaries. Taken together, these findings demonstrate that GSD-Net effectively addresses inter-expert variability while maintaining stable performance across datasets, underscoring its potential for real-world clinical deployment in scenarios where label noise arises from inter-annotator variability.

### 4.5. Ablation Study

To evaluate the contribution of each component, we conducted ablation studies on the Kvasir dataset. Using $\mathcal{L}_{JoCoR}$ (Eq. (6)) as the baseline, we analyzed the impact of the GDA, SGLR, and KT modules. Detailed results are reported in Table 6. The GDA module, which introduces geometric distance–aware reweighting, significantly improved segmentation accuracy across all three types, demonstrating its effectiveness to mitigate annotation noise. The SGLR module further enhances performance by incorporating superpixel-based structural priors to refine boundaries and correct labels, outperforming the LR variant. Building on these gains, the KT module delivers additional improvements by facilitating cross-sample knowledge transfer, which increased data diversity and reduced the influence of noisy labels. Together, these results demon-

**Table 6. Ablation study of component combinations on the Kvasir dataset. LR indicates the label refinement module without superpixel processing, and SGLR denotes the version with structure-guide integration.**

| Set. | $\mathcal{L}_{JoCoR}$ | GDA | LR | SGLR | KT | $S_R$ | $S_E$ | $S_{DE}$ |
|---|---|---|---|---|---|---|---|---|
| 1 | ✓ | | | | | 71.05±0.78 | 73.79±0.79 | 69.23±1.15 |
| 2 | ✓ | ✓ | | | | 76.76±0.64 | 75.91±0.37 | 75.13±0.48 |
| 3 | ✓ | ✓ | ✓ | | | 78.74±0.41 | 76.73±0.29 | 77.26±0.60 |
| 4 | ✓ | ✓ | ✓ | ✓ | | 78.84±0.46 | 78.94±0.62 | 78.56±0.53 |
| 5 | ✓ | ✓ | ✓ | | ✓ | 78.94±0.35 | 78.47±0.46 | 78.78±0.35 |
| 6 | ✓ | ✓ | ✓ | ✓ | ✓ | **80.04±0.42** | **79.39±0.33** | **79.97±0.53** |

strate that each module is effective and complementary in enhancing robustness. Grad-CAM visualizations (Fig. 8) further illustrate that as modules are integrated, the model progressively focuses on relevant features, especially lesion boundaries, providing visual evidence of improved interpretability.

