


**LETTER**

# Revisiting instrument segmentation: Learning from decentralized surgical sequences with various imperfect annotations

Zhou Zheng<sup>1</sup>  | Yuichiro Hayashi<sup>1</sup> | Masahiro Oda<sup>1,2</sup> | Takayuki Kitasaka<sup>3</sup> | Kensaku Mori<sup>1,2,4</sup>

<sup>1</sup>Graduate School of Informatics, Nagoya University, Chikusa-ku, Nagoya, Aichi, Japan

<sup>2</sup>Information Strategy Office, Information and Communications, Nagoya University, Chikusa-ku, Nagoya, Aichi, Japan

<sup>3</sup>School of Information Science, Aichi Institute of Technology, Yagusa-cho, Toyota, Aichi, Japan

<sup>4</sup>Research Center for Medical Bigdata, National Institute of Informatics, Chiyoda-ku, Tokyo, Japan

**Correspondence**

Kensaku Mori, Graduate School of Informatics, Nagoya University, Chikusa-ku, Nagoya, Aichi, Japan.

Email: [kensaku@is.nagoya-u.ac.jp](mailto:kensaku@is.nagoya-u.ac.jp)

**Funding information**

Japan Society for the Promotion of Science (JSPS), KAKENHI, Grant/Award Numbers: 21K19898, 17H00867; Japan Science and Technology Agency (JST), Core Research for Evolutional Science and Technology (CREST), Grant/Award Number: JPMJCR20D5; Japan Science and Technology Agency (JST), Moonshot Research and Development (R&D) Program, Grant/Award Number: JPMJMS2214-07

**Abstract**

This paper focuses on a new and challenging problem related to instrument segmentation. This paper aims to learn a generalizable model from distributed datasets with various imperfect annotations. Collecting a large-scale dataset for centralized learning is usually impeded due to data silos and privacy issues. Besides, local clients, such as hospitals or medical institutes, may hold datasets with diverse and imperfect annotations. These datasets can include scarce annotations (many samples are unlabelled), noisy labels prone to errors, and scribble annotations with less precision. Federated learning (FL) has emerged as an attractive paradigm for developing global models with these locally distributed datasets. However, its potential in instrument segmentation has yet to be fully investigated. Moreover, the problem of learning from various imperfect annotations in an FL setup is rarely studied, even though it presents a more practical and beneficial scenario. This work rethinks instrument segmentation in such a setting and propose a practical FL framework for this issue. Notably, this approach surpassed centralized learning under various imperfect annotation settings. This method established a foundational benchmark, and future work can build upon it by considering each client owning various annotations and aligning closer with real-world complexities.

## 1 | INTRODUCTION

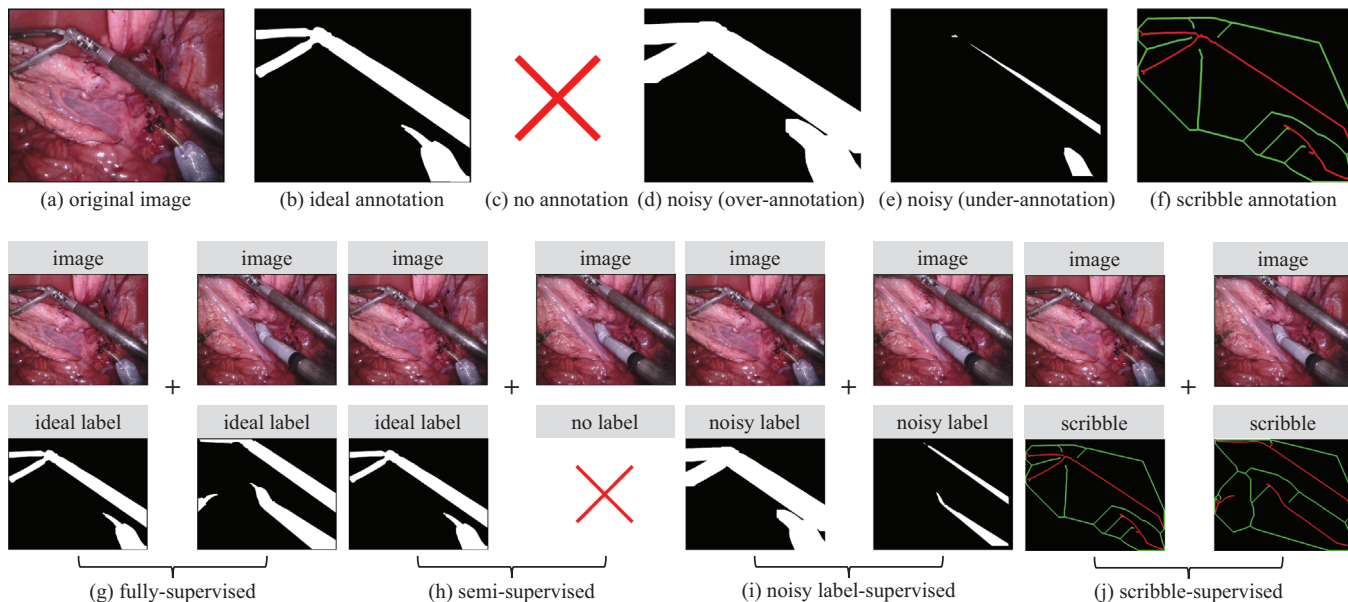
Robust and accurate surgical instrument segmentation serves as the cornerstone for potential applications such as instrument tracking and augmented reality within the domain of robotic minimally invasive surgery. Deep learning-based methodologies, especially the deep convolutional neural networks, have been state-of-the-art solutions for this task, driving a series of significant advancements in this field [1, 2]. Conventionally, establishing a high-performing, generalizable model for instrument segmentation is contingent on centralized learning [3–5]. This process necessitates the collection of large-scale surgical videos/sequences for training. However, practical

implementation is often limited due to stringent privacy and confidentiality-related regulations that restrict the sharing and collecting sensitive surgical data [6]. In light of these limitations of centralized learning, federated learning (FL) is emerging as an appealing alternative, enabling multiple clients or institutes to collaboratively prepare a global model without sharing local datasets. However, the application of FL in the context of instrument segmentation lacks exploration.

Moreover, the problem of learning from various imperfect annotations [7] in an FL setup is hardly investigated, even though it presents a more practical and advantageous scenario. For instance, due to the costly and time-consuming annotation process, a large portion of images might remain unlabelled at

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivs](https://creativecommons.org/licenses/by-nc-nd/4.0/) License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

© 2024 The Authors. *Healthcare Technology Letters* published by John Wiley & Sons Ltd on behalf of The Institution of Engineering and Technology.



**FIGURE 1** Illustration of various annotations and comparison of segmentation paradigms with various annotations. Note that endoscopic images are from EndoVis17 [1]. (a) Original image. (b) Ideal annotation (ground truth). (c) Scarce annotations, where a large portion of images have no annotation. (d,e) Noisy annotations, where we simulate two types of noisy labels, that is, over-annotation of foreground using dilation operation and under-annotation of foreground using erosion operation. (f) Scribble annotation, where only a small part of the foreground (in red) and background (in green) pixels are annotated, leaving a large portion (in black) unlabelled. Note that we further dilate the skeleton lines for better visualization. (g) Fully-supervised segmentation utilizes ideal annotations for learning. (h) Semi-supervised segmentation leverages a combination of a small proportion of labelled data and a large quantity of unlabelled data. (i) Noisy label-supervised segmentation aims to minimize the adverse impact of label noise during model training. (j) Scribble-supervised segmentation capitalizes on sparse supervision signals.

some clients [8, 9], resulting in scarce annotations. In addition, noisy annotations are sometimes unavoidable due to inter-observer variability, as the annotation process is inherently subjective [10, 11]. Some clients might use weak annotations, for example, scribble annotations, to reduce their reliance on densely annotated labels [12, 13]. Examples of the ideal, scarce, noisy<sup>1</sup>, and scribble<sup>2</sup> annotations are illustrated in Figure 1a through Figure 1f. As local clients could face various imperfect annotations, this complexity implies that local clients need to maximally utilize these imperfect annotations to contribute an effective local model for global model preparation.

