

Blessing Few-Shot Segmentation via Semi-Supervised Learning with Noisy Support Images

Run tong Zhang^a, Hongyuan Zhu^b, Hanwang Zhang^c, Chen Gong^d, Joey
Tianyi Zhou^e, Fanman Meng^{a,*}

^a*University of Electronic Science and Technology of China, Chengdu, China*

^b*Institute for Infocomm Research (I²R) & Centre for Frontier AI Research (CFAR),
A*STAR, Singapore*

^c*Nanyang Technological University, Singapore*

^d*Nanjing University of Science and Technology, Nanjing, China*

^e*Centre for Frontier AI Research (CFAR), A*STAR, Singapore*

Abstract

Mainstream few-shot segmentation methods meet performance bottleneck due to the data scarcity of novel classes with insufficient intra-class variations, which results in a biased model primarily favoring the base classes. Fortunately, owing to the evolution of the Internet, an extensive repository of unlabeled images has become accessible from diverse sources such as search engines and publicly available datasets. However, such unlabeled images are not a free lunch. There are noisy inter-class and intra-class samples causing severe feature bias and performance degradation. Therefore, we propose a semi-supervised few-shot segmentation framework named **F4S**, which incorporates a ranking algorithm designed to eliminate noisy samples and select superior pseudo-labeled images, thereby fostering the improvement of few-shot segmentation within a semi-supervised paradigm. The proposed F4S framework can not only enrich the intra-class variations of novel classes during the test phase, but also enhance meta-learning of the network during the training phase. Furthermore, it can be readily implemented with ease on any off-the-shelf few-shot segmentation methods. Additionally, based on a Structural Causal Model (SCM), we further theoretically explain why the proposed method can solve the noise problem: the severe noise effects are removed by

*Corresponding author

Email address: fmmeng@uestc.edu.cn (Fanman Meng)

cutting off the backdoor path between pseudo labels and noisy support images via causal intervention. On PASCAL-5ⁱ and COCO-20ⁱ datasets, we show that the proposed F4S can boost various popular few-shot segmentation methods to new state-of-the-art performances.

Keywords: Few-shot segmentation, Semi-supervised learning, Noisy images, Causal inference

1 Introduction

Few-shot segmentation (FSS) [1] aims to segment the object regions in query images of novel classes using a minimal number (N-shot) of annotated support images. The most common experimental settings for FSS use 1-shot and 5-shot annotated support samples, as shown in Fig. 1 (a) and (b). The primary challenge for FSS is how to effectively utilize the information provided by the N-shot support images. Prototype-based approaches [2, 3, 4, 5, 6] focus on generating representative prototypes from the N-shot support images to accurately characterize the novel classes. In contrast, the metric-based approaches [7, 8, 9] focus on learning a class-agnostic similarity metric that can precisely measure the regions similar to the N-shot support regions in the query image. However, the most significant challenge of few-shot learning is how to maximize the exploration of data distributions under data scarcity [10]. Increasing manually annotated data is the most direct and effective method, but it is extremely time and labor-consuming.

Thanks to semi-supervised learning (SSL), the pseudo-labeling methods have provided a practical solution for the data scarcity issue in few-shot learning tasks, and there is already relevant research work published on this. For example, the method in [11] combines semi-supervised learning with few-shot classification and proposes the PLCM network, which generates and selects good pseudo labels based on loss distribution to enrich the dataset. the method in [12] proposes a semi-supervised few-shot segmentation method in remote sensing cases, which generates pseudo labels on super-pixels of backgrounds for mining latent features to enhance the network’s generalization capacity. The method in [13] combines semi-supervised learning with few-shot object detection and proposes the APLDet network, which utilizes a teacher model adaptively generating pseudo labels to guide the training of a student model.

In this study, we combine semi-supervised learning (SSL) with few-shot

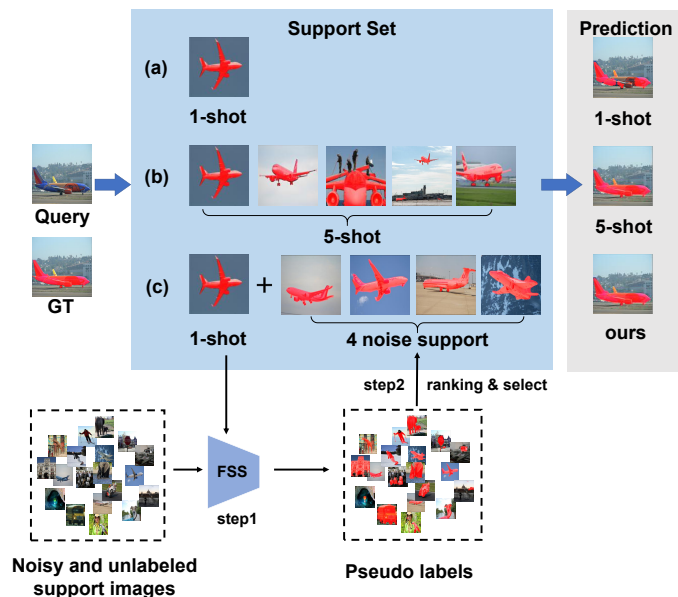


Figure 1: (a) 1-shot setting. (b) 5-shot setting. (c) 1-shot with additional 4 noise support images with pseudo labels. There is a large performance gap between 1-shot and 5-shot. Using 1-shot and 4 noise support can achieve comparable performance to 5-shot without increasing annotation cost.

30 segmentation (FSS) and propose a novel semi-supervised few-shot segmentation
 31 framework named **F4S**. Different from existing method [12] that intro-
 32 duces super-pixels to generate pseudo labels and only enhances the training
 33 phase of FSS, the proposed F4S framework generates pseudo labels of unlabeled
 34 images directly, and quantitatively evaluates the quality of pseudo
 35 labels based on a novel ranking algorithm, and finally enhance both the training
 36 and test phases of any off-the-shelf FSS models. A brief pipeline of F4S
 37 is shown in Fig. 1 (c), which consists of three steps. Firstly, pseudo labels
 38 are generated using a pre-trained FSS model for noisy and unlabeled support
 39 images. Secondly, pseudo labels with high confidence scores are selected as
 40 ground truth to augment the support set. Thirdly, the augmented support
 41 set is utilized to enhance the FSS model in both the training and test phases.

42 However, unlabeled support images are not a free lunch, as there are two
 43 problems that complicate pseudo-label selection (as shown in Fig. 2). 1)
 44 **Noisy Intra-Class Samples:** The noisy intra-class samples contain am-
 45 biguous objects that may strengthen the background and weaken the fore-
 46 ground, e.g., noisy “background” dominates the image as shown in Fig. 2 (a).

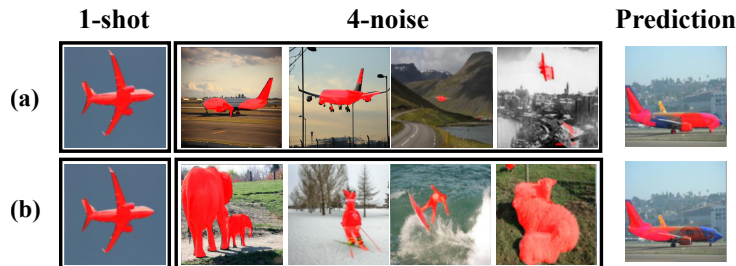


Figure 2: Examples of two basic problems. (a) Noisy intra-class samples as support samples. (b) Noisy inter-class samples as support samples.

47 **2) Noisy Inter-Class Samples:** The noisy inter-class samples introduce ir-
 48 relevant features to the task, which may cause feature bias and thus confuse
 49 the FSS model, e.g., the FSS model is confused by “elephant”, “person” and
 50 “sheep” when segmenting “aeroplane” as shown in Fig. 2 (b). We need to
 51 eliminate the two types of samples.

52 To solve the two basic problems, we propose a ranking algorithm in F4S
 53 to automatically eliminate the noisy intra-class samples and inter-class sam-
 54 ples. This ranking algorithm consists of two terms: an intra-class confidence
 55 term R and an inter-class confidence term T . The term R aims to iden-
 56 tify the noisy intra-class samples by calculating three terms: E_{sc} , E_{imc} and
 57 E_{cyc} . Specifically, E_{sc} measures prediction uncertainty based on binary en-
 58 tropy, E_{imc} identifies different types of errors based on the co-teaching frame-
 59 work [14, 15], and E_{cyc} measures object features completeness based on the
 60 cycle-consistency strategy [16, 17]. Besides, the term T aims to identify the
 61 noisy inter-class samples. It calculates the feature similarities between the
 62 support prototypes and the pseudo labels of noisy images. Finally, a ranking
 63 score E is calculated by weighting R and T , and the top-scored pseudo labels
 64 are treated as new support samples.

