

Rectifying Noisy Labels with Sequential Prior: Multi-Scale Temporal Feature Affinity Learning for Robust Video Segmentation

Beilei Cui^{*,1}[0009-0009-7900-8032], Mingqing Zhang^{*,2}[0000-0002-7214-0569],
Mengya Xu^{*,3}[0000-0002-4338-7079], An Wang¹[0000-0001-5515-0653], Wu
Yuan^{†,2}[0000-0001-9405-519X], and Hongliang Ren^{†,1,3}[0000-0002-6488-1551]

- ¹ Dept. of Electronic Engineering, The Chinese University of Hong Kong, Hong Kong SAR, China
² Dept. of Biomedical Engineering, The Chinese University of Hong Kong, Hong Kong SAR, China
³ Dept. of Biomedical Engineering, National University of Singapore, Singapore
 beileicui@link.cuhk.edu.hk, lars.zhang@link.cuhk.edu.hk,
 mengya@u.nus.edu, wa09@link.cuhk.edu.hk, wyuan@cuhk.edu.hk,
 ren@nus.edu.sg

Abstract. Noisy label problems are inevitably in existence within medical image segmentation causing severe performance degradation. Previous segmentation methods for noisy label problems only utilize a single image while the potential of leveraging the correlation between images has been overlooked. Especially for video segmentation, adjacent frames contain rich contextual information beneficial in cognizing noisy labels. Based on two insights, we propose a Multi-Scale Temporal Feature Affinity Learning (MS-TFAL) framework to resolve noisy-labeled medical video segmentation issues. First, we argue the sequential prior of videos is an effective reference, i.e., pixel-level features from adjacent frames are close in distance for the same class and far in distance otherwise. Therefore, Temporal Feature Affinity Learning (TFAL) is devised to indicate possible noisy labels by evaluating the affinity between pixels in two adjacent frames. We also notice that the noise distribution exhibits considerable variations across video, image, and pixel levels. In this way, we introduce Multi-Scale Supervision (MSS) to supervise the network from three different perspectives by re-weighting and refining the samples. This design enables the network to concentrate on clean samples in a coarse-to-fine manner. Experiments with both synthetic and real-world label noise demonstrate that our method outperforms recent state-of-the-art robust segmentation approaches. Code is available at <https://github.com/BeileiCui/MS-TFAL>.

Keywords: Noisy label learning · Feature affinity · Semantic segmentation.

* Authors contributed equally to this work.

† Corresponding Author.

1 Introduction

Video segmentation, which refers to assigning pixel-wise annotation to each frame in a video, is one of the most vital tasks in medical image analysis. Thanks to the advance in deep learning algorithms based on Convolutional Neural Networks, medical video segmentation has achieved great progress over recent years [9]. But a major problem of deep learning methods is that they are largely dependent on both the quantity and quality of training data [13]. Datasets annotated by non-expert humans or automated systems with little supervision typically suffer from very high label noise and are extremely time-consuming. Even expert annotators could generate different labels based on their cognitive bias [6]. Based on the above, noisy labels are inevitably in existence within medical video datasets causing misguidance to the network and resulting in severe performance degradation. Hence, it is of great importance to design medical video segmentation methods that are robust to noisy labels within training data [4, 18].

Most of the previous noisy label methods mainly focus on classification tasks. Only in recent years, the problem of noise labels in segmentation tasks has been more explored, but still less involved in medical image analysis. Previous techniques for solving noisy label problems in medical segmentation tasks can be categorized in three directions. The first type of method aims at deriving and modeling the general distribution of noisy labels in the form of Noise Transition Matrix (NTM) [3, 8]. Secondly, some researchers develop special training strategies to re-weight or re-sample the data such that the model could focus on more dependable samples. Zhang et al. [19] concurrently train three networks and each network is trained with pixels filtered by the other two networks. Shi et al. [14] use stable characteristics of clean labels to estimate samples' uncertainty map which is used to further guide the network. Thirdly, label refinement is implemented to renovate noisy labels. Li et al. [7] represent the image with superpixels to exploit more advanced information in an image and refine the labels accordingly. Liu et al. [10] use two different networks to jointly determine the error sample, and use each other to refine the labels to prevent error accumulation. Xu et al. [15] utilize the mean-teacher model and Confident learning to refine the low-quality annotated samples.

Despite the amazing performance in tackling noisy label issues for medical image segmentation, almost all existing techniques only make use of the information within a single image. *To this end, we make the effort in exploring the feature affinity relation between pixels from consecutive frames.* The motivation is that the embedding features of pixels from adjacent frames should be close if they belong to the same class, and should be far if they belong to different classes. Hence, if a pixel's feature is far from the pixels of the same class in the adjacent frame and close to the ones of different classes, its label is more likely to be incorrect. Meanwhile, the distribution of noisy labels may vary among different videos and frames, which also motivates us to supervise the network from multiple perspectives.