To this end, we start a first attempt to study a more challenging problem related to instrument segmentation. We aim to learn instrument segmentation from distributed surgical sequences with various imperfect annotations. Instrument segmentation in real-world scenarios inherently introduces complexities, particularly when each client possesses multiple annotation types. These diverse annotations carry unique challenges and demand tailored strategies, elevating the intricacy of the problem. In our study, we proceed with the assumption that each client has just one type of annotation. Such an assumption allows us to focus on a specific and rudimentary condition of this problem. By tackling the challenge from this perspective, we can explore a potential solution and provide an initial baseline. Subsequent research can expand on this foundation by gradually integrating additional complexities,

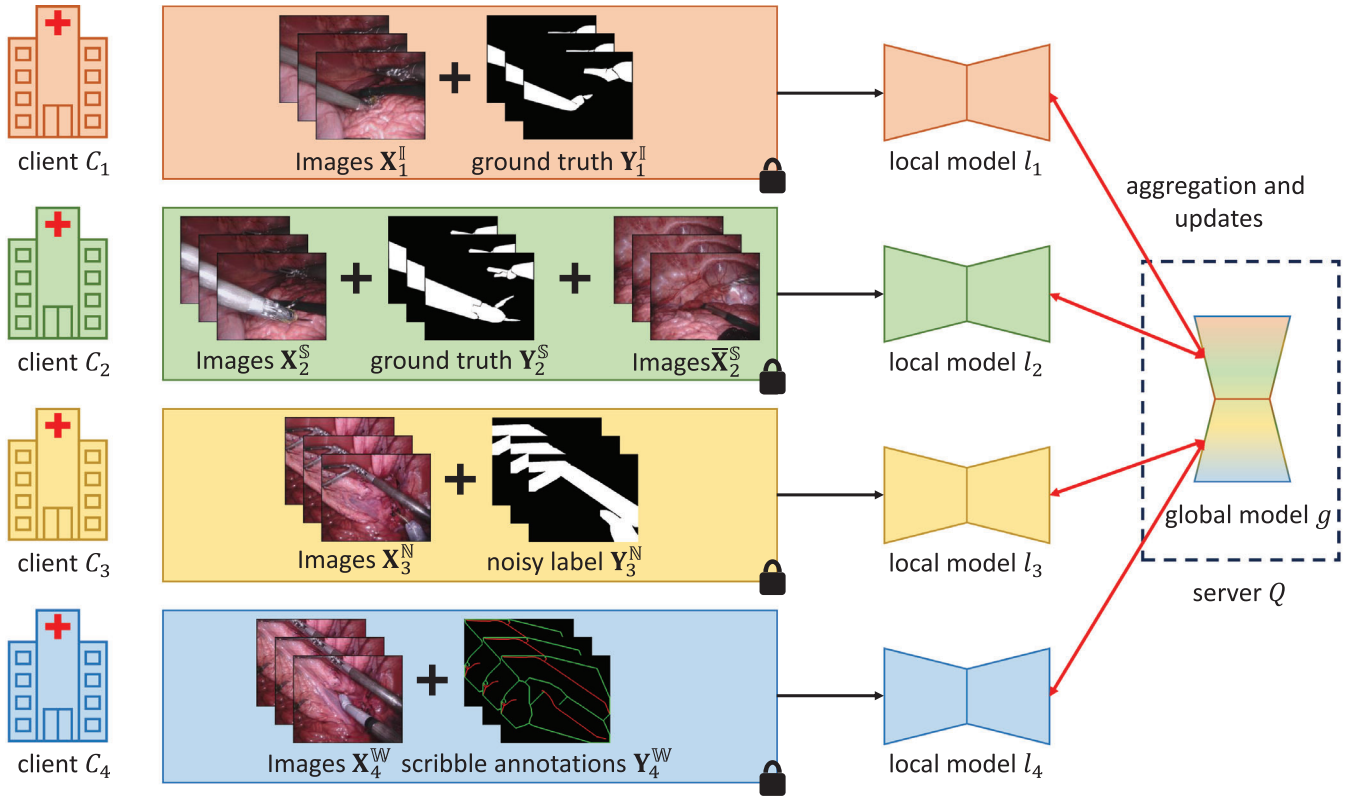
such as accommodating multiple annotation types per client, to strive for a more comprehensive approach that better reflects real-world scenarios.

From the above perspectives, we propose a practical FL framework to handle this problem and demonstrate a potential solution. Specifically, we adopt existing mainstream methodologies, for example, semi-supervised learning [15], noisy-label learning [11], and scribble-supervised learning [16], and show that these specific methods can be integrated to our framework to endow local clients with the ability to deal with imperfect annotations. Furthermore, we maintain a central server to coordinate these local clients and prepare a generalizable model using these distributed datasets with imperfect annotations. Our contributions are as follows:

- We delve into the problem of learning a generalizable model for instrument segmentation using distributed datasets with imperfect annotations. This problem, which adheres more closely to the complexities encountered in practice, remains unexplored and presents greater challenges in instrument segmentation.
- We present a novel and effective framework to address this problem by unifying semi-, noisy label-, and scribble-supervised segmentation and FL. We show that the careful integration of existing advanced techniques can provide successful solutions for instrument segmentation, taking into account both data availability and imperfect annotations.
- We conduct experiments with the EndoVis17 dataset [1] to present our method's efficacy. Our approach surpasses

<sup>1</sup> We simulate noisy annotations with dilation and erosion operations, following the work [10].

<sup>2</sup> We simulate scribble annotations by skeletonizing [14] the ground truth.



**FIGURE 2** Flowchart of proposed framework. This illustration depicts our solution to the novel and practical problem regarding data silos and imperfect annotations within instrument segmentation from surgical sequences, with the assumption that each local client only possesses a dataset with one type of annotation. For clarity, we feature four local clients  $\{C_i \mid i = 1, 2, 3, 4\}$  alongside a central server  $Q$ . Local clients  $\{C_i \mid i = 1, 2, 3, 4\}$  signify distinct data repositories, and their datasets reflect diverse annotations, that is, ideal annotations (ground truth), scarce annotations (a large portion of images are unlabelled), noisy annotations, and scribble annotations. Local clients adopt specific strategies, that is, fully-, semi-, noisy label-, and scribble-supervised learning, to prepare their respective local models  $\{l_i \mid i = 1, 2, 3, 4\}$  by making use of their distinctive datasets  $\mathbb{I}_1 = \{X_1^I, Y_1^I\}$ ,  $\mathbb{S}_2 = \{X_2^S, Y_2^S, \bar{X}_2^S\}$ ,  $\mathbb{N}_3 = \{X_3^N, Y_3^N\}$ , and  $\mathbb{W}_4 = \{X_4^W, Y_4^W\}$ . The central server  $Q$  manages model aggregation and updates with the FedAvg scheme [22], guiding local clients toward a generalizable global model  $g$ .

centralized learning under various imperfect annotation settings, underscoring its potential to tackle this innovative challenge. Our method can serve as an initial baseline for addressing this problem.

## 2 | RELATED WORK

In this paper, different from existing instrument segmentation methods, we unify semi-, noisy label-, and scribble-supervised segmentation and FL to address instrument segmentation while considering data privacy and imperfect annotation-related issues, leading to a novel and more challenging application. In the following, we review related literature on instrument segmentation, learning segmentation from imperfect annotations, and FL.

### 2.1 | Instrument segmentation

Instrument segmentation in robotic surgery [1, 17] plays an important role in enhancing surgical procedures. Over the past few years, it has attracted increasing interest, particularly with

the advent and rise of robotic platforms such as da Vinci®. Deep convolutional neural networks (CNNs) have emerged as the dominant solutions, surpassing traditional schemes by delivering automatic and highly accurate results. There is a line of advanced CNN models [2–5, 18–20] showing promising performance in this task. Nevertheless, most existing methods are introduced in the centralized learning setting where a large-scale dataset is collected for training. In reality, data centralization is often limited due to privacy-related issues, especially for medical data like surgical sequences. By contrast, the approach of learning from distributed datasets, which aligns more with practical scenarios, remains less explored. In our work, we revisit the challenge of instrument segmentation within a practical FL setting.