65 In order to theoretically explain the effectiveness of the ranking algo-
 66 rithm, we design a Structural Causal Model (SCM), which models the rel-
 67 evance of input support samples, noisy support set, and query labels. The
 68 SCM proves that the proposed ranking algorithm can successfully remove the
 69 confounding bias in the noisy support set (*cf.* Sect. 5). We also evaluate the
 70 proposed F4S framework on two popular FSS benchmarks: PASCAL-5ⁱ [1],
 71 and COCO-20ⁱ [18] in Sect. 6. Extensive quantitative and qualitative studies
 72 show that the F4S achieves new SOTA performance compared with existing
 73 fully supervised FSS methods.

74 This paper represents a very substantial extension of our previous confer-
75 ence paper [19]. The main improvements compared with [19] lie in threefold:
76 (i) We have improved the F4S framework by integrating a new term, E_{cyc} ,
77 derived from the cycle-consistency strategy, into the proposed ranking al-
78 gorithm. This enhancement notably boosts the model’s ability to identify
79 noisy samples without increasing its learnable parameters or memory cost,
80 achieving improved performances. (ii) We have added a justification sec-
81 tion (Sect. 5), where we theoretically explain why the proposed method can
82 work successfully based on a Structural Causal Model (SCM), which mod-
83 els the causal relevance of input data, generated pseudo labels, and output
84 predictions. (iii) We have conducted more comprehensive experiments to
85 evaluate the proposed method thoroughly. These experiments include exten-
86 sive evaluations on the PASCAL-5ⁱ dataset, along with additional compar-
87 isons with both inductive and transductive FSS methods, as well as recent
88 semi-supervised FSS methods. Furthermore, we have included visualization
89 results, conducted more comprehensive ablation studies, and performed ad-
90 ditional experimental analysis.

91 Our main contributions are as follows:

- 92 • We incorporate semi-supervised learning into the few-shot segmenta-
93 tion task and propose the **F4S** framework. It can benefit any off-
94 the-shelf few-shot segmentation models by solving the data scarcity
95 problem via introducing pseudo-labeled images, which has less been
96 studied.
- 97 • We design a ranking algorithm including an intra-class confidence score
98 R and an inter-class confidence score T to automatically identify and
99 eliminate the noisy samples in pseudo labels. The designing of R and
100 T are based on the underlying mechanism of FSS models. To the best
101 of our knowledge, this is the first work that quantitatively evaluates
102 the quality of pseudo labels in semi-supervised few-shot segmentation.
- 103 • We offer a theoretical explanation of the ranking algorithm grounded
104 in a Structural Causal Model (SCM). This analysis proves that the
105 proposed method has the capability to mitigate confounding bias within
106 the noisy support set through causal intervention.

107 2. Related Work

108 2.1. Few-shot Segmentation

109 Few-shot segmentation performs semantic segmentation in the few-shot
110 scenario, where only a few support images are given for a new class. Two
111 types of FSS methods, i.e., the prototype-based approaches [20, 2, 3, 4, 21,
112 5, 6] and the metric-based approaches [7, 8, 9], are mainly used to achieve
113 accurate segmentation.

114 The prototype-based approaches try to generate prototypes that describe
115 the class well from the limited training samples. For example, the method in
116 [20] generates foreground and background prototypes via a classifier trained
117 by support images with image-level labels. The method in [2] uses a proto-
118 type alignment strategy to make the prototypes more consistent. Seeing the
119 fact that one single prototype is hard to fully describe the class, some meth-
120 ods [3, 4] try to generate multiple prototypes for each class. For example, the
121 methods in [3] and [4] decompose the single class representation into a set of
122 part-aware prototypes that can describe diverse fine-grained object features
123 more precisely. The methods in [21] and [5] propose a parameter-free based
124 prototype generation method via feature clustering.

125 The metric-based approaches try to learn a class-agnostic similarity met-
126 ric that measures the similarity of region pairs, by which the query region
127 similar to the support region can be obtained. For example, the method in
128 [7] proposes a dense comparison module to calculate the similarity between
129 support features and query features under multiple levels. The method in [8]
130 proposes a multi-scale decoder with attention prior masks to achieve better
131 measurement. Besides, the methods in [22] and [23] provide a fresh insight
132 into the FSS task. The proposed BAM network incorporates an auxiliary
133 base learner into the conventional FSS meta learner to identify and remove
134 the feature-biased problem caused by base-class objects, and thus learn a bet-
135 ter class-agnostic metric function. Moreover, the method in [24] introduces
136 a divide-and-conquer strategy in FSS, which divides coarse results into small
137 regions and conquers the segmentation failures by leveraging the information
138 derived from support image-mask pairs.

139 Different from these existing methods, we generalize few-shot segmenta-
140 tion with more noisy and unlabeled images in both the training and testing
141 phases. Furthermore, we propose a new quality ranking algorithm that can
142 select good support samples from noisy samples accurately.

143 *2.2. Semi-Supervised Learning*

144 Semi-supervised learning [25, 26, 27, 28, 29, 30, 31] trains neural networks
145 on partially labeled datasets, including both labeled and unlabeled data. The
146 labeled data provides discriminative information about classes, while the un-
147 labeled data provides the underlying structure of the input data. Recent
148 works based on semi-supervised learning not only improve the performance
149 of deep neural networks, but also significantly reduce the cost associated
150 with data labeling. For example, the method in [25] generates and selects
151 pseudo labels for unlabeled data that exhibit high confidence above a spe-
152 cific threshold to enhance image classification. The method in [29] utilizes
153 the teacher-student framework, where the teacher model learns to generate
154 good pseudo labels from unlabeled data to benefit the student model for ob-
155 ject detection. The method in [28] proposes a new confidence score based on
156 the loss distribution to select good pseudo labels and benefit few-shot clas-
157 sification. The method in [27] generates and retains pseudo-labeled samples
158 with high confidence of the target domain for adversarial learning to solve
159 the domain adaptation problem. The method in [30] proposes a transfer
160 network, which is trained by pseudo labels and learns to exploit beneficial
161 feature representation knowledge in the extractor to enhance the training of
162 semantic segmentation network. In this paper, we propose a semi-supervised
163 FSS framework to expand the support image set with unlabeled images and
164 their pseudo labels.

165 *2.3. Few-shot Learning with Noisy Samples*

166 Few-shot learning with noisy samples [32, 33, 34, 35, 36] represents a more
167 realistic scenario, where support sets are susceptible to mislabeled samples.
168 Robustness to noisy samples is crucial for practical few-shot learning meth-
169 ods. Some existing works [32, 33] focus on feature similarity to identify
170 and eliminate the noisy samples. For instance, the method proposed in [32]
171 employs soft k-means clustering to detect noise within the support samples,
172 given that the features of noisy samples deviate significantly from the current
173 support set. The method described in [33] utilizes a feature-level similarity
174 assessment to reveal the heterogeneity and homogeneity within support sam-
175 ples.

176 Additionally, designing attention mechanisms is widely utilized for sup-
177 pressing noise. For example, the method in [34] introduces a semantically-
178 conditioned attention mechanism to estimate the importance of training in-
179 stances and bolster the model’s resilience to noise. Similarly, the method

180 outlined in [35] introduces an attention mechanism based on a novel trans-
181 former architecture, to effectively weigh mislabeled samples against cor-
182 rect ones. Moreover, the method described in [36] presents an attention-
183 based contrastive learning model incorporating discrete cosine transform in-
184 put. This model utilizes transformed frequency domain representations ob-
185 tained through discrete cosine transform as input, effectively removing high-
186 frequency components to suppress input noise.

187 Furthermore, recent research effort [37] extends the handling of noisy
188 samples to the few-shot segmentation task. It proposes a noise suppression
189 module to eliminate noisy activations by analyzing the correlation distribu-
190 tion between query and support features. However, [37] only considers the
191 inter-class noisy samples and cannot be generalized to a semi-supervised sce-
192 nario, where both intra-class and inter-class noisy samples abound. There-
193 fore, semi-supervised few-shot segmentation with noisy samples is a more
194 crucial scenario and remains largely unexplored. In this study, we introduce
195 a novel quality ranking algorithm designed to select high-quality support
196 samples from noisy pseudo-labeled data. This approach enhances few-shot
197 segmentation models in a semi-supervised way during both the training and
198 testing phases.

199 2.4. Causal Inference

200 Causal inference [38, 39] aims to formulate tasks in the view of causal-
201 ities and makes the network benefit from causal effects by removing the
202 *confounder*. Recently, a growing number of methods combing with causal in-
203 ference are proposed [40, 41, 42, 43, 44] in computer vision. For example, the
204 method in [40] uses causal inference to solve the semi-supervised semantic
205 segmentation, where the co-occurrence context is considered as a *confounder*
206 making the model hard to distinguish the category boundaries. A context
207 adjustment method with causal intervention is proposed to remove the con-
208 founding bias. The method in [41] treats the pre-trained knowledge as a
209 *confounder* in few-shot learning, and uses causal intervention to remove the
210 negative effect of the pre-trained knowledge. The method in [42] tackles the
211 out-of-distribution (OOD) generalization problem with causality. A causal
212 invariant transformation is proposed to keep the causal features from non-
213 causal features. Similarly, the method in [43] designs a meta-causal learner
214 to capture common causal features from multiple tasks and realize out-of-
215 distribution generalization. In this paper, we propose a structural causal
216 model in Sect. 5.1 to analyze the causalities among support samples, noisy

217 support set, and query labels in our F4S framework, and aim at improving
218 the FSS performance.