Inspired by the motivation above and to better resolve noisy label problems with temporal consistency, we propose Multi-Scale Temporal Feature Affinity

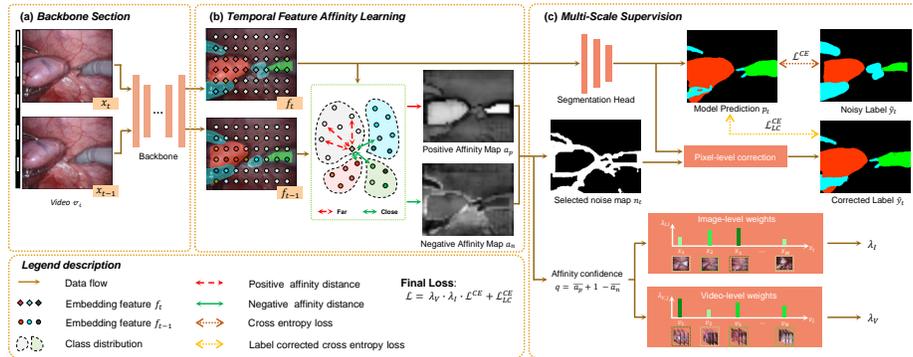


Fig. 1. Illustration of proposed Multi-Scale Temporal Feature Affinity Learning framework. We acquire the embedding feature maps of adjacent frames in the Backbone Section. Then, the temporal affinity is calculated for each pixel in current frame to obtain the positive and negative affinity map indicating possible noisy labels. The affinity maps are then utilized to supervise the network in a multi-scale manner.

Learning (MS-TFAL) framework. Our contributions can be summarized as the following points:

1. In this work, we first propose a novel Temporal Feature Affinity Learning (TFAL) method to evaluate the temporal feature affinity map of an image by calculating the similarity between the same and different classes' features of adjacent frames, therefore indicating possible noisy labels.
2. We further develop a Multi-Scale Supervision (MSS) framework based on TFAL by supervising the network through video, image, and pixel levels. Such a coarse-to-fine learning process enables the network to focus more on correct samples at each stage and rectify the noisy labels, thus improving the generalization ability.
3. Our method is validated on a publicly available dataset with synthetic noisy labels and a real-world label noise dataset and obtained superior performance over other state-of-the-art noisy label techniques.
4. To the best of our knowledge, we are the first to tackle noisy label problems using inter-frame information and discover the superior ability of sequential prior information to resolve noisy label issues.

2 Method

The proposed Multi-Scale Temporal Feature Affinity Learning Framework is illustrated in Fig. 1. We aim to exploit the information from adjacent frames to identify the possible noisy labels, thereby learning a segmentation network robust to label noises by re-weighting and refining the samples. Formally, given an input training image $x_t \in \mathbb{R}^{H \times W \times 3}$, and its adjacent frame x_{t-1} , two feature maps $f_t, f_{t-1} \in \mathbb{R}^{h \times w \times C_f}$ are first generated by a CNN backbone, where h, w

and C_f represent the height, width and channel number. Intuitively, for each pair of features from f_t and f_{t-1} , their distance should be close if they belong to the same class and far otherwise. Therefore for each pixel in f_t , we calculate two affinity relations with f_{t-1} . The first one is called positive affinity, computed by averaging the cosine similarity between one pixel $f_t(i)$ in the current frame and all the same class' pixels as $f_t(i)$ in previous frame. The second one is called negative affinity, computed by averaging the cosine similarity between one pixel $f_t(i)$ in current frame and all the different class' pixels as $f_t(i)$ in previous frame. Then through up-sampling, the Positive Affinity Map a_p and Negative Affinity Map a_n can be obtained, where $a_p, a_n \in \mathbb{R}^{H \times W}$, denote the affinity relation between x_t and x_{t-1} . The positive affinity of clean labels should be high while the negative affinity of clean labels should be low. Therefore, the black areas in a_p and the white areas in a_n are more likely to be noisy labels.

Then we use two affinity maps a_p, a_n to conduct Multi-Scale Supervision training. Multi-scale refers to video, image, and pixel levels. Specifically, for pixel-level supervision, we first obtain thresholds t_p and t_n by calculating the average positive and negative affinity over the entire dataset. The thresholds are used to determine the possible noisy label sets based on positive and negative affinity separately. The intersection of two sets is selected as the final noisy set and relabeled with the model prediction p_t . The affinity maps are also used to estimate the image-level weights λ_I and video-level weights λ_V . The weights enable the network to concentrate on videos and images with higher affinity confidence. Our method is a plug-in module that is not dependent on backbone type and can be applied to both image-based backbones and video-based backbones by modifying the shape of inputs and feature maps.