### 2.2 | Learning segmentation from imperfect annotations

In real-world scenarios, image datasets are often accompanied by imperfect labels, such as scarce, noisy, and scribble annotations [7]. The quality of annotations is a crucial factor influencing the performance of learned models, making it

imperative to leverage these inferior annotations. For instance, semi-supervised segmentation [8, 9, 21] endeavours to make the most use of unlabelled data. Scribble-supervised segmentation [12, 13] makes an effort to learn a model from sparse supervision signals. Besides, other efforts like designing a noise-tolerance loss function [11] and correcting noise labels during training [10] have been studied to handle label noises. Figure 1g through Figure 1j show the comparison among different segmentation diagrams. Although numerous attempts have been made to deal with inferior annotations, handling various imperfect annotations simultaneously within a more pragmatic FL framework is more practical but remains unexplored.

## 2.3 | Federated learning

FL [22–24] is a machine learning mechanism that enables multiple clients or devices to collaboratively learn a statistical model while keeping their data localized, effectively addressing privacy concerns in sensitive domains like healthcare. Despite the widespread application of FL in medical imaging [25–27], learning from various imperfect annotations remains unexplored in FL, where we need to tackle the challenges related to imperfect annotations of local datasets. We posit that investigating this problem is crucial, as it presents a practical and more challenging situation.

## 3 | METHOD

### 3.1 | Overview

Assuming there are  $K$  clients, denoted as  $\{C_i \mid i = 1, 2, 3, \dots, K\}$ , where each client  $C_i$  holds a private dataset. Ideally, each dataset is expected to contain images and the corresponding ideal annotations (ground truth). However, in reality, local datasets may come with imperfect annotations, such as scarce, noisy, and scribble annotations. Our goal is to learn a generalizable model with these distributed datasets with imperfect annotations.

In our study, we represent  $\mathbb{I}_i = \{\mathbf{X}_i^{\mathbb{I}}, \mathbf{Y}_i^{\mathbb{I}}\}$  as the local dataset consisting of images  $\mathbf{X}_i^{\mathbb{I}}$  and the corresponding ideal annotations  $\mathbf{Y}_i^{\mathbb{I}}$ . For the local dataset with scarce annotations, namely, the dataset comprising a small part of labelled images and a large part of unlabelled images, we denote it as  $\mathbb{S}_i = \left\{ \mathbf{X}_i^{\mathbb{S}}, \mathbf{Y}_i^{\mathbb{S}}, \bar{\mathbf{X}}_i^{\mathbb{S}} \right\}$ , where  $\mathbf{X}_i^{\mathbb{S}}$  and  $\mathbf{Y}_i^{\mathbb{S}}$  represent the labelled images and the corresponding ground truth, and  $\bar{\mathbf{X}}_i^{\mathbb{S}}$  indicates the unlabelled images. We further represent  $\mathbb{N}_i = \{\mathbf{X}_i^{\mathbb{N}}, \mathbf{Y}_i^{\mathbb{N}}\}$  as the local dataset containing images  $\mathbf{X}_i^{\mathbb{N}}$  and the related noisy annotations  $\mathbf{Y}_i^{\mathbb{N}}$ ,  $\mathbb{W}_i = \{\mathbf{X}_i^{\mathbb{W}}, \mathbf{Y}_i^{\mathbb{W}}\}$  as the local dataset comprising images  $\mathbf{X}_i^{\mathbb{W}}$  and the related scribble annotations  $\mathbf{Y}_i^{\mathbb{W}}$ .

Following the spirit of FL, we maintain a central server  $\mathcal{Q}$  to coordinate these local clients  $\{C_i \mid i = 1, 2, 3, \dots, K\}$  and their local models  $\{l_i \mid i = 1, 2, 3, \dots, K\}$  to develop a global model  $g$ . Each local model  $l_i$  is with its own weights  $\mathbf{w}_i^l$  and the global model owns the weights  $\mathbf{w}^g$ .

The flowchart of our framework is depicted in Figure 2. For simplicity and ease of understanding, four local clients  $\{C_i \mid i = 1, 2, 3, 4\}$  and one central server  $\mathcal{Q}$  are illustrated. Within the framework, local clients  $\{C_i \mid i = 1, 2, 3, 4\}$  train their respective local models  $\{l_i \mid i = 1, 2, 3, 4\}$  using their private datasets  $\mathbb{I}_1 = \{\mathbf{X}_1^{\mathbb{I}}, \mathbf{Y}_1^{\mathbb{I}}\}$ ,  $\mathbb{S}_2 = \left\{ \mathbf{X}_2^{\mathbb{S}}, \mathbf{Y}_2^{\mathbb{S}}, \bar{\mathbf{X}}_2^{\mathbb{S}} \right\}$ ,  $\mathbb{N}_3 = \{\mathbf{X}_3^{\mathbb{N}}, \mathbf{Y}_3^{\mathbb{N}}\}$ , and  $\mathbb{W}_4 = \{\mathbf{X}_4^{\mathbb{W}}, \mathbf{Y}_4^{\mathbb{W}}\}$ . These datasets come with various annotations, that is, ideal annotations (ground truth), scarce annotations (a large part of samples remain unlabelled), noisy annotations, and scribble annotations. The central server  $\mathcal{Q}$  orchestrates the process of model aggregation and updates, coordinating local clients to develop a generalizable global model  $g$ .

### 3.2 | Local training with various annotations

#### 3.2.1 | Learning from ideal annotations

For a client  $C_i$  holding a dataset  $\mathbb{I}_i = \{\mathbf{X}_i^{\mathbb{I}}, \mathbf{Y}_i^{\mathbb{I}}\}$  that has ideal annotations (ground truth), such as client  $C_1$  in Figure 2, we can simply perform fully-supervised learning to train the local model. During this process, we input images  $\mathbf{X}_i^{\mathbb{I}}$  to the local model, and obtain probability maps  $\hat{\mathbf{Y}}_i^{\mathbb{I}}$ . We can then calculate the supervised loss, for example, the cross-entropy (CE) loss  $\mathcal{L}_{ce}(\hat{\mathbf{Y}}_i^{\mathbb{I}}, \mathbf{Y}_i^{\mathbb{I}})$  between the probability maps  $\hat{\mathbf{Y}}_i^{\mathbb{I}}$  and the ground truth  $\mathbf{Y}_i^{\mathbb{I}}$ . Thus, in this setting, the training objective  $\mathcal{L}_i^{\mathbb{I}}$  can be expressed as

$$\mathcal{L}_i^{\mathbb{I}} = \mathcal{L}_{ce}(\hat{\mathbf{Y}}_i^{\mathbb{I}}, \mathbf{Y}_i^{\mathbb{I}}) = - \sum_{j=1}^J \sum_{p=1}^P \mathbf{Y}_{i,j,p}^{\mathbb{I}} \log(\hat{\mathbf{Y}}_{i,j,p}^{\mathbb{I}}), \quad (1)$$

where  $J$  is the number of classes,  $P$  indicates the number of pixels for each image, and  $\mathbf{Y}_{i,j,p}^{\mathbb{I}}$  and  $\hat{\mathbf{Y}}_{i,j,p}^{\mathbb{I}}$  represent the label value and probability value of the  $p$ th pixel within the  $j$ th channel of the ground truth  $\mathbf{Y}_i^{\mathbb{I}}$  and probability maps  $\hat{\mathbf{Y}}_i^{\mathbb{I}}$ .

#### 3.2.2 | Learning from scarce annotations

A local client  $C_i$  such as  $C_2$  in Figure 2 may hold a dataset  $\mathbb{S}_i = \left\{ \mathbf{X}_i^{\mathbb{S}}, \mathbf{Y}_i^{\mathbb{S}}, \bar{\mathbf{X}}_i^{\mathbb{S}} \right\}$  with scarce annotations. We employ the well-regarded mean-teacher (MT) [15] framework to train the local model  $l_i$ . Specifically, we sustain an additional teacher model  $\bar{l}_i$  alongside the local model.  $\bar{l}_i$  shares the same architecture as  $l_i$  but maintains the exponential moving average (EMA) weights  $\bar{\mathbf{w}}_i^l$  of  $\mathbf{w}_i^l$ . During training, for the labelled data, images  $\mathbf{X}_i^{\mathbb{S}}$  are fed into  $l_i$  to get the probability maps  $\hat{\mathbf{Y}}_i^{\mathbb{S}}$ , and the CE loss  $\mathcal{L}_{ce}(\hat{\mathbf{Y}}_i^{\mathbb{S}}, \mathbf{Y}_i^{\mathbb{S}})$  is then calculated between  $\hat{\mathbf{Y}}_i^{\mathbb{S}}$  and the ground truth  $\mathbf{Y}_i^{\mathbb{S}}$ . For the unlabelled data, images  $\bar{\mathbf{X}}_i^{\mathbb{S}}$  are input into  $l_i$  and  $\bar{l}_i$  to get probability maps  $\bar{\mathbf{Y}}_i^{\mathbb{S}}$  and  $\tilde{\mathbf{Y}}_i^{\mathbb{S}}$ . We then impose a consistency loss, such as the mean square error (MSE) loss  $\mathcal{L}_{mse}(\bar{\mathbf{Y}}_i^{\mathbb{S}}, \tilde{\mathbf{Y}}_i^{\mathbb{S}})$ , between  $\bar{\mathbf{Y}}_i^{\mathbb{S}}$  and  $\tilde{\mathbf{Y}}_i^{\mathbb{S}}$ . Thus, the training