219 3. Formulation

220 We mathematically formulate the conventional few-shot segmentation
221 methods and the proposed F4S for better understanding.

222 **Conventional few-shot segmentation methods:** ① In the training
223 phase, a support set S^{base} including images I_S^{base} and pixel-level annotations
224 M_S^{base} of base classes is given. A few-shot segmentation network N_θ param-
225 eterized by θ need to be trained on $\{I_S^{base}, M_S^{base}\}$ to segment objects from a
226 query set Q^{base} within the meta-learning paradigm. The ground truth M_Q^{base}
227 of Q^{base} is given for loss calculation and backward propagation. ② In the test
228 phase, $\{I_S^{novel}, M_S^{novel}\}$ of novel classes is given, which provides support fea-
229 tures to help network N_θ predict segmentation masks M_Q^{novel} of novel objects
230 from Q^{novel} . Then, an evaluation metric, e.g. mIoU, is adopted to evaluate
231 the performance of N_θ , i.e. $mIoU(\hat{M}_Q^{novel}, M_Q^{novel})$.

232 **The proposed method F4S:** ① Before training, $\{I_S^{base}, M_S^{base}\}$ and a
233 set of noisy unlabeled images $I_{unlabel}$ are given. Pseudo labels P of $I_{unlabel}$
234 are generated by the pretrained network N_θ based on the support features of
235 $\{I_S^{base}, M_S^{base}\}$. ② A ranking algorithm is proposed here to obtain $\{I_{unlabel}^{base}, P^{base}\}$,
236 where the noisy pseudo-labeled samples are eliminated and superior pseudo-
237 labeled samples of base classes are retained. ③ In the training phase, based
238 on $\{I_S^{base}, M_S^{base}, I_{unlabel}^{base}, P^{base}\}$, the network N_θ is retrained within the meta-
239 learning paradigm. ④ Before test, we implement ① and ② again based
240 on $\{I_S^{novel}, M_S^{novel}\}$ to obtain $\{I_{unlabel}^{novel}, P^{novel}\}$ of novel classes. ⑤ In the
241 test phase, based on $\{I_S^{novel}, M_S^{novel}, I_{unlabel}^{novel}, P^{novel}\}$, the network N_θ outputs
242 the predictions \hat{M}_Q^{novel} of the query set Q^{novel} . Then, an evaluation metric
243 $mIoU(\hat{M}_Q^{novel}, M_Q^{novel})$ is utilized to evaluate the performance.

244 4. Method

245 4.1. Overview

246 Fig. 3 (a) shows the proposed F4S framework, which consists of three
247 phases. In phase I, a pretrained FSS network N_θ is used to obtain the
248 pseudo labels of the noisy and unlabeled support images. Various existing
249 FSS models can be employed here.

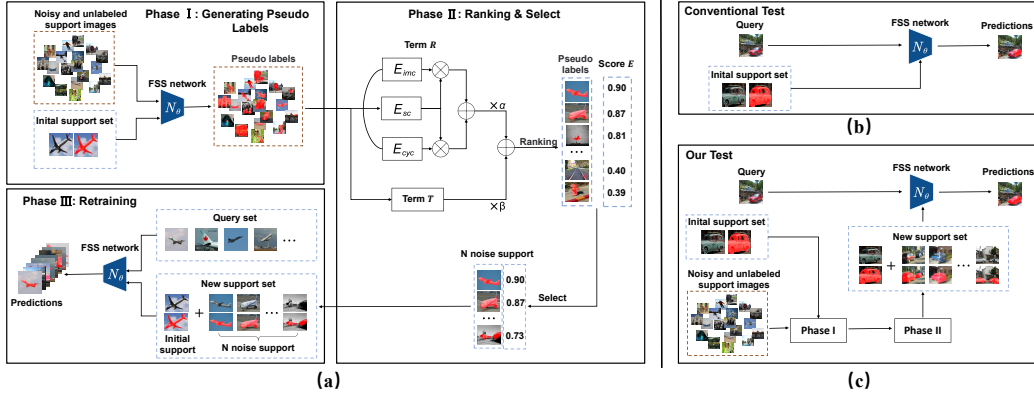


Figure 3: (a) The pipeline of the proposed F4S framework, which consists of three phases. In phase I, a pretrained FSS network N_θ is used to obtain the pseudo labels. Then, in phase II, a ranking algorithm is utilized to calculate quality scores E of pseudo labels and rank them. Finally, in phase III, top-scored pseudo labels are selected as new support samples to retrain N_θ . (b) The pipeline of the conventional FSS test. After retraining N_θ , it is tested on novel classes, e.g., “car”, with an annotated initial support set. (c) The pipeline of our FSS test based on the proposed semi-supervised framework. N_θ is tested on novel classes with a new support set, which is expanded following phase I and phase II.

250 In phase II, the ranking algorithm is utilized to evaluate the pseudo labels.
 251 Specifically, an intra-class confidence term R and an inter-class confidence
 252 term T are calculated for each pseudo label. Then, a final ranking score E
 253 is obtained by simply calculating the weighted sum of R and T :

$$E = \alpha \cdot R + \beta \cdot T \quad (1)$$

254 where α and β are weighting coefficients. Afterwards, the top k scored pseudo
 255 labels are selected to form a new annotation set:

$$\mathcal{S}_{new}^{base} \leftarrow \mathcal{S}^{base} + \{(X_1, \hat{Y}_{X_1}), (X_2, \hat{Y}_{X_2}), \dots, (X_k, \hat{Y}_{X_k})\} \quad (2)$$

256 where \mathcal{S}^{base} indicates the initial annotation set of base classes in the training
 257 phase, \hat{Y}_X indicates the pseudo label of image X .

258 Finally, in phase III, the new annotation set \mathcal{S}_{new}^{base} is used to retrain N_θ and
 259 get better predictions. More details of the intra-class confidence term R and
 260 the inter-class confidence term T are introduced in Sect. 4.2 and Sect. 4.3,
 261 respectively. Besides, in order to enhance the inference of FSS models, we
 262 further propose a new test process based on F4S in Sect. 4.4.

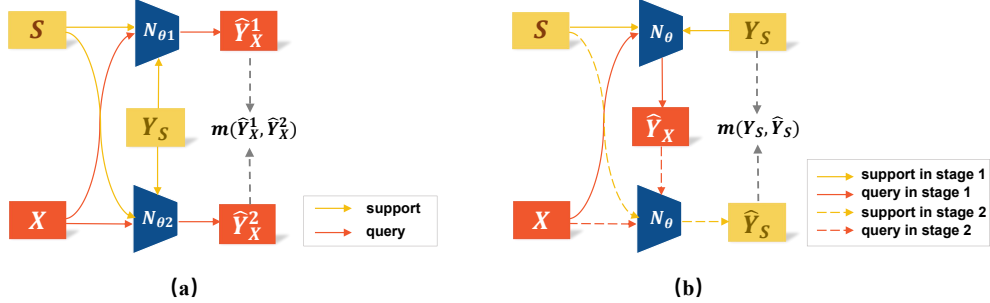


Figure 4: (a) The pipeline of E_{imc} . The unlabeled image X is processed by two FSS models $N_{\theta 1}$, $N_{\theta 2}$, with a given support sample $\{S, Y_S\}$. Then, a metric $m(\cdot, \cdot)$ is calculated between the two output \hat{Y}_X^1 , \hat{Y}_X^2 . (b) The pipeline of E_{cyc} , which consists of two stages. In stage 1, a FSS model N_θ makes prediction \hat{Y}_X of the unlabeled image X based on a given support sample $\{S, Y_S\}$. In stage 2, N_θ makes prediction \hat{Y}_S of S based on $\{X, \hat{Y}_X\}$. Finally, a metric $m(\cdot, \cdot)$ is calculated between Y_S and \hat{Y}_S .

263 4.2. Intra-Class Confidence Term R

264 The term R aims to identify the noisy intra-class samples. The calculation
 265 of R is shown in Eq. 3:

$$R = E_{sc} \times (E_{imc} + E_{cyc}) \quad (3)$$

266 where the segmentation confidence term E_{sc} estimates the prediction uncer-
 267 tainty of pseudo labels, the instance mask consensus term E_{imc} identifies
 268 different types of errors in pseudo labels, and the cyclic mask consensus term
 269 E_{cyc} identifies pseudo labels with incomplete object features. Now, we intro-
 270 duce the three terms E_{sc} , E_{imc} , and E_{cyc} in detail.