## 5. Discussion

In medical image segmentation, label noise is often unavoidable due to intra- and inter-observer variability, indistinct lesion boundaries, and coarse delineations that arise when lesions are annotated in a rough or imprecise manner (Fig. 1) which can impair the network's ability to capture object features. To address these challenges, we propose GSD-Net, a framework designed to enhance robustness under noisy annotations. To validate its effectiveness, we first validated the framework on four public datasets with simulated noise, employing not only the widely used morphology-based erosion and dilation strategy ($S_{DE}$) but also Markov-based boundary perturbation strategies ($S_R$, $S_E$)
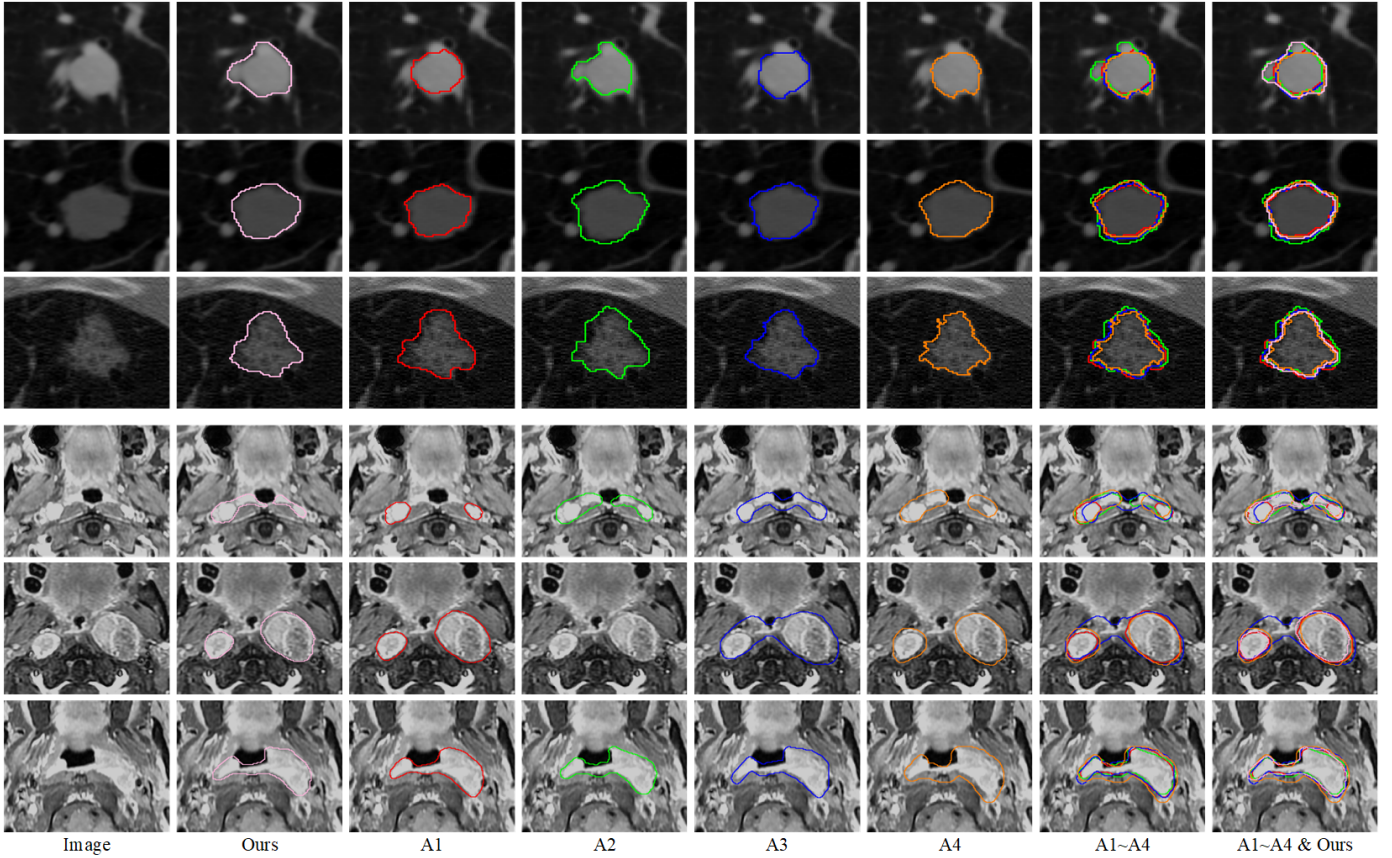
**Fig. 7. Qualitative comparisons on the LIDC dataset (upper) and the MMIS-2024 dataset (lower). A1~A4 denote annotation from different experts.**

to more accurately simulate the irregular boundaries observed in clinical practice (Fig. 5). Moreover, in large-scale medical datasets, collaborative annotations from multiple experts make inter-observer variability a frequent and realistic source of label noise. Results on the LIDC and MMIS-2024 datasets (Table 5 and Fig. 7) highlight this challenge: conventional models exhibit notable performance drops under cross-expert evaluation, whereas GSD-Net maintains stable performance. These findings demonstrate that our framework more effectively addresses inter-observer variability than existing methods, thereby improving its reliability in real-world multi-rater scenarios. These results underscore the clinical value of GSD-Net, as its robustness to inter-observer variability enables more consistent and reliable segmentation across different radiologists, thereby facilitating reproducible diagnosis and treatment planning.