objective  $\mathcal{L}_i^{\mathbb{S}}$  in this setting is written as

$$\begin{aligned} \mathcal{L}_i^{\mathbb{S}} &= \mathcal{L}_{ce}(\hat{\mathbf{Y}}_i^{\mathbb{S}}, \mathbf{Y}_i^{\mathbb{S}}) + \alpha \mathcal{L}_{mse}(\hat{\mathbf{Y}}_i^{\mathbb{S}}, \tilde{\mathbf{Y}}_i^{\mathbb{S}}) = \\ &= - \sum_{j=1}^J \sum_{p=1}^P \mathbf{Y}_{i,j,p}^{\mathbb{S}} \log(\hat{\mathbf{Y}}_{i,j,p}^{\mathbb{S}}) + \alpha \sum_{j=1}^J \sum_{p=1}^P \left| \hat{\mathbf{Y}}_{i,j,p}^{\mathbb{S}} - \tilde{\mathbf{Y}}_{i,j,p}^{\mathbb{S}} \right|^2, \end{aligned} \quad (2)$$

where  $J$  indicates the number of classes and  $P$  denotes the number of pixels for each image, namely,  $\mathbf{Y}_{i,j,p}^{\mathbb{S}}$ ,  $\hat{\mathbf{Y}}_{i,j,p}^{\mathbb{S}}$ ,  $\tilde{\mathbf{Y}}_{i,j,p}^{\mathbb{S}}$ , and  $\tilde{\mathbf{Y}}_{i,j,p}^{\mathbb{S}}$  respectively, indicate the value of the  $p$ th pixel at the  $j$ th channel in  $\mathbf{Y}_i^{\mathbb{S}}$ ,  $\hat{\mathbf{Y}}_i^{\mathbb{S}}$ ,  $\tilde{\mathbf{Y}}_i^{\mathbb{S}}$ , and  $\tilde{\mathbf{Y}}_i^{\mathbb{S}}$ .  $\alpha$  is a trade-off hyperparameter whose value progressively increases from 0 to 0.1, guided by a time-dependent function  $\alpha(t) = 0.1 \exp \left[ -5 \left( 1 - \frac{t}{T} \right)^2 \right]$ . Here,  $t$  represents the current training epoch, and  $T$  denotes the maximum training epoch. Besides, a strong augmentation method, CutMix [28], is introduced during training to enhance the performance further.

### 3.2.3 | Learning from noisy annotations

A client  $C_i$ , such as  $C_3$  in Figure 2, may have a dataset  $\mathbb{N}_i = \{\mathbf{X}_i^{\mathbb{N}}, \mathbf{Y}_i^{\mathbb{N}}\}$  with noisy annotations. To train a local model from this noisy dataset, we employ the noise-robust Dice (NR-Dice) loss  $\mathcal{L}_{nd}$  [11], a generalization of dice loss and mean absolute error (MAE) loss. Assuming the probability maps for images  $\mathbf{X}_i^{\mathbb{N}}$  output by the local model are denoted by  $\hat{\mathbf{Y}}_i^{\mathbb{N}}$ , we optimize the training objective  $\mathcal{L}_i^{\mathbb{N}}$  of local model by combining the CE loss  $\mathcal{L}_{ce}$  and the NR-Dice loss  $\mathcal{L}_{nd}$ , defined as

$$\begin{aligned} \mathcal{L}_i^{\mathbb{N}} &= \mathcal{L}_{ce}(\hat{\mathbf{Y}}_i^{\mathbb{N}}, \mathbf{Y}_i^{\mathbb{N}}) + \beta \mathcal{L}_{nd}(\hat{\mathbf{Y}}_i^{\mathbb{N}}, \mathbf{Y}_i^{\mathbb{N}}) = - \sum_{j=1}^J \sum_{p=1}^P \mathbf{Y}_{i,j,p}^{\mathbb{N}} \\ &= - \sum_{j=1}^J \sum_{p=1}^P \mathbf{Y}_{i,j,p}^{\mathbb{N}} \log(\hat{\mathbf{Y}}_{i,j,p}^{\mathbb{N}}) + \beta \sum_{j=1}^J \frac{\sum_{p=1}^P \left| \mathbf{Y}_{i,j,p}^{\mathbb{N}} - \hat{\mathbf{Y}}_{i,j,p}^{\mathbb{N}} \right|^{\gamma}}{\sum_{p=1}^P \left( \mathbf{Y}_{i,j,p}^{\mathbb{N}} \right)^2 + \sum_{p=1}^P \left( \hat{\mathbf{Y}}_{i,j,p}^{\mathbb{N}} \right)^2 + \epsilon}, \end{aligned} \quad (3)$$

in which  $J$  is the number of classes,  $P$  is the number of pixels for each image, and  $\mathbf{Y}_{i,j,p}^{\mathbb{N}}$  and  $\hat{\mathbf{Y}}_{i,j,p}^{\mathbb{N}}$  represent the label value and probability value of the  $p$ th pixel within the  $j$ th channel of the noisy labels  $\mathbf{Y}_i^{\mathbb{N}}$  and probability maps  $\hat{\mathbf{Y}}_i^{\mathbb{N}}$ .  $\beta$  and  $\gamma \in [1, 2]$  are hyperparameters and  $\epsilon$  is a small constant to avoid zero-division. We set  $\gamma$  to 1.5, and  $\epsilon$  to  $10^{-5}$ . We gradually decrease the value of  $\beta$  from 0.1 to 0 with a time-based function  $\beta(t) = 0.1 \left\{ 1 - \exp \left[ -5 \left( 1 - \frac{t}{T} \right)^2 \right] \right\}$ , where  $t$  indicates the current training epoch and  $T$  is the maximum training epoch.

### 3.2.4 | Learning from scribble annotations

For a client  $C_i$  like  $C_4$  in Figure 2 that possesses a dataset  $\mathbb{W}_i = \{\mathbf{X}_i^{\mathbb{W}}, \mathbf{Y}_i^{\mathbb{W}}\}$ , in which scribble annotations provide much fewer supervision signals than the ideal annotations, we adopt the partial cross-entropy (pCE) [16] for training, similar to the work [12] for weakly-supervised instrument segmentation. Let the probability maps yielded by the local model for images  $\mathbf{X}_i^{\mathbb{W}}$  be  $\hat{\mathbf{Y}}_i^{\mathbb{W}}$ . We calculate the training objective  $\mathcal{L}_i^{\mathbb{W}}$  as the pCE loss  $\mathcal{L}_{pce}(\hat{\mathbf{Y}}_i^{\mathbb{W}}, \mathbf{Y}_i^{\mathbb{W}})$  between  $\hat{\mathbf{Y}}_i^{\mathbb{W}}$  and  $\mathbf{Y}_i^{\mathbb{W}}$ , written as

$$\mathcal{L}_i^{\mathbb{W}} = \mathcal{L}_{pce}(\hat{\mathbf{Y}}_i^{\mathbb{W}}, \mathbf{Y}_i^{\mathbb{W}}) = - \sum_{j=1}^J \sum_{b \in H} \mathbf{Y}_{i,j,b}^{\mathbb{W}} \log(\hat{\mathbf{Y}}_{i,j,b}^{\mathbb{W}}), \quad (4)$$

where  $J$  is the number of classes and  $H$  is the set of indexes of the annotated pixels.  $\hat{\mathbf{Y}}_{i,j,b}^{\mathbb{W}}$  and  $\mathbf{Y}_{i,j,b}^{\mathbb{W}}$  denote the probability value and label value of the  $b$ th annotated pixel in the  $j$ th channel of probability maps  $\hat{\mathbf{Y}}_i^{\mathbb{W}}$  and scribble annotations  $\mathbf{Y}_i^{\mathbb{W}}$ .