2.1 Temporal Feature Affinity Learning

The purpose of this section is to estimate the affinity between pixels in the current frame and previous frame, thus indicating possible noisy labels. Specifically, in addition to the aforementioned feature map $f_t, f_{t-1} \in \mathbb{R}^{h \times w \times C_f}$, we obtain the down-sampled labels with the same size of feature map $\tilde{y}'_t, \tilde{y}'_{t-1} \in \mathbb{R}^{h \times w \times \mathcal{C}}$, where \mathcal{C} means the total class number. We derive the positive and negative label maps with binary variables: $M_p, M_n \subseteq \{0, 1\}^{hw \times hw}$. The value corresponds to pixel (i, j) is determined by the label as:

$$M_p(i, j) = \mathbb{1} \left[\tilde{y}'_t(i) = \tilde{y}'_{t-1}(j) \right], \quad M_n(i, j) = \mathbb{1} \left[\tilde{y}'_t(i) \neq \tilde{y}'_{t-1}(j) \right] \quad (1)$$

where $\mathbb{1}(\cdot)$ is the indicator function. $M_p(i, j) = 1$ when i th label in \tilde{y}'_t and j th label in \tilde{y}'_{t-1} are the same class, while $M_p(i, j) = 0$ otherwise; and M_n vice versa. The value of cosine similarity map $S \in \mathbb{R}^{hw \times hw}$ corresponds to pixel (i, j) is determined by: $S(i, j) = \frac{f_t(i) \cdot f_{t-1}(j)}{\|f_t(i)\| \times \|f_{t-1}(j)\|}$. We then use the average cosine similarity of a pixel with all pixels in the previous frame belonging to the same or different class to represent its positive or negative affinity:

$$a_{p,f}(i) = \frac{\sum_{j=1}^{hw} S(i,j) M_p(i,j)}{\sum_{j=1}^{hw} M_p(i,j)}, \quad a_{n,f}(i) = \frac{\sum_{j=1}^{hw} S(i,j) M_n(i,j)}{\sum_{j=1}^{hw} M_n(i,j)} \quad (2)$$

where $a_{p,f}, a_{n,f} \in \mathbb{R}^{h \times w}$ means the positive and negative map with the same size as the feature map. With simple up-sampling, we could obtain the final affinity maps $a_p, a_n \in \mathbb{R}^{H \times W}$, indicating the positive and negative affinity of pixels in the current frame.

2.2 Multi-Scale Supervision

The feature map is first connected with a segmentation head generating the prediction p . Besides the standard cross entropy loss $\mathcal{L}^{CE}(p, \tilde{y}) = -\sum_i^{HW} \tilde{y}(i) \log p(i)$, we applied a label corrected cross entropy loss $\mathcal{L}_{LC}^{CE}(p, \hat{y}) = -\sum_i^{HW} \hat{y}(i) \log p(i)$ to train the network with pixel-level corrected labels. We further use two weight factors λ_I and λ_V to supervise the network in image and video levels. The specific descriptions are explained in the following sections.

Pixel-Level Supervision. Inspired by the principle in Confident Learning [12], we use affinity maps to denote the confidence of labels. if a pixel $x(i)$ in an image has both small enough positive affinity $a_p(i) \leq t_p$ and large enough negative affinity $a_n(i) \geq t_n$, then its label $\tilde{y}(i)$ can be suspected as noisy. The threshold t_p, t_n are obtained empirically by calculating the average positive and negative affinity, formulated as $t_p = \frac{1}{|A_p|} \sum_{a_p \in A_p} \bar{a}_p$, $t_n = \frac{1}{|A_n|} \sum_{a_n \in A_n} \bar{a}_n$, where \bar{a}_p, \bar{a}_n means the average value of positive and negative affinity over an image. The noisy pixels set can therefore be defined by:

$$\tilde{x} := \{x(i) \in x : a_p(i) \leq t_p\} \cap \{x(i) \in x : a_n(i) \geq t_n\}. \quad (3)$$

Then we update the pixel-level label map \hat{y} as:

$$\hat{y}(i) = \mathbb{1}(x(i) \in \tilde{x}) p(i) + \mathbb{1}(x(i) \notin \tilde{x}) \tilde{y}(i), \quad (4)$$

where $p(i)$ is the prediction of network. Through this process, we only replace those pixels with both low positive affinity and large negative affinity.