Extensive experiments across six datasets spanning endoscopy, X-ray, US, MRI, and CT modalities show that GSD-Net achieves superior performance under diverse noise conditions, demonstrating strong generalizability and stability across different modalities and noise types. Furthermore, Grad-CAM visualizations (Fig. 8) provide qualitative evidence that GSD-Net progressively focuses on more relevant and anatomically meaningful structures, supporting its interpretability. Another notable strength is that GSD-Net serves as a training-level enhancement independent of specific network architectures, ensuring broad applicability and generalization. To further validate this property, we evaluated the framework on two alterna-

**Table 7. Performance of GSD-Net integrated with UNet++ and ViT U-Net.**

| Backbones | $S_R$ | $S_E$ | $S_{DE}$ |
|---|---|---|---|
| Kvasir Dataset | | | |
| UNet++ | 78.95±0.60 | 78.78±0.50 | 78.67±0.39 |
| ViT U-Net | 82.37±0.38 | 80.68±0.30 | 82.65±0.28 |
| U-Net | 80.04±0.42 | 79.39±0.33 | 79.97±0.53 |
| Shenzhen Dataset | | | |
| UNet++ | 93.84±0.06 | 90.58±0.16 | 94.86±0.57 |
| ViT U-Net | 93.86±0.08 | 90.14±0.17 | 94.45±0.10 |
| U-Net | 93.31±0.09 | 89.25±0.25 | 94.68±0.07 |
| BU_SUC Dataset | | | |
| UNet++ | 77.33±2.43 | 84.00±0.41 | 89.68±0.16 |
| ViT U-Net | 83.28±0.63 | 86.21±0.21 | 90.06±0.09 |
| U-Net | 80.27±0.68 | 84.59±0.41 | 89.30±0.25 |

tive backbones: UNet++ Zhou et al. (2020) and a U-Net variant with a Mix Vision Transformer encoder (ViT U-Net) Xie et al. (2021). As shown in Table 7, GSD-Net consistently improved segmentation performance for both backbones, underscoring its robustness and adaptability across diverse segmentation models.