## 3.3 | Federated learning module

Following the standard FL paradigm, we involve a central server  $\mathcal{Q}$  for model aggregation and updating. Concretely, at each global round  $r$ , local clients  $\{C_i \mid i = 1, 2, 3, \dots, K\}$  upload local model weights  $\{\mathbf{w}_i^r \mid i = 1, 2, 3, \dots, K\}$  to the central server  $\mathcal{Q}$  after local training, then the central server  $\mathcal{Q}$  aggregates local weights to update the global model weight  $\mathbf{w}^g$ . Afterward, local clients download the global model weights, assign them to their local models, and then fine-tune them with their local datasets. By repeating this procedure until convergence, we can develop a generalizable global model  $g$ . In our study, we adopted the FedAvg scheme [22] to update  $\mathbf{w}^g$ , written as

$$\mathbf{w}_{r+1}^g = \sum_{i=1}^K \frac{U_i}{\sum_{i=1}^K U_i} \mathbf{w}_{i,r}^r, \quad (5)$$

where  $U_i$  denotes the dataset size of client  $C_i$ .

## 4 | EXPERIMENTS AND RESULTS

### 4.1 | Experimental setup

#### 4.1.1 | Dataset and metric

We validated our method on the publicly available endoscopic dataset EndoVis17, provided by the 2017 robotic instrument segmentation challenge [1]. The EndoVis17 dataset consists of ten sequences from abdominal porcine procedures along with the corresponding ground truth for binary, multi-class, and multi-part segmentation tasks. All images are in  $1920 \times 1080$  pixel resolution. Our study focused on the binary

**TABLE 1** Detailed information on dataset setup for local clients.  $\mathcal{I}$ : ideal annotation type,  $\mathcal{S}$ : scarce annotation type,  $\mathcal{N}$ : noisy annotation type,  $\mathcal{W}$ : scribble annotation type,  $\xi$ : labelled sample ratio, and  $\delta$ : noisy label ratio.

Client	Annotation type	Sequence	Training sample	Testing sample	Notation
$C_1$	$\mathcal{I}$	{1, 2}	First 225 frames of {1, 2}	Last 75 frames of {1, 2}	—
$C_2$	$\mathcal{S}$	{3, 4}	First 225 frames of {3, 4}	{1, 2, 3, 4, 5, 6, 7, 8}	$\xi = \{0.2, 0.4\}$
$C_3$	$\mathcal{N}$	{5, 6}	First 225 frames of {5, 6}		$\delta = \{0.2, 0.4, 0.6, 0.8\}$
$C_4$	$\mathcal{W}$	{7, 8}	First 225 frames of {7, 8}		—

segmentation task. We utilized the former eight sequences, denoted as {1, 2, 3, 4, 5, 6, 7, 8}. As suggested by the work [1], we used the first 225 frames from each sequence for training and the remaining 75 frames for testing. We employed the intersection over union (IoU) as the evaluation metric.

#### 4.1.2 | Problem simulation

For the FL framework, we established a central server  $Q$  and four local clients  $\{C_i \mid i = 1, 2, 3, 4\}$ . Concretely, we allocated the first 225 frames of sequences {1, 2}, {3, 4}, {5, 6}, and {7, 8} to  $C_1$ ,  $C_2$ ,  $C_3$ , and  $C_4$ , respectively. We constructed the test dataset using the last 75 frames of sequences {1, 2, 3, 4, 5, 6, 7, 8}. To mimic scenarios with scarce annotations<sup>3</sup>, we treated 20% and 40% of the data as labelled data, leaving the remaining 80% and 60% as unlabelled data for  $C_2$ . To generate noisy annotations<sup>4</sup> for  $C_3$ , we utilized dilation and erosion operations to create noisy annotations for 20%, 40%, 60%, and 80% of the samples. The operations were carried out with a  $5 \times 5$  all-ones matrix kernel. The iteration number  $\tau$  was randomly sampled from the range of [15, 20] for each image. For  $C_4$ , we used the skeletonization method [14] on the ground truth to obtain scribble annotations. Detailed information on data setup for local clients is described in Table 1. For the sake of simplicity, we represent the ideal, scarce, noisy, and scribble annotation types as  $\mathcal{I}$ ,  $\mathcal{S}$ ,  $\mathcal{N}$ , and  $\mathcal{W}$ .

#### 4.1.3 | Implementation details

We preprocessed the data by cropping all images from the pixel at (320, 28) to achieve a pixel resolution of  $1280 \times 1024$ , as described in the approach [2]. Subsequently, we resized all images to  $320 \times 256$  by a factor of 16. We adopted the U-Net model [18] as the backbone for training. We trained local models using the Adam optimizer with a learning rate of  $10^{-4}$ . The batch size was set to 16. We ran the local training for 1 epoch and repeated the global round 200 times. We employed horizontal and vertical flips (with a probability of 0.5 to flip) as data augmentation strategies. Alongside our method, we also executed experiments for **standalone** (local

training) and **centralization** (centralized training) for comparison. For **standalone**, we conducted local training with 100 epochs for local clients using their private datasets. For **centralization**, we pooled all local datasets, constructed individual data loaders for these datasets, and performed centralized training by iterating through these data loaders with 200 epochs. Each data loader with a specific annotation type corresponded to a specific scheme from Section 3.2. In our study, all experiments were conducted three times with different random seeds.

## 4.2 | Experimental results

### 4.2.1 | Quantitative results

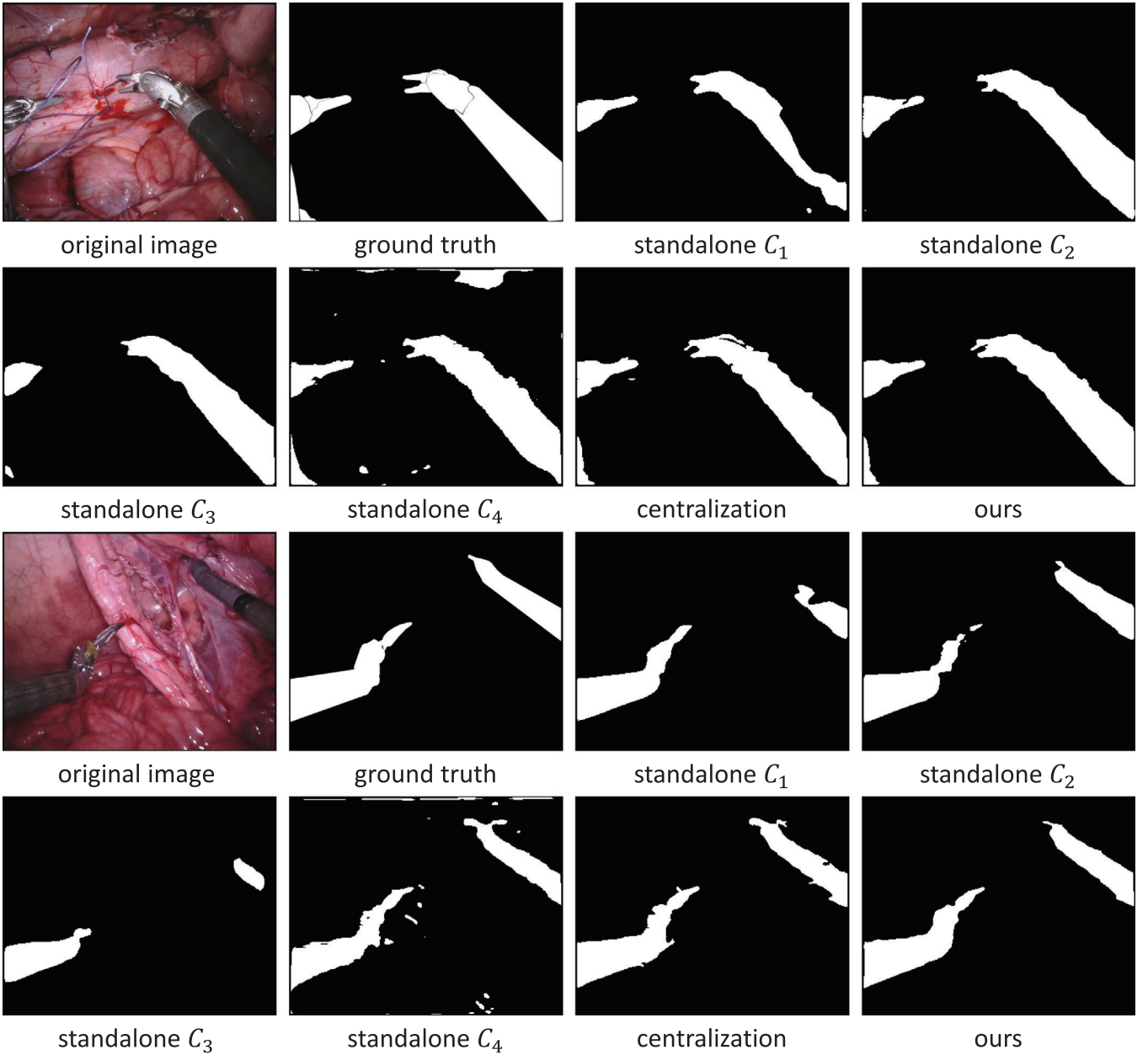
We first show the quantitative results of **standalone** for each client, as shown in Table 2.  $C_1$  trains the local model using its dataset with ideal annotations and achieves an IoU score of 71.02%. For  $C_2$ ,  $C_3$ , and  $C_4$  holding datasets with imperfect annotations, the upper-bound results obtained with ideal annotations are also reported for reference. We can note that imperfect annotations can negatively affect the model performance compared to ideal annotations.  $C_2$  makes use of its dataset with scarce annotations to train its local model and realizes IoU scores of 79.40%, 79.65%, and 79.29%, respectively, when using ideal annotations, scarce annotations with labelled data ratios of 0.2 and 0.4.  $C_3$  brings IoU scores of 81.56%, 78.45%, 76.59%, 69.56%, and 73.82%, respectively, under the settings of training with ideal annotations, noisy label ratios of 0.2, 0.4, 0.6, and 0.8.  $C_4$  learns its local model from scribble annotations and achieves an accuracy of 64.09% in IoU, while the upper-bound accuracy obtained by training with ideal annotations is 73.87% in IoU.