271 **Segmentation Confidence Term E_{sc} .** This term is calculated by
 272 adopting a binary-entropy-based function to measure the prediction uncer-
 273 tainty:

$$E_{sc} = -\frac{1}{N} \sum_i H(i) + B \quad (4)$$

274 where i indicates a pixel position, $H(\cdot)$ is the binary entropy function, N is
 275 the total number of pixels, and B is a bias term to ensure $E_{sc} \in [0, 1]$. The
 276 formulation of $H(x)$ is shown in Eq. 5, where $p(i)$ is the logit at pixel position
 277 i .

$$H(x) = -p(i)\log(p(i)) - (1 - p(i))\log(1 - p(i)) \quad (5)$$

278 **Instance Mask Consensus Term E_{imc} .** This term is motivated by the
 279 co-teaching theory [14, 15], which proves that two diverged networks can filter

280 different types of errors. Therefore, if two diverged few-shot segmentation
 281 networks output similar predictions to the same wild image, the predictions
 282 contain less error and have high confidence. The pipeline of getting E_{imc} is
 283 shown in Fig. 4 (a) and its calculation is:

$$E_{imc} = m(\hat{Y}_X^1, \hat{Y}_X^2) \quad (6)$$

284 where \hat{Y}_X^1 and \hat{Y}_X^2 are predictions of the same unlabeled image X from two
 285 diverged networks N_{θ_1} and N_{θ_2} . $m(\cdot, \cdot)$ indicates a segmentation metric score,
 286 e.g., mIoU.

287 **Cyclic Mask Consensus Term E_{cyc} .** Inspired by the cycle-consistency
 288 strategy of [16], we design a cyclic pipeline in FSS to estimate the segmenta-
 289 tion confidence. The detailed pipeline is shown in Fig. 4 (b). Specifically, it
 290 consists of two stages: in stage 1, a FSS model N_θ makes a prediction \hat{Y}_X of
 291 the unlabeled image X based on the annotated support sample $\{S, Y_S\}$; in
 292 stage 2, based on $\{X, \hat{Y}_X\}$, N_θ makes a prediction \hat{Y}_S of the support image
 293 S . Finally, the E_{cyc} can be calculated by:

$$E_{cyc} = m(Y_S, \hat{Y}_S) \quad (7)$$

294 4.3. Inter-Class Confidence Term T

295 The term T aims to identify the noisy inter-class samples based on the
 296 feature similarities between the support prototypes and the pseudo labels.
 297 First, the prototype of class c of the initial support set $\mathcal{S}^c = \{S_1^c, S_2^c, \dots, S_n^c\}$
 298 are calculated by:

$$\mathcal{P}^c = \frac{1}{n} \sum_{i=1}^n \sigma(\mathcal{F}_{S_i^c}, Y_{S_i^c}) \quad (8)$$

299 where $\mathcal{F}_{S_i^c} \in \mathbb{R}^{C \times H \times W}$ is the feature map of support S_i^c of class c , $Y_{S_i^c}$ is the
 300 manual annotation, $\sigma(\cdot)$ is the masked global average pooling, and $\mathcal{P}^c \in \mathbb{R}^C$
 301 is the prototype of class c . Then, the term T can be calculated by:

$$T = s(\mathcal{P}^c, \sigma(\mathcal{F}_X, \hat{Y}_X)) \quad (9)$$

302 where $\mathcal{F}_X \in \mathbb{R}^{C \times H \times X}$ is the feature map of X , \hat{Y}_X is the generated pseudo
 303 label, $s(\cdot, \cdot)$ is a similarity metric, e.g., cosine similarity.

304 4.4. A New Test Process based on F4S

305 To enhance the inference of FSS models, we can further expand the initial
 306 support set of novel classes via F4S in the test phase, of which the pipeline
 307 is shown in Fig. 3 (c). Specifically, different from the conventional FSS test
 308 (Fig. 3 (b)), where only a small annotated support set \mathcal{S}^{novel} of novel classes
 309 is utilized, our test enriches \mathcal{S}^{novel} following the pipeline of phase I and phase
 310 II of the proposed F4S to obtain a new support set $\mathcal{S}_{new}^{novel}$:

$$\mathcal{S}_{new}^{novel} \leftarrow \mathcal{S}^{novel} + \{(X_1, \hat{Y}_{X_1}), (X_2, \hat{Y}_{X_2}), \dots, (X_k, \hat{Y}_{X_k})\} \quad (10)$$

311 Then, the query images will be segmented with the new support set $\mathcal{S}_{new}^{novel}$
 312 to get better predictions.

313 **5. Justification**

314 5.1. Structural Causal Model

315 We construct a causal graph to formulate the causalities among the se-
 316 lected support sample, query prediction, and the noisy support set, which is
 317 shown in Fig. 5 (a). The causal graph consists of four nodes: X indicates the
 318 selected support sample; Y is the query label; D indicates the noisy support
 319 set, which includes the noisy intra-class and inter-class samples and acts as
 320 the *confounder* in the causal graph; M is the transformed representation of
 321 X in the low-dimensional manifold embedded in the latent high-dimension
 322 space via FSS model [40]. The directed path between two nodes indicates
 323 the causalities : cause \rightarrow effect. Next, we detail the rationale of Fig. 5 (a).

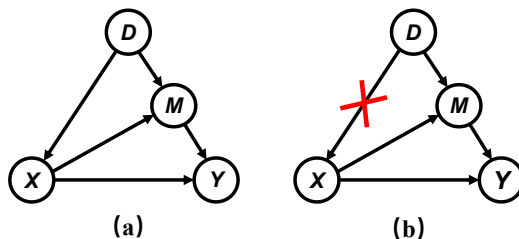


Figure 5: (a) The causal graph for FSS. The *confounder* D degrades FSS via $X \leftarrow D \rightarrow M \rightarrow Y$, i.e., noisy intra-class and inter-class samples in D are mistakenly selected as support samples X causing serious feature bias and bad query predictions of Y . (b) The revised causal graph of our F4S, where the proposed ranking algorithm in F4S can cut off the path towards X by $do(X)$, and thus ensures the selected support samples are noiseless.

324 $D \rightarrow X$. The support sample X is sampled from the noisy support set D .

325 $X \rightarrow Y$. The support sample X provides object cues to predict query
 326 label Y . However, this latent relevance between X and Y cannot be obtained
 327 directly, and therefore a FSS model $f(\cdot)$ is needed here to learn a transformed
 328 representation M between X and Y .

329 $D \rightarrow M$. The transformed representation M is a subset of that of D due
 330 to that the FSS model $f(\cdot)$ is trained on D .

331 $X \rightarrow M \rightarrow Y$. The support sample X leads to the transformed rep-
 332 resentation M via FSS model, i.e., $M = f(X)$, and M contributes to the
 333 prediction of Y , i.e., $P(Y|X, M)$. X with less noise leads to better M , and
 334 finally benefits the prediction of Y .

335 Based on the causal graph, one can see that the *confounder* D degrades
 336 $P(Y|X)$ via the backdoor path $X \leftarrow D \rightarrow M \rightarrow Y$. Removing the backdoor
 337 path is the key challenge for improving F4S performance. Next, we show
 338 how to remove the confounding effect by causal intervention $P(Y|do(X))$.

339 5.2. Causal Intervention via Backdoor Adjustment

340 In this section, we propose to use the causal intervention $P(Y|do(X))$,
 341 which can remove the confounding effect by $do(\cdot)$ to get a better prediction
 342 of label Y . The key idea is to cut off the path $D \rightarrow X$ (Fig. 5(b)) via
 343 backdoor adjustment [38], i.e., identifying and eliminating noisy intra-class
 344 and inter-class samples when sampling X from D . Following [45, 38], we
 345 have:

$$\begin{aligned}
 P(Y|do(X)) &= \sum_{D=\{d_0, d_1\}} P(Y|X, M = f(X, D))P(D) \\
 &= P(Y|X, f(X, D = d_0))P(D = d_0) \\
 &\quad + P(Y|X, f(X, D = d_1))P(D = d_1) \\
 &= P(Y|X, f(X, D = d_0)) \cdot \alpha \\
 &\quad + P(Y|X, f(X, D = d_1)) \cdot \beta
 \end{aligned} \tag{11}$$

346 where the noisy support set D includes two types of noisy samples: d_0 in-
 347 dicates the noisy intra-class samples, and d_1 indicates the noisy inter-class
 348 samples. $P(D = d_0)$ and $P(D = d_1)$ indicate the ratio of d_0 and d_1 in D .
 349 For simplicity, they are set as two constants: α and β , respectively. Next,
 350 we estimate $P(Y|X, f(X, D = d_0))$ and $P(Y|X, f(X, D = d_1))$.