In addition, GSD-Net exhibits limited sensitivity to the hyperparameter $\tau$, which is introduced in the small-loss strategy to filter noisy samples. As demonstrated in the ablation study on the Kvasir dataset (Fig. 9), variations in $\tau$ led to minor perfor-
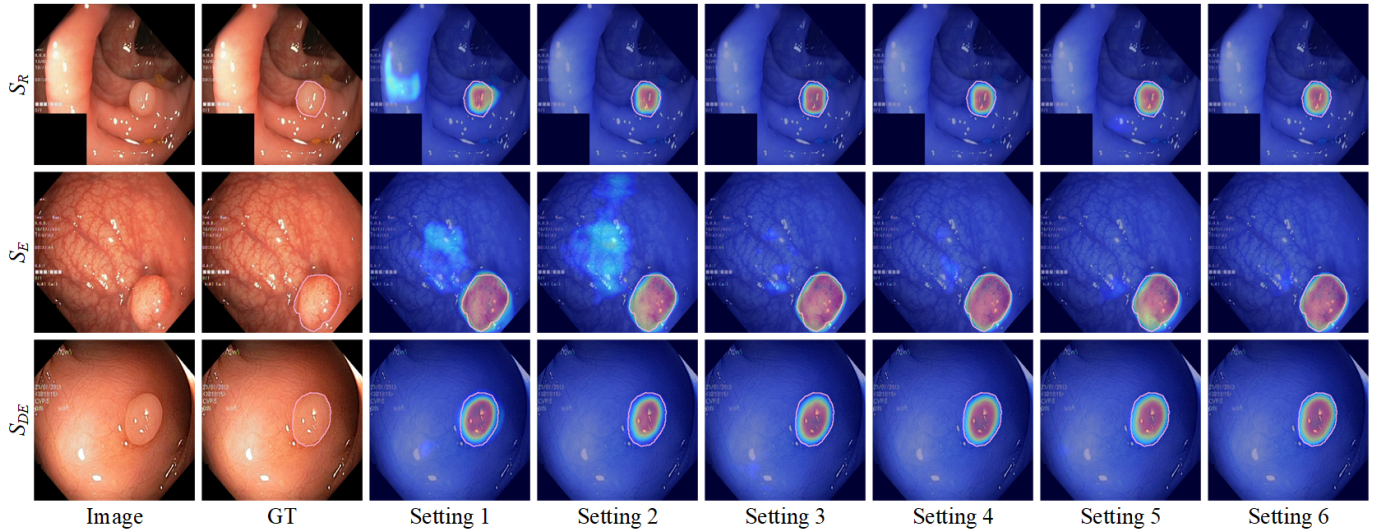
**Fig. 8. Grad-CAM visualizations under different ablation settings (Settings 1-6), consistent with the configurations shown in Table 6. The pink contours denote the boundary of ground truth.**
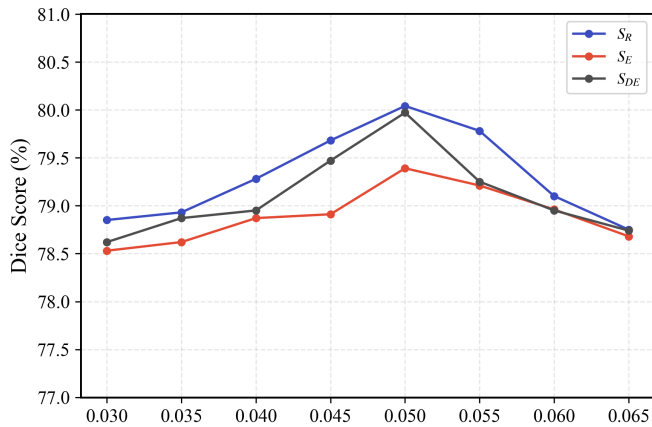


**Fig. 9. Effect of hyperparameters $\tau$ on segmentation performance on the Kvasir dataset.**

**Table 8. Performance across different proportions of clean labels trained with U-Net, and $S_{DE}$ noisy labels trained with GSD-Net.**

| Annotation | GT (Clean Labels) | | | | $S_{DE}$ (Noisy) |
|---|---|---|---|---|---|
| Ratio | 100% | 50% | 25% | 12.5% | 100% |
| Kvasir | **80.62±1.39** | 75.99±2.28 | 68.82±1.28 | 55.93±4.04 | 79.97±0.53 |
| Shenzhen | **95.24±0.30** | 94.93±0.39 | 94.57±0.28 | 94.20±0.58 | 94.68±0.07 |
| BU_SUC | **89.95±1.00** | 88.26±0.78 | 87.40±0.74 | 83.57±1.35 | 89.30±0.25 |

of striking a balance between annotation precision and dataset scale, thereby mitigating reliance on exhaustive expert labeling and leveraging coarse delineations to improve efficiency while maintaining reliable segmentation under noisy supervision.

Despite these promising results, several limitations remain. First, the robustness of the proposed framework requires further validation in real-world clinical settings involving multiple centers and imaging devices. Second, the current evaluation is restricted to binary segmentation tasks, as existing noise simulation strategies exhibit limited scalability to multi-class scenarios, where inter-class overlaps are more prevalent. Future work will focus on developing multi-class, multi-center datasets with authentic clinical annotations to more comprehensively evaluate and extend the robustness and applicability of the proposed framework.

mance changes, and the model maintained strong segmentation accuracy even under suboptimal values. These results highlight the framework's robustness to hyperparameter selection and its reliability under imperfect noise-level estimation.