Table 3 shows the comparison between **ours** and **centralization** under various data settings. Firstly, we can observe that **ours** and **centralization** consistently improves **standalone** under various dataset settings, providing evidence that taking advantage of more data helps improve model accuracy and generalization ability.

Interestingly, when comparing **centralization** to **ours**, we note that **centralization** surpasses **ours** by about 0.53% in IoU under the setting where both train models with datasets with ideal annotations, but realizes much worse performance under the other settings where imperfect

<sup>3</sup> We define  $\xi$  as the ratio of labelled images. For instance,  $\xi = 0.2$  signifies that 20% of the images are labelled, leaving the remaining 80% unlabelled.

<sup>4</sup>  $\delta$  is used to represent the ratio of noisy annotations. For example,  $\delta = 0.2$  denotes that 20% of the images contain noisy annotations, while the other 80% have perfect annotations.



**FIGURE 3** Qualitative visualization results. Segmentation results are from the settings of  $C_1$ :  $\mathcal{I}$ ,  $C_2$ :  $\mathcal{S}$  ( $\xi = 0.2$ ),  $C_3$ :  $\mathcal{N}$  ( $\delta = 0.8$ ), and  $C_4$ :  $\mathcal{W}$ .  $\mathcal{I}$ : ideal annotation type,  $\mathcal{S}$ : scarce annotation type,  $\mathcal{N}$ : noisy annotation type,  $\mathcal{W}$ : scribble annotation type,  $\xi$ : labelled sample ratio, and  $\delta$ : noisy label ratio.

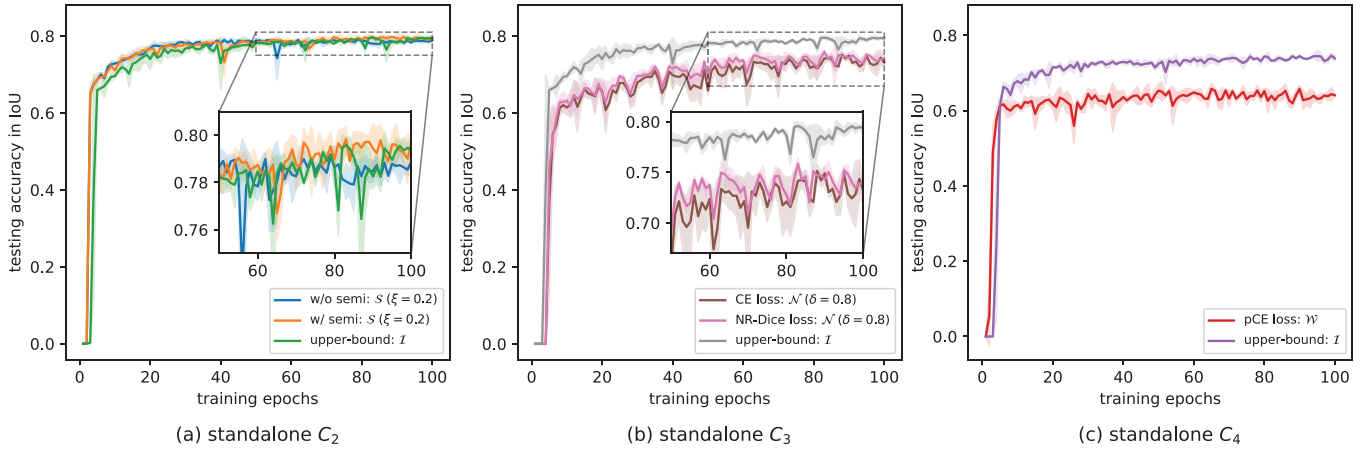
annotations are involved. The case of **centralization** outperforming **ours** using ideal annotations is as anticipated since we believe **centralization** should benefit from the consolidation of all available data, allowing the model to learn a more comprehensive representation of the data. However, the decrease in performance of **centralization** with various imperfect annotations presents an intriguing phenomenon.

We hypothesize that the inconsistency of annotation types and their corresponding learning strategies could be a potential explanation. In the centralized setup, the model might be overly influenced by the majority trend in the data, which

could lead to a “rich get richer” effect, often ignoring or misrepresenting the minority class or outliers. This effect can boost performance when ideal annotations are used due to the high consistency of the data. However, with the introduction of various imperfect annotations, the data becomes inconsistent and diverse. Moreover, each type of imperfect annotation requires a different learning strategy, leading to inconsistent strategies for different data loaders. The single centralized model may not adapt well to the diverse learning strategies, thus struggling with inconsistent or conflicting information, which results in a decline in performance. Contrarily, due to its distributed spirit, **ours** with FL might be more robust







**FIGURE 4** Ablation study on effectiveness of adopted methods by local clients in handling various imperfect annotations.  $I$ : ideal annotation type,  $S$ : scarce annotation type,  $N$ : noisy annotation type,  $W$ : scribble annotation type,  $\xi$ : labelled sample ratio, and  $\delta$ : noisy label ratio. IoU: intersection over union. CE: cross-entropy. NR-Dice: noise-robust Dice. pCE: partial cross-entropy. w/o: without. w/: with.

against such inconsistencies. In **ours**, each local client learns a model based on its local data, and these local models are then combined by a global model. This process allows each client to focus on their specific data subset, accommodating the local characteristics and imperfections of the data. Besides, the aggregation of these models might provide a better balance between different learning strategies, leading to enhanced performance.

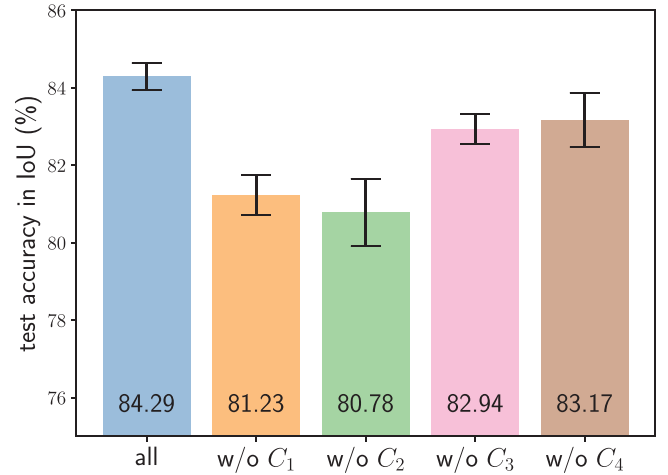
#### 4.2.2 | Qualitative results

In addition to the quantitative results, we present a qualitative analysis of different methods, as shown in Figure 3. We show-case two examples from the settings of  $C_1$ :  $I$ ,  $C_2$ :  $S$  ( $\xi = 0.2$ ),  $C_3$ :  $N$  ( $\delta = 0.8$ ), and  $C_4$ :  $W$ . We can note that our method produces results considerably closer to the ground truth, even under this challenging data situation. It demonstrates that our method can effectively learn from distributed datasets with imperfect annotations and generate accurate predictions.