Image-Level Supervision. Even in the same video, different frames may contain different amounts of noisy labels. Hence, we first define the affinity confidence value as: $q = \bar{a}_p + 1 - \bar{a}_n$. The average affinity confidence value is therefore denoted as: $\bar{q} = t_p + 1 - t_n$. Finally, we define the image-level weight as:

$$\lambda_I = e^{2(q-\bar{q})}. \quad (5)$$

$\lambda_I > 1$ if the sample has large affinity confidence and $\lambda_I < 1$ otherwise, therefore enabling the network to concentrate more on the clean samples.

Video-Level Supervision. We assign different weights to different videos such that the network can learn from more correct videos in the early stage. We first define the video affinity confidence as the average affinity confidence of

all the frames: $q_v = \frac{1}{|V|} \sum_{x \in V} q_x$. Supposing there are N videos in total, we use $k \in \{1, 2, \dots, N\}$ to represent the ranking of video affinity confidence from small to large, which means $k = 1$ and $k = N$ denote the video with lowest and highest affinity confidence separately. Video-level weight is thus formulated as:

$$\lambda_V = \begin{cases} \theta_l, & \text{if } k < \frac{N}{3} \\ \theta_l + \frac{3k-N}{N}(\theta_u - \theta_l), & \text{if } \frac{N}{3} \leq k \leq \frac{2N}{3} \\ \theta_u, & \text{if } k > \frac{2N}{3} \end{cases} \quad (6)$$

where θ_l and θ_u are the presetted lower-bound and upper-bound of weight.

Combining the above-defined losses and weights, we obtain the final loss as: $\mathcal{L} = \lambda_V \lambda_I \mathcal{L}^{CE} + \mathcal{L}_{LC}^{CE}$, which supervise the network in a multi-scale manner. These losses and weights are enrolled in training after initialization in an order of video, image, and pixel enabling the network to enhance the robustness and generalization ability by concentrating on clean samples from rough to subtle.

3 Experiments

3.1 Dataset Description and Experiment Settings

EndoVis 2018 Dataset and Noise Patterns. EndoVis 2018 Dataset is from the MICCAI robotic instrument segmentation dataset of endoscopic vision challenge 2018 [1]. It is officially divided into 15 videos with 2235 frames for training and 4 videos with 997 frames for testing separately. The dataset contains 12 classes including different anatomy and robotic instruments. Each image is resized into 256×320 in pre-processing. To better simulate manual noisy annotations within a video, we first randomly select a ratio of α of videos and in each selected video, we divide all frames into several groups in a group of $3 \sim 6$ consecutive frames. Then for each group of frames, we randomly apply dilation, erosion, affine transformation, or polygon noise to each class [7, 16, 18, 19]. We investigated our algorithms in several noisy settings with α being $\{0.3, 0.5, 0.8\}$. Some examples of data and noisy labels are shown in supplementary.

Rat Colon Dataset. For real-world noisy dataset, we have collected rat colon OCT images using 800nm ultra-high resolution endoscopic spectral domain OCT. We refer readers to [17] for more details. Each centimeter of rat colon imaged corresponds to 500 images with 6 class layers of interest. We select 8 sections with 2525 images for training and 3 sections with 1352 images for testing. The labels of test set were annotated by professional endoscopists as ground truth while the training set was annotated by non-experts. Each image is resized into 256×256 in pre-processing. Some dataset examples are shown in supplementary.

Implementation Details. We adopt Deeplabv3+ [2] as our backbone network for fair comparison. The framework is implemented with PyTorch on two Nvidia 3090 GPUs. We adopt the Adam optimizer with an initial learning rate of $1e - 4$ and weight decay of $1e - 4$. Batch size is set to 4 with a maximum of 100

<https://endovissub2018-roboticscenese segmentation.grand-challenge.org/>

Table 1. Comparison of other methods and our models on EndoVis 2018 Dataset under different ratios of noise. The best results are **highlighted**.