In practical applications, GSD-Net shows strong potential for reducing annotation costs. Annotations produced by junior doctors or medical students with limited training may contain noise. However, when processed through the proposed framework, they can still yield competitive performance. As shown in Table 8, under the mild noise setting ($S_{DE}$) the model achieves results comparable to those obtained with fully clean labels. In contrast, reducing the proportion of clean annotations to 12.5% results in a substantial performance drop on some datasets. These results suggest that a small amount of expert annotation alone is insufficient, whereas larger volumes of coarse annotations from less experienced annotators, combined with noise-robust learning, offer greater potential for competitive performance. Overall, the results underscore the importance

## 6. Conclusion

In this study, we proposed GSD-Net, a noise-robust framework for medical image segmentation. Building on the small-loss strategy and integrating geometric distance–aware weighting, structure-guided label refinement, and knowledge transfer mechanisms, the framework effectively suppresses label noise and improves segmentation reliability. Extensive experiments on six public datasets showed consistent gains over state-of-the-art methods under both simulated and real-world noise, underscoring the robustness and practical potential of GSD-Net for accurate segmentation under imperfect supervision.

## 7. Acknowledgments

## References

Achanta, R., Shaji, A., Smith, K., Lucchi, A., Fua, P., Süsstrunk, S., 2012. SLIC Superpixels Compared to State-of-the-Art Superpixel Methods. IEEE transactions on pattern analysis and machine intelligence 34, 2274–2282.

Armato III, S.G., McLennan, G., Bidaut, L., McNitt-Gray, M.F., Meyer, C.R., Reeves, A.P., Zhao, B., Aberle, D.R., Henschke, C.I., Hoffman, E.A., et al., 2011. The Lung Image Database Consortium (LIDC) and Image Database Resource Initiative (IDRI): A Completed Reference Database of Lung Nodules on CT Scans. Medical Physics 38, 915–931.

Arpit, D., Jastrzębski, S., Ballas, N., Krueger, D., Bengio, E., Kanwal, M.S., Maharaj, T., Fischer, A., Courville, A., Bengio, Y., et al., 2017. A Closer Look at Memorization in Deep Networks, in: International Conference on Machine Learning, PMLR. pp. 233–242.

Bakas, S., Sako, C., Akbari, H., Bilello, M., Sotiras, A., Shukla, G., Rudie, J., Flores Santamaria, N., Fathi Kazerooni, A., Pati, S., et al., 2021. Multi-parametric magnetic resonance imaging (mpMRI) scans for de novo Glioblastoma (GBM) patients from the University of Pennsylvania Health System (UPENN-GBM). The Cancer Imaging Archive .

Barron, J.T., 2019. A General and Adaptive Robust Loss Function, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4331–4339.

Candemir, S., Jaeger, S., Palaniappan, K., Musco, J.P., Singh, R.K., Xue, Z., Karargyris, A., Antani, S., Thoma, G., McDonald, C.J., 2013. Lung Segmentation in Chest Radiographs Using Anatomical Atlases With Nonrigid Registration. IEEE Transactions on Medical Imaging 33, 577–590.

Cepeda, S., García-García, S., Arrese, I., Herrero, F., Escudero, T., Zamora, T., Sarabia, R., 2023. The Río Hortega University Hospital Glioblastoma dataset: A comprehensive collection of preoperative, early postoperative and recurrence MRI scans (RHUH-GBM). Data in Brief 50, 109617.

Çiçek, Ö., Abdulkadir, A., Lienkamp, S.S., Brox, T., Ronneberger, O., 2016. 3D U-Net: Learning Dense Volumetric Segmentation from Sparse Annotation, in: Ourselin, S., Joskowicz, L., Sabuncu, M.R., Unal, G., Wells, W. (Eds.), Medical Image Computing and Computer-Assisted Intervention – MICCAI 2016, Springer International Publishing, Cham. pp. 424–432.

Fang, C., Wang, Q., Cheng, L., Gao, Z., Pan, C., Cao, Z., Zheng, Z., Zhang, D., 2023. Reliable Mutual Distillation for Medical Image Segmentation Under Imperfect Annotations. IEEE Transactions on Medical Imaging 42, 1720–1734.