### 4.3 | Analytical studies

#### 4.3.1 | Effectiveness in handling various imperfect annotations

We first evaluated the efficacy of adopted methods, that is, the MT framework [15], NR-Dice loss [11], and pCE loss [16] in handling imperfect annotations, since local clients should contribute effective local models for global model preparation. We show the results of **standalone** for  $C_2$ ,  $C_3$ , and  $C_4$  in Figure 4, where the testing accuracy in IoU of each training epoch is recorded. We can observe that the adopted methods present successful solutions. Specifically, the MT scheme improves the baseline when training with very scarce annotations ( $\xi = 0.2$ ), as shown in Figure 4a. Besides, the NR-Dice loss shows better noise tolerance than the CE loss in a very high noise level



**FIGURE 5** Ablation study on the contribution of each client for global model performance under the settings of  $C_1$ :  $I$ ,  $C_2$ :  $S$  ( $\xi = 0.2$ ),  $C_3$ :  $N$  ( $\delta = 0.8$ ), and  $C_4$ :  $W$ .  $I$ : ideal annotation type.  $S$ : scarce annotation type.  $N$ : noisy annotation type.  $W$ : scribble annotation type.  $\xi$ : labeled sample ratio.  $\delta$ : noisy label ratio. IoU: intersection over union. w/o: without.

setting ( $\delta = 0.8$ ), as illustrated in Figure 4b. In addition, the pCE loss presents its potential solution in learning with scribble annotations, as depicted in Figure 4c. It achieves an accuracy of about 64.09% in IoU, which is about 9.78% lower than the upper-bound accuracy.

#### 4.3.2 | Contribution of each client

We then analyzed the contribution of each client for global model performance under the settings of  $C_1$ :  $I$ ,  $C_2$ :  $S$  ( $\xi = 0.2$ ),  $C_3$ :  $N$  ( $\delta = 0.8$ ), and  $C_4$ :  $W$ . As shown in Figure 5, we can observe that combining all clients achieves the best performance while excluding any single client leads to accuracy degradation. The global model's performance decreases by approximately 3.06%, 3.51%, 1.35%, and 1.12% in IoU when excluding  $C_1$ ,  $C_2$ ,  $C_3$ , and  $C_4$ , respectively.

**TABLE 4** Ablation study on impact of various noise levels simulated with different ranges for dilation/erosion iteration  $\tau$  and different values of noisy label ratio  $\delta$ .  $\mathcal{I}$ : ideal annotation type. We chose  $\tau \in [15, 20]$  for noisy label simulation in our study.

Ranges	$\delta = 0.2$	$\delta = 0.4$	$\delta = 0.6$	$\delta = 0.8$
$\tau \in [1, 5]$	$81.05 \pm 0.87$	$81.00 \pm 0.34$	$80.32 \pm 1.05$	$80.40 \pm 0.68$
$\tau \in [5, 10]$	$79.98 \pm 1.21$	$78.72 \pm 1.11$	$78.49 \pm 0.43$	$78.36 \pm 1.29$
$\tau \in [10, 15]$	$79.42 \pm 1.26$	$77.70 \pm 1.12$	$73.94 \pm 0.85$	$74.82 \pm 1.76$
$\tau \in [15, 20]$	$78.45 \pm 1.32$	$76.59 \pm 1.69$	$69.56 \pm 0.97$	$73.82 \pm 1.82$
Upper-bound: $\mathcal{I}$	$81.56 \pm 0.26$			

This demonstrates that each client contributes to the overall performance.

### 4.3.3 | Choice of noise level for noisy annotation simulation

We further studied the impact of various noise levels simulated with different ranges for dilation/erosion iteration  $\tau$  and distinct values of noisy label ratio  $\delta$ . Larger values for both  $\tau$  and  $\delta$  lead to an increased level of noise. As expected, our experimental results, shown in Table 4, demonstrate a decline in model accuracy in correlation with increasing noise levels. It is well-accepted that a small degree of dilation/erosion, for example, with  $\tau \in [1, 5]$ , can better mimic the uncertainty and ambiguity often experienced during the annotation process, especially near the boundaries of objects. Interestingly, our experiments indicate that the model’s performance remains relatively stable even with smaller iteration ranges, suggesting that the model is resilient to a certain degree of annotation noise in binary instrument segmentation since it is a relatively easy task. However, in this study, we intentionally opted for a more extensive iteration range for our noisy simulations, that is,  $\tau \in [15, 20]$ . We believe that this decision allows us to thoroughly evaluate and highlight the robustness of our method under more challenging conditions where the noise level is substantially increased, thus providing a more rigorous and convincing demonstration of its capabilities.

## 5 | DISCUSSION AND CONCLUSION

We aim to learn a generalizable model for instrument segmentation from decentralized surgical videos with various imperfect annotations. In practice, surgical datasets are usually highly siloed with individual hospitals or medical institutions due to privacy concerns. Besides, these datasets often come with various imperfect annotations. Most existing methods ignore these realities and focus on centralized and well-annotated data, limiting their applicability. This problem relates to data availability and quality issues in the real world and calls for more robust, widely applicable models for robotic surgery.

Our method unifies the semi-, noisy label-, and scribble-supervised segmentation paradigms and the FL scheme into a single framework. Our method handles data privacy issues and effectively learns from imperfectly annotated data, accommo-

dating diverse annotation scenarios. Our method outperformed **standalone** and **centralization** under various imperfect annotation settings, demonstrating its successful solution for this new and challenging problem. We posit that the variation in annotation types and their respective learning strategies might account for the diminished performance in centralized learning. However, this hypothesis warrants deeper exploration, including examining intermediate training patterns and offering a more detailed analysis of the relationship between different strategies and performance outcomes.

Despite its efficacy, our method has several limitations that warrant future investigation. Firstly, a gap exists between our method and its real-world implementations. One aspect of this gap stems from our use of simulated imperfect annotations, such as noisy annotations and scribble annotations. Although imperfect annotation simulation is widely applied in existing methods [10, 29], authenticating our model with real-world imperfect annotations will further fortify the validity of our approach. Additionally, our method relied on the simplified assumption that each client possessed only a single type of annotation. Future research should address the challenges associated with each client having multiple annotation types to better reflect the intricacies of real-world scenarios, thus ensuring a deeper understanding and a more comprehensive solution. Secondly, our method relied on the conventional FedAvg scheme [22]. However, considering that each client possesses distinct datasets and diverse types of imperfect annotations, the performance of local models can exhibit significant variance. This raises the question of whether naive aggregation with FedAvg in our study would inadvertently impact the robustness and fairness of the global model [30, 31]. Therefore, it is essential to investigate the effects of local model performance variance on the global model in our context and to ascertain whether more tailored aggregation strategies could yield improved performance. Thirdly, our method was only evaluated with binary instrument segmentation. To assess the generalization ability of our method, further validation with other instrument segmentation tasks and additional surgical datasets is also required.

### AUTHOR CONTRIBUTIONS

**Zhou Zheng**: Conceptualization; data curation; methodology; software; validation; writing—original draft; writing—review and editing. **Yuichiro Hayashi**: Methodology; writing—review and editing. **Masahiro Oda**: Methodology; supervision; writing—review and editing. **Takayuki Kitasaka**: Methodology; writing—review and editing. **Kensaku Mori**: Conceptualization; methodology; project administration; resources; supervision; writing—review and editing.

### ACKNOWLEDGEMENTS

This work was supported by the Japan Society for the Promotion of Science (JSPS) KAKENHI Grant Numbers 21K19898 and 17H00867, the Japan Science and Technology Agency (JST) Core Research for Evolutional Science and Technology (CREST) Grant Number JPMJCR20D5, and the Japan Science and Technology Agency (JST) Moonshot Research and Development (R&D) Program Grant Number JPMJMS2214-07, Japan.