Data	Method	$mIOU(\%)$	Sequence $mIOU(\%)$				$Dice(\%)$
			Seq 1	Seq 2	Seq 3	Seq 4	
Clean	Deeplabv3+ [2]	53.98	54.07	51.46	72.35	38.02	64.30
Noisy, $\alpha = 0.3$	Deeplabv3+ [2]	50.42	49.90	48.50	67.18	36.10	60.60
	STswin (22') [5]	50.29	49.96	48.67	66.52	35.99	60.62
	RAUNet (19') [11]	50.36	44.97	48.06	68.90	39.53	60.61
	JCAS (22') [4]	48.65	48.77	46.60	64.83	34.39	58.97
	VolMin (21') [8]	47.64	45.60	45.31	64.01	35.63	57.42
	MS-TFAL(Ours)	52.91	49.48	51.60	71.08	39.52	62.91
Noisy, $\alpha = 0.5$	Deeplabv3+ [2]	42.87	41.72	42.96	59.54	27.27	53.02
	STswin (22') [5]	44.48	40.78	45.22	60.50	31.45	54.99
	RAUNet (19') [11]	46.74	46.16	43.08	63.00	34.73	57.44
	JCAS (22') [4]	45.24	41.90	44.06	61.13	33.90	55.22
	VolMin (21') [8]	44.02	42.68	46.26	59.67	27.47	53.59
	MS-TFAL(Ours)	50.34	49.15	50.17	67.37	34.67	60.50
Noisy, $\alpha = 0.8$	Deeplabv3+ [2]	33.35	27.57	35.69	45.30	24.86	42.22
	STswin (22') [5]	32.27	28.92	34.48	42.97	22.72	42.61
	RAUNet (19') [11]	33.25	30.23	34.95	44.99	22.88	43.67
	JCAS (22') [4]	35.99	28.29	38.06	51.00	26.66	44.75
	VolMin (21') [8]	33.85	28.40	39.38	43.76	23.90	42.63
	MS-TFAL(Ours)	41.36	36.33	41.65	59.57	27.88	51.01
Noisy, $\alpha = 0.5$	w/ V	47.80	44.73	48.71	66.87	30.91	57.45
	w/ V & I	48.72	43.20	48.44	66.34	36.93	58.54
	Same frame	48.99	46.77	49.58	64.70	34.94	59.31
	Any frame	48.69	46.25	48.92	65.56	34.08	58.89
	MS-TFAL(Ours)	50.34	49.15	50.17	67.37	34.67	60.50

epochs for both Datasets. θ_l and θ_u are set to 0.4 and 1 separately. The video, image, and pixel level supervision are involved from the 16th, 24th, and 40th epoch respectively. The segmentation performance is assessed by $mIOU$ and $Dice$ scores.

3.2 Experiment Results on EndoVis 2018 Dataset

Table 1 presents the comparison results under different ratios of label noises. We evaluate the performance of backbone trained with clean labels, two state-of-the-art instrument segmentation network [5, 11], two noisy label learning techniques [4, 8], backbone [2] and the proposed MS-TFAL. We re-implement [4, 8] with the same backbone [2] for a fair comparison. Compared with all other methods, MS-TFAL shows the minimum performance gap with the upper bound (Clean) for both $mIOU$ and $Dice$ scores under all ratios of noises demonstrating the robustness of our method. As noise increases, the performance of all baselines decreases significantly indicating the huge negative effect of noisy labels. It is noteworthy that when the noise ratio rises from 0.3 to 0.5 and from 0.5 to 0.8, our method only drops 2.57% $mIOU$ with 2.41% $Dice$ and 8.98% $mIOU$ with 9.49% $Dice$, both are the minimal performance degradation, which further demonstrates the robustness of our method against label noise. In the extreme

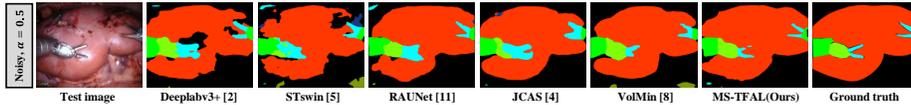


Fig. 2. Comparison of qualitative segmentation results on EndoVis18 Dataset.

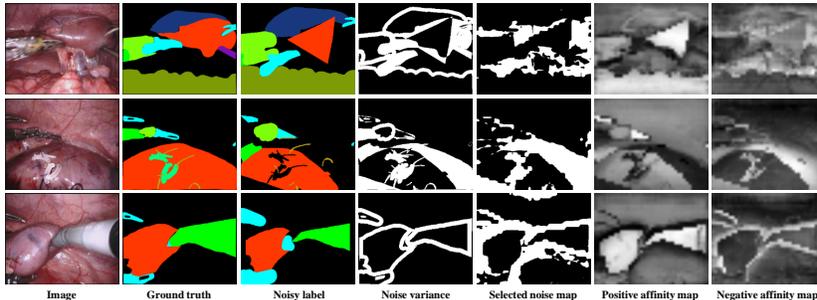


Fig. 3. Illustration of Noise variance and feature affinity. Selected noisy label (Fifth column) means the noise map selected with Equation (3).

noise setting ($\alpha = 0.8$), our method achieves 41.36% *mIOU* and 51.01% *Dice* and outperforms second best method 5.37% *mIOU* and 6.26% *Dice*. As shown in Fig. 2, we provide partial qualitative results indicating the superiority of MS-TFAL over other methods in the qualitative aspect. More qualitative results are shown in supplementary.