Gonzalez-Jimenez, A., Lionetti, S., Gottfrois, P., Gröger, F., Navarini, A., Pouly, M., 2025. Robust T-loss for medical image segmentation. Medical Image Analysis , 103735.

Han, B., Yao, Q., Yu, X., Niu, G., Xu, M., Hu, W., Tsang, I., Sugiyama, M., 2018. Co-teaching: Robust Training of Deep Neural Networks with Extremely Noisy Labels. Advances in Neural Information Processing Systems 31.

Iqbal, A., Sharif, M., 2024. Memory-efficient transformer network with feature fusion for breast tumor segmentation and classification task. Engineering Applications of Artificial Intelligence 127, 107292.

Jaeger, S., Karargyris, A., Candemir, S., Folio, L., Siegelman, J., Callaghan, F., Xue, Z., Palaniappan, K., Singh, R.K., Antani, S., et al., 2013. Automatic Tuberculosis Screening Using Chest Radiographs. IEEE Transactions on Medical Imaging 33, 233–245.

Jha, D., Smedsrud, P.H., Riegler, M.A., Halvorsen, P., De Lange, T., Johansen, D., Johansen, H.D., 2019. Kvasir-SEG: A Segmented Polyp Dataset, in: International Conference on Multimedia Modeling, Springer. pp. 451–462.

Jin, S., Lu, W., Monkam, P., 2022. Deep Neural Network-Based Noisy Pixel Estimation for Breast Ultrasound Segmentation, in: 2022 IEEE International Conference on Image Processing (ICIP), IEEE. pp. 1776–1780.

Ko, Y., Moon, S., Baek, J., Shim, H., 2021. Rigid and non-rigid motion artifact reduction in X-ray CT using attention module. Medical Image Analysis 67, 101883.

Kohl, S., Romera-Paredes, B., Meyer, C., De Fauw, J., Ledsam, J.R., Maier-Hein, K., Eslami, S., Jimenez Rezende, D., Ronneberger, O., 2018. A Probabilistic U-Net for Segmentation of Ambiguous Images. Advances in Neural Information Processing Systems 31.

Kuang, H., Wang, Y., Tan, X., Yang, J., Sun, J., Liu, J., Qiu, W., Zhang, J., Zhang, J., Yang, C., et al., 2025. Lw-ctrans: A lightweight hybrid network of cnn and transformer for 3d medical image segmentation. Medical Image Analysis 102, 103545.

Lee, H.J., Kim, J.U., Lee, S., Kim, H.G., Ro, Y.M., 2020. Structure Boundary Preserving Segmentation for Medical Image With Ambiguous Boundary, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

Li, S., Gao, Z., He, X., 2021. Superpixel-Guided Iterative Learning from Noisy Labels for Medical Image Segmentation, in: Medical Image Computing and Computer Assisted Intervention – MICCAI 2021, Springer. pp. 525–535.

Li, X.C., Xia, X., Zhu, F., Liu, T., Zhang, X.Y., Liu, C.L., 2023. Dynamics-aware loss for learning with label noise. Pattern Recognition 144, 109835.

Liu, S., Liu, K., Zhu, W., Shen, Y., Fernandez-Granda, C., 2022. Adaptive Early-Learning Correction for Segmentation From Noisy Annotations, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2606–2616.

Liu, T., Bai, Q., Torigian, D.A., Tong, Y., Udupa, J.K., 2024. Vsmtrans: A hybrid paradigm integrating self-attention and convolution for 3d medical image segmentation. Medical Image Analysis 98, 103295.

Luo, X., Liao, W., He, Y., Tang, F., Wu, M., Shen, Y., Huang, H., Song, T., Li, K., Zhang, S., et al., 2023. Deep learning-based accurate delineation of primary gross tumor volume of nasopharyngeal carcinoma on heterogeneous magnetic resonance imaging: A large-scale and multi-center study. Radiotherapy and Oncology 180, 109480.