## CONFLICT OF INTEREST STATEMENT

The authors declare no conflicts of interest.

## DATA AVAILABILITY STATEMENT

The data that support the findings of this study are available from the corresponding author upon reasonable request.

## ORCID

Zhou Zheng  <https://orcid.org/0009-0001-5593-2281>

## REFERENCES

- Allan, M., Shvets, A., Kurmann, T., Zhang, Z., Duggal, R., Su, Y.H., et al.: 2017 robotic instrument segmentation challenge. arXiv:190206426 (2019)
- Shvets, A.A., Rakhlin, A., Kalinin, A.A., Iglovikov, V.I.: Automatic instrument segmentation in robot-assisted surgery using deep learning. In: Proceedings of 2018 17th IEEE International Conference on Machine Learning and Applications, pp. 624–628. IEEE, Piscataway, NJ (2018)
- Iglovikov, V., Shvets, A.: Terausnet: U-NET with VGG11 encoder pre-trained on imagenet for image segmentation. arXiv:180105746 (2018)
- Jin, Y., Cheng, K., Dou, Q., Heng, P.A.: Incorporating temporal prior from motion flow for instrument segmentation in minimally invasive surgery video. In: Proceedings of Medical Image Computing and Computer Assisted Intervention—MICCAI 2019, pp. 11768, 440–448. Springer, Cham (2019)
- Pakhomov, D., Premachandran, V., Allan, M., Azizian, M., Navab, N.: Deep residual learning for instrument segmentation in robotic surgery. In: Proceedings of Machine Learning in Medical Imaging, pp. 566–573. Springer, Cham (2019)
- Maier Hein, L., Vedula, S.S., Speidel, S., Navab, N., Kikinis, R., Park, A., et al.: Surgical data science for next-generation interventions. *Nat. Biomed. Eng.* 1(9), 691–696 (2017)
- Tajbakhsh, N., Jeyaseelan, L., Li, Q., Chiang, J.N., Wu, Z., Ding, X.: Embracing imperfect datasets: A review of deep learning solutions for medical image segmentation. *Med. Image Anal.* 63, 101693 (2020)
- Zhao, Z., Jin, Y., Gao, X., Dou, Q., Heng, P.A.: Learning motion flows for semi-supervised instrument segmentation from robotic surgical video. In: Proceedings of Medical Image Computing and Computer Assisted Intervention—MICCAI 2020, pp. 679–689. Springer, Cham (2020)
- Yang, H., Shan, C., Bouwman, A., Dekker, L.R.C., Kolen, A.F., de With, P.H.N.: Medical instrument segmentation in 3D us by hybrid constrained semi-supervised learning. *IEEE J. Biomed. Health. Inf.* 26(2), 762–773 (2022)
- Xue, C., Deng, Q., Li, X., Dou, Q., Heng, P.A.: Cascaded robust learning at imperfect labels for chest X-ray segmentation. In: Proceedings of Medical Image Computing and Computer Assisted Intervention—MICCAI 2020, pp. 579–588. Springer, Cham (2020)
- Wang, G., Liu, X., Li, C., Xu, Z., Ruan, J., Zhu, H., et al.: A noise-robust framework for automatic segmentation of Covid-19 pneumonia lesions from CT images. *IEEE Trans. Med. Imaging* 39(8), 2653–2663 (2020)
- Yang, Z., Simon, R., Linte, C.: A weakly supervised learning approach for surgical instrument segmentation from laparoscopic video sequences. In: Proceedings of Medical Imaging 2022: Image-Guided Procedures, Robotic Interventions, and Modeling, pp. 120341U. SPIE, Bellingham, MA (2022)
- Fuentes Hurtado, F., Kadhodamohammadi, A., Flouty, E., Barbarisi, S., Luengo, I., Stoyanov, D.: Easylabels: Weak labels for scene segmentation in laparoscopic videos. *Int. J. Comput. Assisted Radiol. Surg.* 14, 1247–1257 (2019)
- Zhang, T.Y., Suen, C.Y.: A fast parallel algorithm for thinning digital patterns. *Commun. ACM* 27(3), 236–239 (1984)
- Tarvainen, A., Valpola, H.: Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. arXiv:170301780 (2018)
- Tang, M., Djelouah, A., Perazzi, F., Boykov, Y., Schroers, C.: Normalized cut loss for weakly-supervised CNN segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1818–1827. IEEE, Piscataway, NJ (2018)
- Ross, T., Reinke, A., Full, P.M., Wagner, M., Kenngott, H., Apitz, M., et al.: Robust medical instrument segmentation challenge 2019. arXiv:200310299 (2020)
- Ronneberger, O., Fischer, P., Brox, T.: U-NET: Convolutional networks for biomedical image segmentation. In: Proceedings of Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015, pp. 9351, 234–241. Springer, Cham (2015)
- Wang, A., Islam, M., Xu, M., Ren, H.: Rethinking surgical instrument segmentation: A background image can be all you need. In: Proceedings of Medical Image Computing and Computer Assisted Intervention—MICCAI 2022, pp. 355–364. Springer, Singapore (2022)
- Yang, L., Gu, Y., Bian, G., Liu, Y.: Tmf-net: A transformer-based multiscale fusion network for surgical instrument segmentation from endoscopic images. *IEEE Trans. Instrum. Meas.* 72, 1–15 (2022)
- French, G., Laine, S., Aila, T., Mackiewicz, M., Finlayson, G.: Semi-supervised semantic segmentation needs strong, varied perturbations. arXiv:190601916 (2019)
- McMahan, B., Moore, E., Ramage, D., Hampson, S., Arcas, yB.A.: Communication-efficient learning of deep networks from decentralized data. In: Proceedings of Artificial Intelligence and Statistics, pp. 1273–1282. PMLR, Microtome Publishing, Brookline, MA (2017)
- Kairouz, P., McMahan, H.B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A.N., et al.: Advances and open problems in federated learning. *Found. Trends Mach. Learn.* 14(1–2), 1–210 (2021)
- Zhang, C., Xie, Y., Bai, H., Yu, B., Li, W., Gao, Y.: A survey on federated learning. *Knowledge-Based Syst.* 216, 106775 (2021)
- Li, W., Milletari, F., Xu, D., Rieke, N., Hancox, J., Zhu, W., et al.: Privacy-preserving federated brain tumour segmentation. In: Proceedings of Machine Learning in Medical Imaging, pp. 133–141. Springer, Cham (2019)
- Rieke, N., Hancox, J., Li, W., Milletari, F., Roth, H.R., Albarqouni, S., et al.: The future of digital health with federated learning. *npj Digital Med.* 3(1), 119 (2020)
- Shen, C., Wang, P., Roth, H.R., Yang, D., Xu, D., Oda, M., et al.: Multi-task federated learning for heterogeneous pancreas segmentation. In: Proceedings of Clinical Image-Based Procedures, Distributed and Collaborative Learning, Artificial Intelligence for Combating COVID-19 and Secure and Privacy-Preserving Machine Learning, pp. 12969, 101–110. Springer, Cham (2021)
- Yun, S., Han, D., Oh, S.J., Chun, S., Choe, J., Yoo, Y.: Cutmix: Regularization strategy to train strong classifiers with localizable features. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 6023–6032. IEEE, Piscataway, NJ (2019)
- Gao, F., Hu, M., Zhong, M.E., Feng, S., Tian, X., Meng, X., et al.: Segmentation only uses sparse annotations: Unified weakly and semi-supervised learning in medical images. *Med. Image Anal.* 80, 102515 (2022)
- Li, T., Hu, S., Beirami, A., Smith, V.: Ditto: Fair and robust federated learning through personalization. In: Proceedings of International Conference on Machine Learning, pp. 6357–6368. PMLR, Microtome Publishing, Brookline, MA (2021)
- Jiang, M., Roth, H.R., Li, W., Yang, D., Zhao, C., Nath, V., et al.: Fair federated medical image segmentation via client contribution estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 16302–16311. IEEE, Piscataway, NJ (2023)

**How to cite this article:** Zheng, Z., Hayashi, Y., Oda, M., Kitasaka, T., Mori, K.: Revisiting instrument segmentation: Learning from decentralized surgical sequences with various imperfect annotations. *Healthc. Technol. Lett.* 11, 146–156 (2024). <https://doi.org/10.1049/htl2.12068>