Ablation Studies. We further conduct two ablation studies on our multi-scale components and choice of frame for feature affinity under noisy dataset with $\alpha = 0.5$. With only video-level supervision (w/ V), *mIOU* and *Dice* are increased by 4.93% and 4.43% compared with backbone only. Then we apply both video and image level supervision (w/ V & I) and gain an increase of 0.92% *mIOU* and 1.09% *Dice*. Pixel-level supervision is added at last forming the complete Multi-Scale Supervision results in another improvement of 1.62% *mIOU* and 1.96% *Dice* verifying the effectiveness in attenuating noisy label issues of individual components. For the ablation study of the choice of frame, we compared two different attempts with ours: conduct TFAL with the same frame and any frame in the dataset (Ours is adjacent frame). Results show that using adjacent frame has the best performance compared to the other two choices.

Visualization of Temporal Affinity. To prove the effectiveness of using affinity relation we defined to represent the confidence of label, we display comparisons between noise variance and selected noise map in Fig. 3. Noise variance (Fourth column) represents the incorrect label map and the Selected noise map (Fifth column) denotes the noise map we select with Equation (3). We can observe that the noisy labels we affirm have a high overlap degree with the true noise labels, which demonstrates the validity of our TFAL module.

Table 2. Comparison of other methods and our models on Rat Colon Dataset.

Method	Deeplabv3+ [2]	STswin [5]	RAUNet [11]	JCAS [4]	VolMin [8]	MS-TFAL(Ours)
<i>mIOU</i> (%)	68.46	68.21	68.24	68.15	68.81	71.05
<i>Dice</i> (%)	75.25	77.70	77.39	77.50	77.89	80.17

3.3 Experiment Results on Rat Colon Dataset

The comparison results on real-world noisy Rat Colon Dataset are presented in Table 2. Our method outperforms other methods consistently on both *mIOU* and *Dice* scores, which verifies the superior robustness of our method on real-world label noise issues. Qualitative results are shown in supplementary.

4 Discussion and Conclusion

In this paper, we propose a robust MS-TFAL framework to resolve noisy label issues in medical video segmentation. Different from previous methods, we first introduce the novel TFAL module to use affinity between pixels from adjacent frames to represent the confidence of label. We further design MSS framework to supervise the network from multiple perspectives. Our method can not only identify noise in labels, but also correct them in pixel-wise with rich temporal consistency. Extensive experiments under both synthetic and real-world label noise data demonstrate the excellent noise resilience of MS-TFAL.

Acknowledgements

This work was supported by Hong Kong Research Grants Council (RGC) Collaborative Research Fund (C4026-21G), General Research Fund (GRF 14211420 & 14203323), Shenzhen-Hong Kong-Macau Technology Research Programme (Type C) STIC Grant SGDX20210823103535014 (202108233000303).

References

1. Allan, M., Kondo, S., Bodenstedt, S., Leger, S., Kadkhodamohammadi, R., Luengo, I., Fuentes, F., Flouty, E., Mohammed, A., Pedersen, M., et al.: 2018 robotic scene segmentation challenge. arXiv preprint arXiv:2001.11190 (2020)
2. Chen, L.C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H.: Encoder-decoder with atrous separable convolution for semantic image segmentation. In: Proceedings of the European conference on computer vision (ECCV). pp. 801–818 (2018)
3. Guo, X., Yang, C., Li, B., Yuan, Y.: Metacorrection: Domain-aware meta loss correction for unsupervised domain adaptation in semantic segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3927–3936 (2021)
4. Guo, X., Yuan, Y.: Joint class-affinity loss correction for robust medical image segmentation with noisy labels. In: Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part IV. pp. 588–598. Springer (2022)