Ma, X., Zhang, Z., Ji, Z., Huang, K., Su, N., Yuan, S., Chen, Q., 2023. Adjustable Robust Transformer for High Myopia Screening in Optical Coherence Tomography, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 504–514.

Mehrnia, M., Kholmovski, E., Katsaggelos, A., Kim, D., Passman, R., Elbaz, M.S., 2024. Novel Self-Calibrated Threshold-Free Probabilistic Fibrosis Signature Technique for 3D Late Gadolinium Enhancement MRI. IEEE transactions on Biomedical Engineering .

Menze, B.H., Jakab, A., Bauer, S., Kalpathy-Cramer, J., Farahani, K., Kirby, J., Burren, Y., Porz, N., Slotboom, J., Wiest, R., et al., 2014. The Multimodal Brain Tumor Image Segmentation Benchmark (BRATS). IEEE Transactions on Medical Imaging 34, 1993–2024.

Northcutt, C., Jiang, L., Chuang, I., 2021. Confident Learning: Estimating Uncertainty in Dataset Labels. Journal of Artificial Intelligence Research 70, 1373–1411.

Paquin, A.L., Chaib-draa, B., Giguère, P., . Symmetrization of loss functions for robust training of neural networks in the presence of noisy labels .

Qiu, L., Zhao, L., Hou, R., Zhao, W., Zhang, S., Lin, Z., Teng, H., Zhao, J., 2023. Hierarchical multimodal fusion framework based on noisy label learning and attention mechanism for cancer classification with pathology and genomic features. Computerized Medical Imaging and Graphics 104, 102176.

Ronneberger, O., Fischer, P., Brox, T., 2015. U-Net: Convolutional Networks for Biomedical Image Segmentation, in: Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015, Springer. pp. 234–241.

Schmidt, A., Morales-Alvarez, P., Molina, R., 2023. Probabilistic Modeling of Inter- and Intra-observer Variability in Medical Image Segmentation, in: Proceedings of the IEEE/CVF international conference on computer vision, pp. 21097–21106.

Shen, W., Peng, Z., Wang, X., Wang, H., Cen, J., Jiang, D., Xie, L., Yang,

X., Tian, Q., 2023. A Survey on Label-Efficient Deep Image Segmentation: Bridging the Gap Between Weak Supervision and Dense Prediction. IEEE Transactions on Pattern Analysis and Machine Intelligence 45, 9284–9305.

Shi, J., Wu, J., 2021. Distilling Effective Supervision for Robust Medical Image Segmentation with Noisy Labels, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 668–677.

Shi, J., Zhang, K., Guo, C., Yang, Y., Xu, Y., Wu, J., 2024. A survey of label-noise deep learning for medical image analysis. Medical Image Analysis 95, 103166.

Sindhwani, V., Niyogi, P., Belkin, M., 2005. A Co-Regularization Approach to Semi-supervised Learning with Multiple Views, in: Proceedings of ICML Workshop on Learning with Multiple Views, Citeseer. pp. 74–79.

Singla, R., Ringstrom, C., Hu, R., Lessoway, V., Reid, J., Rohling, R., Nguan, C., 2022. Speckle and Shadows: Ultrasound-specific Physics-based Data Augmentation for Kidney Segmentation, in: International Conference on Medical Imaging with Deep Learning, PMLR. pp. 1139–1148.

Stirenko, S., Kochura, Y., Alienin, O., Rokovyi, O., Gordienko, Y., Gang, P., Zeng, W., 2018. Chest X-Ray Analysis of Tuberculosis by Deep Learning with Segmentation and Augmentation, in: 2018 IEEE 38th International Conference on Electronics and Nanotechnology (ELNANO), IEEE. pp. 422–428.

Suter, Y., Knecht, U., Valenzuela, W., Notter, M., Hewer, E., Schucht, P., Wiest, R., Reyes, M., 2022. The LUMIERE dataset: Longitudinal Glioblastoma MRI with expert RANO evaluation. Scientific Data 9, 768.