351 5.2.1. Estimation of $P(Y|X, f(X, D = d_0))$

352 Following [46], we implement the sampling process from the intervened
 353 distribution to get $P(Y = y|X = x, f(X = x, D = d_0))$, abbreviated as

354 $P(y|x, f(x, d_0))$. It represents the probability of predicting the label $Y = y$
 355 under the condition of input $X = x$ with intra-class noise $D = d_0$. Intuitively,
 356 less intra-class noise d_0 leads to a higher probability P to predict the correct
 357 label $Y = y$, which can be reflected by a segmentation metric score. To this
 358 end, we can get:

$$P(y|x, f(x, d_0)) \propto m(y, \hat{y}) \quad (12)$$

359 where \hat{y} is the prediction of label y , $m(\cdot, \cdot)$ indicates a segmentation metric
 360 score, e.g., mIoU.

361 However, the label $Y = y$ is unavailable since the noisy support set is
 362 not annotated, and thus $m(y, \hat{y})$ can not be calculated. Fortunately, the
 363 proposed intra-class confidence score R (Eq. 3) can estimate the credibility
 364 of prediction \hat{y} in a blind way, i.e., without annotated label y . Therefore, we
 365 can further obtain:

$$P(y|x, f(x, d_0)) \propto m(y, \hat{y}) \propto R \quad (13)$$

366 In this way, the proposed intra-class confidence term R can estimate the
 367 target $P(Y|X, f(X, D = d_0))$ due to its correlation of metric score $m(\cdot, \cdot)$.

368 5.2.2. **Estimation of $P(Y|X, f(X, D = d_1))$**

369 Implementing the sampling process from the intervened distribution, we
 370 can get the term $P(y|x, f(x, d_1))$, which represents the probability of pre-
 371 dicting the label $Y = y$ based on input $X = x$ with inter-class noise $D = d_1$.
 372 Intuitively, less inter-class noise d_1 leads to higher probability P to predict
 373 label $Y = y$, which can be reflected by the similarity between class prototype
 374 \mathcal{P} and input noisy support sample x . Therefore, we have:

$$P(y|x, f(x, d_1)) \propto s(\mathcal{P}, f(x_s)) \quad (14)$$

375 where \mathcal{P} is the class-specific prototype, $f(x_s)$ is the feature map of the in-
 376 put support sample x , $s(\cdot, \cdot)$ is a similarity metric, e.g., cosine similarity.
 377 Combining Eq. 14 with Eq. 9, we get:

$$P(y|x, f(x, d_1)) \propto T \quad (15)$$

378 In this way, the proposed inter-class confidence term T can estimate the
 379 target $P(Y|X, f(X, D = d_1))$ based on the feature similarities.

380 Finally, combining Eq. 13 with Eq. 15, we can rewrite Eq. 11:

$$P(Y|do(X)) \propto R \cdot \alpha + T \cdot \beta = E \quad (16)$$

381 Therefore, the proposed ranking mechanism can successfully remove the con-
 382 founding effect in the noisy support set D following the causal intervention
 383 $P(Y|do(X))$.

384 6. Experiment

385 6.1. Setup

386 **Datasets.** We evaluate our method on PASCAL-5ⁱ [1] and COCO-20ⁱ
 387 [18] datasets and use the unlabeled 123,403 images in COCO2017 [47] for
 388 conducting experiments. Specifically, following the setup in [1], 20 categories
 389 in the PASCAL VOC 2012 dataset [48] are partitioned into 4 folds (i.e., fold-
 390 0, fold-1, fold-2, and fold-3) and each fold contains 5 categories. Following the
 391 setups in [18], 80 categories in the COCO dataset [47] are also divided into 4
 392 folds and each fold contains 20 categories. The experiments are conducted in
 393 a cross-validation manner and the validation episode is set to 1000 for each
 394 fold.

395 **Evaluation metrics.** Following previous works [3, 4, 21, 49], we adopt
 396 mean intersection over union (mIoU) and foreground-background IoU (FB-
 397 IoU) as our evaluation metrics. The mIoU metric is computed by averaging
 398 IoU of all classes: $mIoU = \frac{1}{n} \sum_{i=1}^n IoU_i$. The FB-IoU metric is calculated
 399 by averaging IoU of foreground and background: $mIoU = \frac{1}{2}(IoU_F + IoU_B)$.

400 **Implementation details.** All of our experiments are conducted on two
 401 NVIDIA Titan XP GPUs and Intel Core i9-9900k CPU @ 3.60GHz \times 16.
 402 Our code is constructed on PyTorch. We build our F4S framework based on
 403 the open-sourced code of methods in [8, 50]. In Sect. 4.2, multiple backbones
 404 are adopted as the two diverged networks $N_{\theta_1}, N_{\theta_2}$. The detailed settings of
 405 $N_{\theta_1}, N_{\theta_2}$ are shown in Table 1. The publicly released pretrained models in
 406 methods [8, 50] are used directly. For the PFENet (VGG16) on PASCAL-5ⁱ
 407 and PFENet (ResNet101) on COCO-20ⁱ, we train the models following the
 408 official settings in [8]. We set $m(\cdot, \cdot)$ to mIoU score in Sect. 4.2 and set $s(\cdot, \cdot)$
 409 to cosine similarity in Sect. 4.3. The feature maps $\mathcal{F} \in \mathbb{R}^{C \times H \times W}$ in Sect. 4.3
 410 are extracted from the last convolutional layer of the backbone. α and β
 411 in Eq. 1 are set to 0.3 and 0.7, respectively. In the training phase, pseudo
 412 labels with $E \geq 0.65$ are selected as new annotations of base classes. In the
 413 test phase, top 4 scored pseudo labels are introduced into the support set
 414 of novel classes. In phase III, the retraining setting strictly follows the base
 415 model [8, 50].

Table 1: The diverged networks in E_{imc} .

Method	N_{θ_1}	N_{θ_2}
HSNet [50]	ResNet50	ResNet101
	VGG16	ResNet101
PFENet [8]	VGG16	ResNet50
	VGG16	ResNet101

416 *6.2. Quantitative Results*

417 We evaluate the proposed F4S on PASCAL-5ⁱ [1] and COCO-20ⁱ datasets
 418 and compare the metric scores with recent FSS methods [2, 8, 50, 51, 53, 54].
 419 Table 2 shows the mIoU and FB-IoU values of our method and the existing
 420 methods under 1-shot settings on PASCAL-5ⁱ and COCO-20ⁱ datasets, where
 421 “F4S (HSNet)” indicates that F4S is implemented on the HSNet [50]. Here,
 422 the F4S is set to 1-shot/5-shot with 4 noise support (as shown in Fig. 1 (c))
 423 and our evaluation has two test ways: the conventional test in Fig. 3 (b)
 424 and our test in Fig. 3 (c), which are annotated as “†” and “‡” in Table 2,
 425 respectively.

426 Compared with the baseline (HSNet), we can observe that on the PASCAL-
 427 5ⁱ dataset, “F4S (HSNet) †” achieves mIoU improvements of 1.6%, 0.8%, and
 428 0.3% on three backbones under 1-shot, and achieves mIoU improvements of
 429 0.7%, 0.6%, and 0.5% under 5-shot. Meanwhile, on the COCO-20ⁱ dataset,
 430 “F4S (HSNet) †” also achieves further improvements of mIoU and FB-IoU on
 431 different backbones under 1-shot and 5-shot. These results demonstrate that
 432 the proposed F4S can benefit FSS models from the unlabeled support images
 433 in the retraining phase (Fig. 3 (a)) without noise disturbance. Besides, fol-
 434 lowing our test (Fig. 3 (c)), “F4S (HSNet) ‡” achieves mIoU improvements of
 435 8.2%, 6.8%, and 6.1% on three backbones on PASCAL-5ⁱ, and mIoU improve-
 436 ments of 10.8%, and 10.2% on two backbones on COCO-20ⁱ under 1-shot.
 437 Moreover, there are also remarkable performance improvements achieved by
 438 “F4S (HSNet) ‡” under 5-shot. These quantitative results verify that ex-
 439 tending the support set with unlabeled support images via F4S can directly
 440 benefit the inference of FSS models in the test phase.

441 We also compare the proposed method with recent transductive and in-
 442 ductive methods. In Table 2, one can observe that the proposed method
 443 “F4S (HSNet) ‡” with different backbones obtains new state-of-the-art per-
 444 formances. On PASCAL-5ⁱ and with ResNet101 backbone, our 1-shot and
 445 5-shot results of “F4S (HSNet)‡” respectively achieve 3.7% and 0.9% of mIoU

Table 2: Performance of the proposed F4S on PASCAL-5ⁱ and COCO-20ⁱ datasets. “†” is the results of the conventional test. “‡” is the results of our test based on the F4S. “Oracle” is the 5-shot performance. “±0.1” is the standard deviation of repeating 5 times.