5. Jin, Y., Yu, Y., Chen, C., Zhao, Z., Heng, P.A., Stoyanov, D.: Exploring intra- and inter-video relation for surgical semantic scene segmentation. *IEEE Transactions on Medical Imaging* **41**(11), 2991–3002 (2022)
6. Karimi, D., Dou, H., Warfield, S.K., Gholipour, A.: Deep learning with noisy labels: Exploring techniques and remedies in medical image analysis. *Medical image analysis* **65**, 101759 (2020)
7. Li, S., Gao, Z., He, X.: Superpixel-guided iterative learning from noisy labels for medical image segmentation. In: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part I* 24. pp. 525–535. Springer (2021)
8. Li, X., Liu, T., Han, B., Niu, G., Sugiyama, M.: Provably end-to-end label-noise learning without anchor points. In: Meila, M., Zhang, T. (eds.) *Proceedings of the 38th International Conference on Machine Learning. Proceedings of Machine Learning Research*, vol. 139, pp. 6403–6413. PMLR (18–24 Jul 2021), <https://proceedings.mlr.press/v139/li211.html>
9. Litjens, G., Kooi, T., Bejnordi, B.E., Setio, A.A.A., Ciompi, F., Ghafoorian, M., van der Laak, J.A., van Ginneken, B., Sánchez, C.I.: A survey on deep learning in medical image analysis. *Medical Image Analysis* **42**, 60–88 (2017). <https://doi.org/https://doi.org/10.1016/j.media.2017.07.005>, <https://www.sciencedirect.com/science/article/pii/S1361841517301135>
10. Liu, L., Zhang, Z., Li, S., Ma, K., Zheng, Y.: S-cuda: self-cleansing unsupervised domain adaptation for medical image segmentation. *Medical Image Analysis* **74**, 102214 (2021)
11. Ni, Z.L., Bian, G.B., Zhou, X.H., Hou, Z.G., Xie, X.L., Wang, C., Zhou, Y.J., Li, R.Q., Li, Z.: Raunet: Residual attention u-net for semantic segmentation of cataract surgical instruments. In: *International Conference on Neural Information Processing*. pp. 139–149. Springer (2019)
12. Northcutt, C., Jiang, L., Chuang, I.: Confident learning: Estimating uncertainty in dataset labels. *Journal of Artificial Intelligence Research* **70**, 1373–1411 (2021)
13. Shen, D., Wu, G., Suk, H.I.: Deep learning in medical image analysis. *Annual Review of Biomedical Engineering* **19**(1), 221–248 (2017). <https://doi.org/10.1146/annurev-bioeng-071516-044442>, <https://doi.org/10.1146/annurev-bioeng-071516-044442>, PMID: 28301734
14. Shi, J., Wu, J.: Distilling effective supervision for robust medical image segmentation with noisy labels. In: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part I* 24. pp. 668–677. Springer (2021)
15. Xu, Z., Lu, D., Luo, J., Wang, Y., Yan, J., Ma, K., Zheng, Y., Tong, R.K.Y.: Anti-interference from noisy labels: Mean-teacher-assisted confident learning for medical image segmentation. *IEEE Transactions on Medical Imaging* **41**(11), 3062–3073 (2022)
16. Xue, C., Deng, Q., Li, X., Dou, Q., Heng, P.A.: Cascaded robust learning at imperfect labels for chest x-ray segmentation. In: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part VI* 23. pp. 579–588. Springer (2020)
17. Yuan, W., Feng, Y., Chen, D., Gharibani, P., Chen, J.D.Z., Yu, H., Li, X.: In vivo assessment of inflammatory bowel disease in rats with ultrahigh-resolution colonoscopic oct. *Biomed. Opt. Express* **13**(4), 2091–2102 (Apr 2022). <https://doi.org/10.1364/boe.447121>

[//doi.org/10.1364/BOE.453396](https://doi.org/10.1364/BOE.453396), <https://opg.optica.org/boe/abstract.cfm?URI=boe-13-4-2091>

18. Zhang, M., Gao, J., Lyu, Z., Zhao, W., Wang, Q., Ding, W., Wang, S., Li, Z., Cui, S.: Characterizing label errors: confident learning for noisy-labeled image segmentation. In: Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part I 23. pp. 721–730. Springer (2020)
19. Zhang, T., Yu, L., Hu, N., Lv, S., Gu, S.: Robust medical image segmentation from non-expert annotations with tri-network. In: Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part IV 23. pp. 249–258. Springer (2020)

5 Supplementary

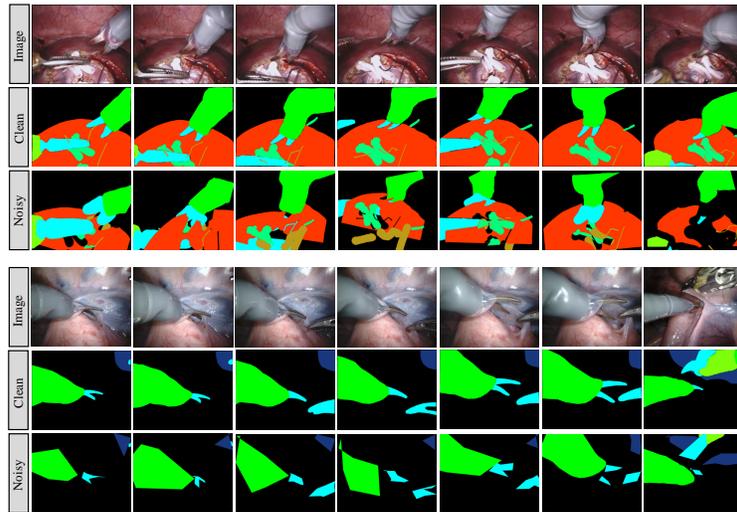


Fig. 4. Examples of EndoVis 2018 Dataset [1] including images , clean and noisy labels.