Wang, Y., Ma, X., Chen, Z., Luo, Y., Yi, J., Bailey, J., 2019. Symmetric Cross Entropy for Robust Learning With Noisy Labels, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 322–330.

Wei, H., Feng, L., Chen, X., An, B., 2020. Combating Noisy Labels by Agreement: A Joint Training Method with Co-Regularization, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 13726–13735.

Wu, Y., Luo, X., Xu, Z., Guo, X., Ju, L., Ge, Z., Liao, W., Cai, J., 2024. Diversified and Personalized Multi-rater Medical Image Segmentation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 11470–11479.

Xia, X., Liu, T., Han, B., Gong, C., Wang, N., Ge, Z., Chang, Y., 2020. Robust early-learning: Hindering the memorization of noisy labels, in: International Conference on Learning Representations.

Xiao, R., Dong, Y., Wang, H., Feng, L., Wu, R., Chen, G., Zhao, J., 2022. Promix: Combating label noise via maximizing clean sample utility. arXiv preprint arXiv:2207.10276 .

Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J.M., Luo, P., 2021. Seg-Former: Simple and Efficient Design for Semantic Segmentation with Transformers. Advances in Neural Information Processing Systems 34, 12077–12090.

Xu, Z., Lu, D., Luo, J., Wang, Y., Yan, J., Ma, K., Zheng, Y., Tong, R.K.Y., 2022. Anti-Interference From Noisy Labels: Mean-Teacher-Assisted Confident Learning for Medical Image Segmentation. IEEE Transactions on Medical Imaging 41, 3062–3073.

Yang, S., Yang, E., Han, B., Liu, Y., Xu, M., Niu, G., Liu, T., 2022. Estimating Instance-dependent Bayes-label Transition Matrix using a Deep Neural Network, in: International Conference on Machine Learning, PMLR. pp. 25302–25312.

Yao, J., Zhang, Y., Zheng, S., Goswami, M., Prasanna, P., Chen, C., 2023. Learning to Segment from Noisy Annotations: A Spatial Correction Approach. arXiv preprint arXiv:2308.02498 .

Yi, R., Huang, Y., Guan, Q., Pu, M., Zhang, R., 2021. Learning From Pixel-Level Label Noise: A New Perspective for Semi-Supervised Semantic Segmentation. IEEE Transactions on Image Processing 31, 623–635.

Yu, S., Chen, M., Zhang, E., Wu, J., Yu, H., Yang, Z., Ma, L., Gu, X., Lu, W., 2020. Robustness study of noisy annotation in deep learning based medical image segmentation. Physics in Medicine & Biology 65, 175007.

Zhang, D., Han, J., Cheng, G., Yang, M.H., 2021. Weakly Supervised Object Localization and Detection: A Survey. IEEE Transactions on Pattern Analysis and Machine Intelligence 44, 5866–5885.

Zhang, T., Yu, L., Hu, N., Lv, S., Gu, S., 2020. Robust Medical Image Segmentation from Non-expert Annotations with Tri-network, in: Medical Image Computing and Computer Assisted Intervention–MICCAI 2020, Springer. pp. 249–258.

Zhang, Z., Sabuncu, M., 2018. Generalized Cross Entropy Loss for Training Deep Neural Networks with Noisy Labels. Advances in Neural Information Processing Systems 31.

Zhao, X., Li, Z., Luo, X., Li, P., Huang, P., Zhu, J., Liu, Y., Zhu, J., Yang, M., Chang, S., et al., 2024. Ultrasound Nodule Segmentation Using Asymmetric Learning With Simple Clinical Annotation. IEEE Transactions on Circuits and Systems for Video Technology .

Zhou, Z., Siddiquee, M.M.R., Tajbakhsh, N., Liang, J., 2020. UNet++: Redesigning Skip Connections to Exploit Multiscale Features in Image Segmentation. IEEE Transactions on Medical Imaging 39, 1856–1867.