Dataset	Backbone	Method	Type	1-shot		5-shot	
				mIoU	FB-IoU	mIoU	FB-IoU
PASCAL-5 ⁱ	VGG16	PFENet [8]	inductive	58.0	72.0	59.0	72.3
		HSNet [50]	inductive	59.7	73.4	64.1	76.6
		HPA [51]	inductive	61.5	75.2	66.2	79.3
		DCP [24]	inductive	62.6	75.6	67.8	80.7
		BAM [22]	inductive	64.4	77.3	68.8	81.1
		BAM* [23]	inductive	65.3	77.5	69.6	81.3
		F4S (HSNet)†	inductive	61.3 (± 0.3)	74.4 (± 0.2)	64.8 (± 0.2)	76.9 (± 0.2)
		F4S (HSNet)‡	inductive	67.9 (± 0.2)	79.2 (± 0.1)	68.2 (± 0.3)	79.7 (± 0.3)
	ResNet50	RePRI [49]	transductive	59.1	-	66.8	-
		PFENet [8]	inductive	60.8	73.3	61.9	73.9
		HSNet [50]	inductive	64.0	76.7	69.5	80.6
		HPA [51]	inductive	64.8	76.4	68.9	81.1
		CDFS [52]	transductive	65.3	-	70.8	-
		DCP [24]	inductive	66.1	77.6	70.3	81.5
BAM [22]		inductive	67.8	79.7	70.9	82.2	
BAM* [23]		inductive	68.3	80.3	71.8	83.1	
ResNet101	F4S (HSNet)†	inductive	64.8 (± 0.2)	77.2 (± 0.2)	70.1 (± 0.2)	81.0 (± 0.2)	
	F4S (HSNet)‡	inductive	70.8 (± 0.2)	81.5 (± 0.1)	72.0 (± 0.3)	82.3 (± 0.2)	
	PFENet [8]	inductive	60.1	72.9	61.4	73.5	
	DCAMA [53]	inductive	64.6	77.6	68.3	80.8	
	HPA [51]	inductive	65.6	76.6	68.9	80.4	
	HSNet [50]	inductive	66.2	77.6	70.4	80.6	
	DCP [24]	inductive	67.3	78.5	71.5	82.7	
	BAM [22]	inductive	68.6	80.2	72.5	84.1	
	F4S (HSNet)†	inductive	66.5 (± 0.2)	78.2 (± 0.2)	70.9 (± 0.3)	81.1 (± 0.2)	
	F4S (HSNet)‡	inductive	72.3 (± 0.1)	82.3 (± 0.1)	73.4 (± 0.2)	82.6 (± 0.3)	
COCO-20 ⁱ	ResNet50	RePRI [49]	transductive	34.0	-	42.1	-
		HSNet [50]	inductive	39.2	68.2	46.9	70.7
		CDFS [52]	transductive	42.0	-	49.8	-
		DCAMA [53]	inductive	43.3	69.5	48.3	71.7
		HPA [51]	inductive	43.4	68.2	50.0	71.2
		DCP [24]	inductive	45.5	-	50.9	-
		BAM [22]	inductive	46.2	-	51.2	-
		BAM* [23]	inductive	46.9	72.3	51.9	74.7
	ResNet101	F4S (HSNet)†	inductive	40.9 (± 0.3)	69.1 (± 0.2)	49.0 (± 0.4)	71.9 (± 0.5)
		F4S (HSNet)‡	inductive	50.0 (± 0.4)	72.6 (± 0.5)	52.0 (± 0.3)	74.0 (± 0.3)
		PFENet [8]	inductive	38.5	63.0	42.7	65.8
		HSNet [50]	inductive	41.2	69.1	49.5	72.4
		DCAMA [53]	inductive	43.5	69.9	51.9	73.3
		HPA [51]	inductive	45.8	68.4	52.4	74.0
ResNet101	BAM* [23]	inductive	48.5	69.9	52.7	74.1	
	F4S (HSNet)†	inductive	42.8 (± 0.2)	69.8 (± 0.2)	51.2 (± 0.5)	73.3 (± 0.4)	
F4S (HSNet)‡	inductive	51.4 (± 0.2)	73.3 (± 0.3)	54.1 (± 0.4)	75.5 (± 0.4)		

* indicates the improved version of the base method.

Table 3: Performance of the proposed F4S without the retraining phase on PASCAL-5ⁱ and COCO-20ⁱ datasets. “Oracle” is the 5-shot performance. “ ± 0.1 ” is the standard deviation of repeating 5 times.

Dataset	Backbone	Method	Type	1-shot		5-shot	
				mIoU	FB-IoU	mIoU	FB-IoU
PASCAL-5 ⁱ	VGG16	PFENet [8]	inductive	58.0	72.0	59.0	72.3
		HSNet [50]	inductive	59.7	73.4	64.1	76.6
		HPA [51]	inductive	61.5	75.2	66.2	79.3
		DCP [24]	inductive	62.6	75.6	67.8	80.7
		BAM [22]	inductive	64.4	77.3	68.8	81.1
		BAM* [23]	inductive	65.3	77.5	69.6	81.3
		F4S (PFENet)‡	inductive	59.8 (± 0.2)	72.1 (± 0.2)	60.3 (± 0.3)	72.5 (± 0.3)
		F4S (HSNet)‡	inductive	66.5 (± 0.2)	78.4 (± 0.1)	67.1 (± 0.2)	78.9 (± 0.3)
	ResNet50	RePRI [49]	transductive	59.1	-	66.8	-
		PFENet [8]	inductive	60.8	73.3	61.9	73.9
		HSNet [50]	inductive	64.0	76.7	69.5	80.6
		HPA [51]	inductive	64.8	76.4	68.9	81.1
		CDFS [52]	transductive	65.3	-	70.8	-
		DCP [24]	inductive	66.1	77.6	70.3	81.5
BAM [22]		inductive	67.8	79.7	70.9	82.2	
BAM* [23]		inductive	68.3	80.3	71.8	83.1	
COCO-20 ⁱ	ResNet50	F4S (PFENet)‡	inductive	62.4 (± 0.2)	73.3 (± 0.2)	62.9 (± 0.3)	73.5 (± 0.2)
		F4S (HSNet)‡	inductive	70.6 (± 0.2)	81.4 (± 0.1)	71.7 (± 0.3)	82.0 (± 0.3)
		PFENet [8]	inductive	60.1	72.9	61.4	73.5
		DCAMA [53]	inductive	64.6	77.6	68.3	80.8
		HPA [51]	inductive	65.6	76.6	68.9	80.4
		HSNet [50]	inductive	66.2	77.6	70.4	80.6
		DCP [24]	inductive	67.3	78.5	71.5	82.7
		BAM [22]	inductive	68.6	80.2	72.5	84.1
	ResNet101	F4S (HSNet)‡	inductive	72.1 (± 0.1)	82.1 (± 0.1)	72.6 (± 0.3)	82.2 (± 0.3)
		RePRI [49]	transductive	34.0	-	42.1	-
		HSNet [50]	inductive	39.2	68.2	46.9	70.7
		CDFS [52]	transductive	42.0	-	49.8	-
		DCAMA [53]	inductive	43.3	69.5	48.3	71.7
		HPA [51]	inductive	43.4	68.2	50.0	71.2
DCP [24]		inductive	45.5	-	50.9	-	
BAM [22]		inductive	46.2	-	51.2	-	
ResNet101	BAM* [23]	inductive	46.9	72.3	51.9	74.7	
	F4S (HSNet)‡	inductive	49.7 (± 0.4)	72.2 (± 0.2)	51.0 (± 0.5)	72.9 (± 0.4)	
	PFENet [8]	inductive	38.5	63.0	42.7	65.8	
	HSNet [50]	inductive	41.2	69.1	49.5	72.4	
	DCAMA [53]	inductive	43.5	69.9	51.9	73.3	
	HPA [51]	inductive	45.8	68.4	52.4	74.0	
	BAM* [23]	inductive	48.5	69.9	52.7	74.1	
	F4S (PFENet)‡	inductive	41.5 (± 0.2)	63.8 (± 0.2)	43.3 (± 0.3)	66.4 (± 0.4)	
	F4S (HSNet)‡	inductive	51.1 (± 0.4)	73.1 (± 0.5)	52.4 (± 0.4)	74.5 (± 0.4)	

* indicates the improved version of the base method.