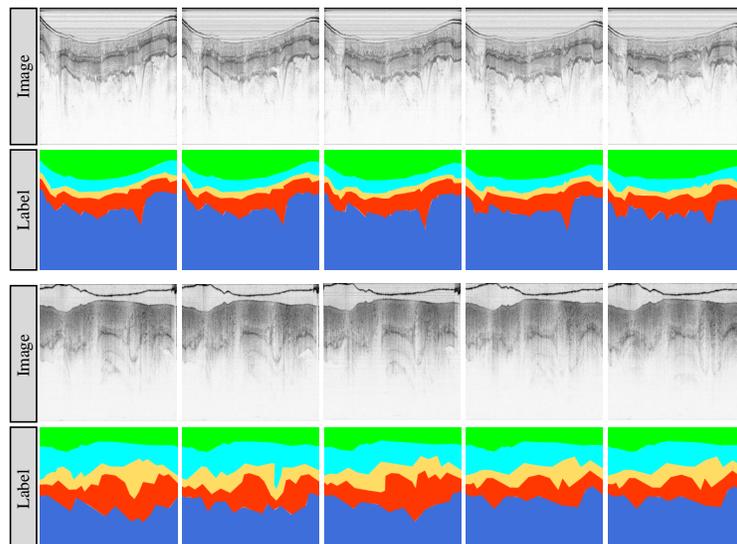


Fig. 5. Examples of Rat Colon Dataset.

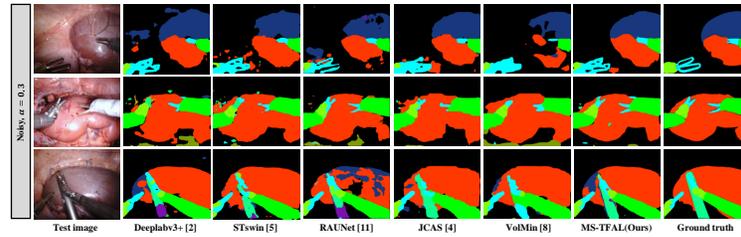


Fig. 6. Qualitative segmentation results on EndoVis 2018 Dataset ($\alpha = 0.3$).

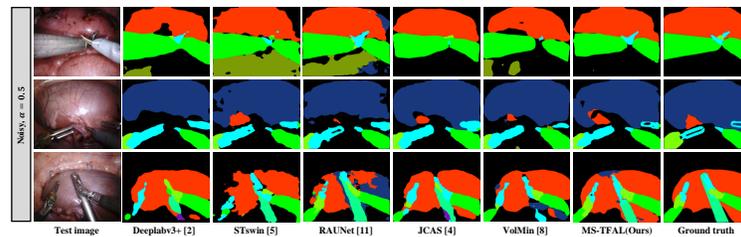


Fig. 7. Qualitative segmentation results on EndoVis 2018 Dataset ($\alpha = 0.5$).

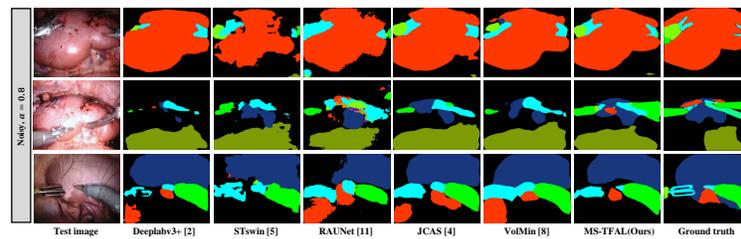


Fig. 8. Qualitative segmentation results on EndoVis 2018 Dataset ($\alpha = 0.8$).

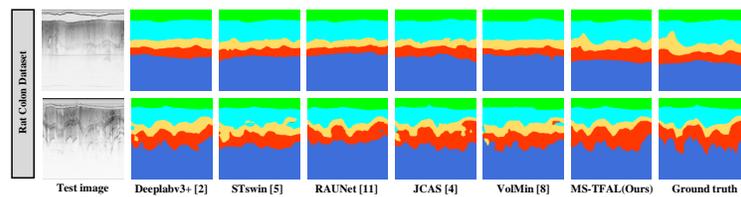


Fig. 9. Qualitative segmentation results on Rat Colon Dataset.