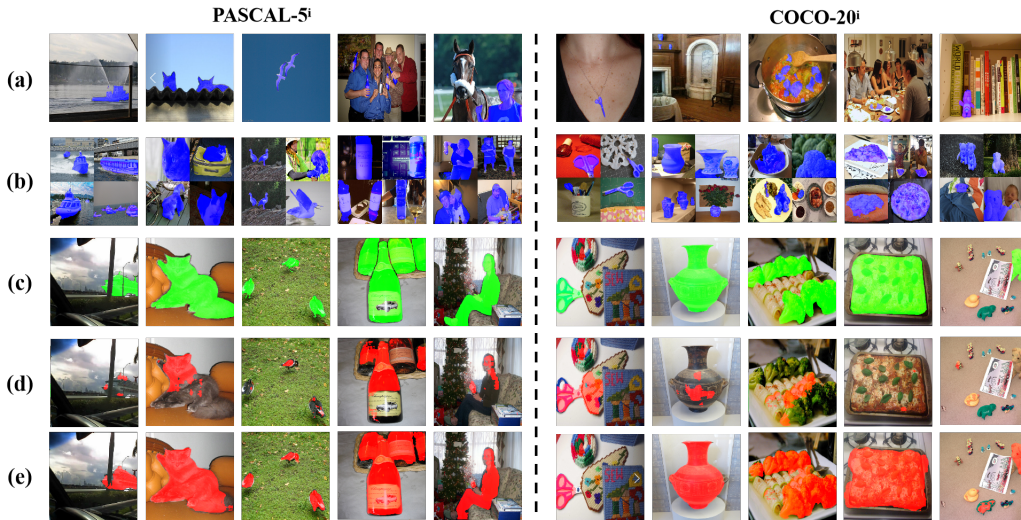


Figure 6: Qualitative results of the proposed F4S and its baseline. The left panel is from PASCAL-5ⁱ, and the right panel is from COCO-20ⁱ. From top to bottom: (a) 1-shot support images with ground truth, (b) 4 noise support images with pseudo labels via F4S, (c) query images with ground truth, (d) baseline predictions, (e) F4S predictions.

446 improvements over BAM [22]. On COCO-20ⁱ and with ResNet101 backbone,
 447 “F4S (HSNet)_‡” also outperforms recent methods with a sizable margin as well,
 448 achieving 2.9% and 1.4% of mIoU improvements over BAM* [23]. These
 449 results verify the superiority of the proposed method in the few-shot segmen-
 450 tation task.

451 Furthermore, we also evaluate F4S in the test phase directly without the
 452 retraining phase to save the training cost. Two popular FSS models, i.e.,
 453 HSNet [50] and PFENet [8], are adopted to implement F4S. The quanti-
 454 tative results are shown in Table 3. One can observe that on PASCAL-5ⁱ
 455 dataset and under the 1-shot setting, “F4S (PFENet)” achieves mIoU im-
 456 provements of 1.8%, and 1.6% on VGG16 and ResNet50 backbones compared
 457 with PFENet performance (baseline), and “F4S (HSNet)” achieves mIoU im-
 458 provements of 6.8%, 6.6%, and 5.9% on three different backbones compared
 459 with HSNet performance (baseline). On COCO-20ⁱ dataset, “F4S (HSNet)”
 460 and “F4S (PFENet)” also obtain superior performance compared with the
 461 baseline. These quantitative results prove that the proposed F4S can benefit
 462 the inference of FSS models directly without extra training.

463 It is worth noting that in both Table 2 and Table 3, the performance of

Table 4: Performance comparison with recent semi-supervised few-shot segmentation methods on PASCAL-5ⁱ and COCO-20ⁱ datasets.

Dataset	Backbone	Method	1-shot		5-shot	
			mIoU	FB-IoU	mIoU	FB-IoU
PASCAL-5 ⁱ	ResNet50	CLRS [12]	56.4	-	67.7	-
		UaFSS [55]	67.0	79.2	68.9	80.2
		F4S (HSNet)†	64.8 (± 0.2)	77.2 (± 0.2)	70.1 (± 0.2)	81.0 (± 0.2)
		F4S (HSNet)‡	70.8 (± 0.2)	81.5 (± 0.1)	72.0 (± 0.3)	82.3 (± 0.2)
	ResNet101	CLRS [12]	64.3	-	68.2	-
		UaFSS [55]	68.5	79.4	69.5	79.4
		F4S (HSNet)†	66.5 (± 0.2)	78.2 (± 0.2)	70.9 (± 0.3)	81.1 (± 0.2)
		F4S (HSNet)‡	72.3 (± 0.1)	82.3 (± 0.1)	73.4 (± 0.2)	82.6 (± 0.3)
COCO-20 ⁱ	ResNet50	CLRS [12]	33.0	-	36.3	-
		UaFSS [55]	41.3	68.9	46.4	70.9
		F4S (HSNet)†	40.9 (± 0.3)	69.1 (± 0.2)	49.0 (± 0.4)	71.9 (± 0.5)
		F4S (HSNet)‡	50.0 (± 0.4)	72.6 (± 0.5)	52.0 (± 0.3)	74.0 (± 0.3)
	ResNet101	UaFSS [55]	43.6	69.9	46.8	70.7
		F4S (HSNet)†	42.8 (± 0.2)	69.8 (± 0.2)	51.2 (± 0.5)	73.3 (± 0.4)
		F4S (HSNet)‡	51.4 (± 0.2)	73.3 (± 0.3)	54.1 (± 0.4)	75.5 (± 0.4)

464 F4S (1-shot with 4 noise support) surprisingly surpasses the 5-shot perfor-
465 mance of HSNet in some cases. This can be attributed to two aspects. First,
466 the training of models is enhanced due to the additional support features
467 from noisy and unlabeled support images introduced by F4S. Second, the
468 annotated support samples in “Oracle” are randomly sampled from datasets
469 and may include noisy intra-class samples, while the proposed F4S guarantees
470 the exclusion of such noisy intra-class samples.

471 Finally, we also compare the proposed method with recent semi-supervised
472 methods [12, 55] to show the superior performance in Table 4. One can
473 see that on PASCAL-5ⁱ dataset and with ResNet50 backbone, the proposed
474 “F4S (HSNet)†” achieves 3.8% of mIoU improvement in 1-shot setting and
475 3.1% of mIoU improvement in 5-shot setting over UaFSS[55]. Besides, with
476 ResNet101 backbone, the proposed method also outperforms recent methods
477 with a sizable margin as well, achieving 3.8% (1-shot) and 3.9% (5-shot) of
478 mIoU improvements over UaFSS[55]. Besides, on COCO-20ⁱ dataset and
479 with ResNet50 and ResNet101 backbones, the 1-shot and 5-shot results of
480 “F4S (HSNet)‡” are also superior to both UaFSS[55] and CLRS[12] with a
481 remarkable margin.

482 *6.3. Qualitative Results*

483 Fig. 6 shows the qualitative results of “F4S (HSNet)” with ResNet101
484 backbone on PASCAL-5ⁱ and COCO-20ⁱ datasets. As can be noticed, (e) F4S
485 predictions include more complete and accurate object regions compared with
486 the (d) baseline, and are close to the (c) ground truth, which demonstrates
487 that the proposed F4S achieves a comparable performance to 5-shot without
488 increasing annotation cost.

489 *6.4. Ablation study*

490 We conduct a series of ablation studies to investigate the effectiveness of
491 each component in the proposed F4S and the results are shown in Table 5.
492 Without loss of generality, the ablation study experiments are performed on
493 “F4S(HSNet)” with ResNet101 backbone on COCO-20ⁱ dataset. In Table 5,
494 one can observe that when only with the E_{sc} , E_{imc} , or E_{cyc} , the proposed
495 method achieves mIoU improvement of 0.4%, 0.7% ,and 0.6% respectively,
496 and their combination leads to 2.3% mIoU improvement. Then, when only
497 using the inter-class confidence term T , the proposed method achieves mIoU
498 improvements of 8.9%, and FB-IoU improvements of 2.6%. Next, with the
499 existence of T , each component (E_{sc} , E_{imc} , and E_{cyc}) of the intra-class con-
500 fidence term R contributes further mIoU improvements to different extents,
501 which are shown in the 7th to 9th rows. Finally, the full combination of R
502 and T achieves the best mIoU of 51.4% and FB-IoU of 73.3%. The ablation
503 studies prove the effectiveness of both R and T in the F4S.

504 We notice that T contributes to larger mIoU improvement while R pro-
505 vides limited improvement. The reason is that the feature bias caused by
506 inter-class noise is greater than intra-class noise, which explains the greater
507 performance improvement of T . However, this does not mean that intra-class
508 noise can be ignored. The results in the 2nd to 5th rows of Table 5 show that R
509 is also essential for eliminating intra-class noise to improve FSS performance.

510 *6.5. Analysis*

511 *6.5.1. Computational analysis*

512 In Table 6, the 1st row shows the computational complexity of the base
513 model HSNet, which is regarded as the baseline. The 2nd row shows the
514 computational complexity of the proposed method in whole stages, including
515 generating (Stage I) and selecting (Stage II) pseudo labels. The 3rd to 5th
516 rows show the computational complexity of each stage respectively.

Table 5: Ablation study of F4S with different design choices. The results represent the mean metric scores of running 5 times. “ ± 0.1 ” indicates the standard deviation of running 5 times.

R			T	Fold-0	Fold-1	Fold-2	Fold-3	mean	FB-IoU
E_{sc}	E_{imc}	E_{cyc}							
				37.2	44.1	42.4	41.3	41.2	69.1
✓				37.9	45.7	41.8	41.1	41.6 (± 0.4)	69.3 (± 0.3)
	✓			38.5	44.6	42.3	42.0	41.9 (± 0.3)	69.8 (± 0.4)
		✓		38.7	45.1	41.8	41.7	41.8 (± 0.5)	69.6 (± 0.6)
✓	✓	✓		39.7	47.0	44.4	42.8	43.5 (± 0.7)	70.6 (± 0.6)
			✓	47.1	53.4	50.3	49.7	50.1 (± 0.4)	71.7 (± 0.5)
✓			✓	46.7	56.2	50.8	48.7	50.6 (± 0.8)	72.0 (± 0.4)
	✓		✓	47.6	55.8	49.6	49.0	50.5 (± 0.6)	71.9 (± 0.3)
		✓	✓	47.6	55.6	51.3	49.6	51.0 (± 0.4)	72.4 (± 0.4)
✓	✓	✓	✓	46.6	56.7	51.5	50.7	51.4 (± 0.2)	73.3 (± 0.3)

Table 6: Computational complexity of F4S compared with the baseline.

Method	Stage			Learnable Params \downarrow	FPS \uparrow	FLOPS(G) \downarrow
	I	II	III			
HSNet (<i>baseline</i>)	-	-	-	2.6M	16.33	20.56
F4S (HSNet)	✓	✓	✓	2.6M	5.08	81.66
	✓			0	15.80	20.52
		✓		0	8.51	40.62
			✓	2.6M	16.45	20.52

517 Specifically, in stage I (3rd row), the trained models of HSNet are officially
518 provided to generate pseudo labels. Therefore, there are no learnable params
519 in this stage, and the FPS and FLOPs are also close to the baseline. In stage
520 II (4th row), a diverged network N_{θ_2} is adopted here to compute E_{imc} in Eq. 6
521 and the base network N_{θ} is utilized to compute E_{cyc} in Eq. 7. Therefore, the
522 FLOPS increases to 40.62G and the FPS decreases to 8.51. In stage III (5th
523 row), F4S (HSNet) is retrained with pseudo labels. Therefore, the learnable
524 params is 2.6M, which is the same as the baseline. Besides, the FPS and
525 FLOPs of F4S (HSNet) are 16.45 and 20.52G, respectively, which are also
526 close to the baseline (16.33 and 20.56G).

527 Here we emphasize that although the proposed method has a high com-
528 putational complexity in whole stages (2nd row), the stage I and stage II only
529 need to be performed once before the training and testing stages, and do not
530 affect the computational complexity of the training and testing stages (5th
531 row). Therefore, in the actual testing process, the computational complexity

Table 7: Performance scores of different weight values. The results represent the mean metric scores of running 5 times. “ ± 0.1 ” indicates the standard deviation of running 5 times.

α	β	Fold-0	Fold-1	Fold-2	Fold-3	mean	FB-IoU
0.5	0.5	72.3	74.7	68.3	70.4	71.4 (± 0.1)	81.5 (± 0.1)
0.4	0.6	72.5	75.0	69.6	70.1	71.8 (± 0.2)	81.9 (± 0.1)
0.2	0.8	72.2	74.5	69.5	71.9	72.0 (± 0.1)	82.0 (± 0.1)
0.3	0.7	72.3	75.4	71.1	70.6	72.3 (± 0.1)	82.3 (± 0.1)

Table 8: Precomputed α and β on PASCAL-5ⁱ dataset.

	fold-1					fold-2				
	aeroplane	bicycle	bird	boat	bottle	bus	car	cat	chair	cow
α	0.14	0.19	0.20	0.22	0.11	0.27	0.25	0.16	0.10	0.26
β	0.86	0.81	0.80	0.78	0.89	0.73	0.75	0.84	0.90	0.74
	fold-3					fold-4				
	diningtable	dog	horse	motorbike	person	pottedplant	sheep	sofa	train	tvmonitor
α	0.21	0.27	0.26	0.18	0.34	0.12	0.29	0.17	0.30	0.24
β	0.79	0.73	0.74	0.82	0.66	0.88	0.71	0.83	0.70	0.76

532 of the inference remains unchanged compared to the baseline.

533 6.5.2. Weights settings

534 Table. 7 shows the quantitative scores when α and β in Eq. 1 are set
 535 to different values. The experiments are conducted on “F4S(HSNet)” with
 536 ResNet101 backbone on PASCAL-5ⁱ. One can observe that when $\alpha = 0.3$
 537 and $\beta = 0.7$, the best quantitative scores (72.3% mIoU and 82.3% FB-IoU)
 538 are obtained. Besides, we also find that by using different α and β , the quan-
 539 titative scores fluctuate within a narrow range ($<1.0\%$), which demonstrates
 540 the stability of the proposed F4S to α and β .

541 Moreover, we conduct experiments of precomputed α and β to obtain the
 542 “oracle” performance. The α and β indicate the ratio of intra- and inter-class
 543 samples in the noisy unlabeled images. Therefore, we count the quantity of
 544 intra- and inter-class samples of each class. We conduct experiments on
 545 PASCAL-5ⁱ dataset and the precomputed α and β of each class are shown
 546 in Table 8 and the “Oracle” results are shown in Table 9.

547 In Table 9, one can observe that with the precomputed α and β , the
 548 “Oracle” results of the proposed method achieve 73.4% (1-shot) and 73.9% (5-
 549 shot) of mIoU with ResNet101 backbone, which outperform “F4S (HSNet) ‡”
 550 with a sizable margin (1.1% and 1.1%). Besides, with VGG16 and ResNet50
 551 backbones, the “Oracle” results also achieve remarkable mIoU improvements.

Table 9: “Oracle” performance by precomputed α and β on PASCAL-5ⁱ dataset.

Backbone	Method	1-shot		5-shot	
		mIoU	FB-IoU	mIoU	FB-IoU
VGG16	F4S (HSNet) †	61.3 (\pm 0.3)	74.4 (\pm 0.3)	64.8 (\pm 0.2)	76.9 (\pm 0.2)
	F4S (HSNet) ‡	67.9 (\pm 0.2)	79.2 (\pm 0.1)	68.2 (\pm 0.3)	79.7 (\pm 0.3)
	Oracle	68.2	79.4	68.6	80.2
ResNet50	F4S (HSNet) †	64.8 (\pm 0.2)	77.2 (\pm 0.2)	70.1 (\pm 0.2)	81.0 (\pm 0.2)
	F4S (HSNet) ‡	70.8 (\pm 0.2)	81.5 (\pm 0.2)	72.0 (\pm 0.3)	82.2 (\pm 0.2)
	Oracle	71.9	82.4	72.5	83.0
ResNet101	F4S (HSNet) †	66.5 (\pm 0.2)	78.2 (\pm 0.2)	70.9 (\pm 0.3)	81.1 (\pm 0.2)
	F4S (HSNet) ‡	72.3 (\pm 0.2)	82.3 (\pm 0.2)	72.8 (\pm 0.2)	82.6 (\pm 0.3)
	Oracle	73.4	83.0	73.9	83.3

552 These results verify the effectiveness of precomputed α and β .

553 6.5.3. Statistical analysis of term R

554 To further investigate the terms E_{sc} , E_{imc} , E_{cyc} in the intra-class confi-
555 dence term R , we sample the image X from the annotated PASCAL-5ⁱ to
556 calculate $m(Y_X, \hat{Y}_X)$, where the ground truth Y_X is available and $m(\cdot, \cdot)$ is
557 set to mIoU score. Then, we calculate E_{sc} , E_{imc} , E_{cyc} following Sect. 4.2. In
558 Fig. 7, we plot the scatter graphs of (a) E_{sc} and $m(Y_X, \hat{Y}_X)$, (b) E_{imc} and
559 $m(Y_X, \hat{Y}_X)$, (c) E_{cyc} and $m(Y_X, \hat{Y}_X)$, (d) R and $m(Y_X, \hat{Y}_X)$ on the 4 folds of
560 PASCAL-5ⁱ. As can be noticed in Fig. 7 (a)-(c), there is a positive correlation
561 between $m(Y_X, \hat{Y}_X)$ and E_{sc} , E_{imc} , E_{cyc} . In Fig. 7 (d), the score R combin-
562 ing the three components contributes to better scatter dots distribution: the
563 dots mainly follow the line $y = x$, which presents a better positive correlation
564 between R and $m(Y_X, \hat{Y}_X)$. Therefore, the results of the scatter graphs prove
565 that the intra-class confidence term R can estimate the credibility of pseudo
566 labels, i.e., $m(Y_X, \hat{Y}_X)$, and thus identify the noisy intra-class samples.

567 6.5.4. F4S performance change with different numbers of unlabelled examples

568 We have investigated the F4S performance change with different numbers
569 of unlabelled examples. We choose “F4S (HSNet)” with ResNet101 backbone
570 as the model to conduct the experiments. Here, Table 10 and Table 11 show
571 the results on PASCAL-5ⁱ and COCO-20ⁱ datasets, respectively.

572 In Table 10 and Table 11, the “baseline” indicates the F4S performance
573 under the 1-shot setting without any additional unlabelled examples in the
574 test phase. The “+ N examples” indicates the F4S performance with addi-
575 tional unlabelled N examples, which are pseudo-labelled and selected by F4S.

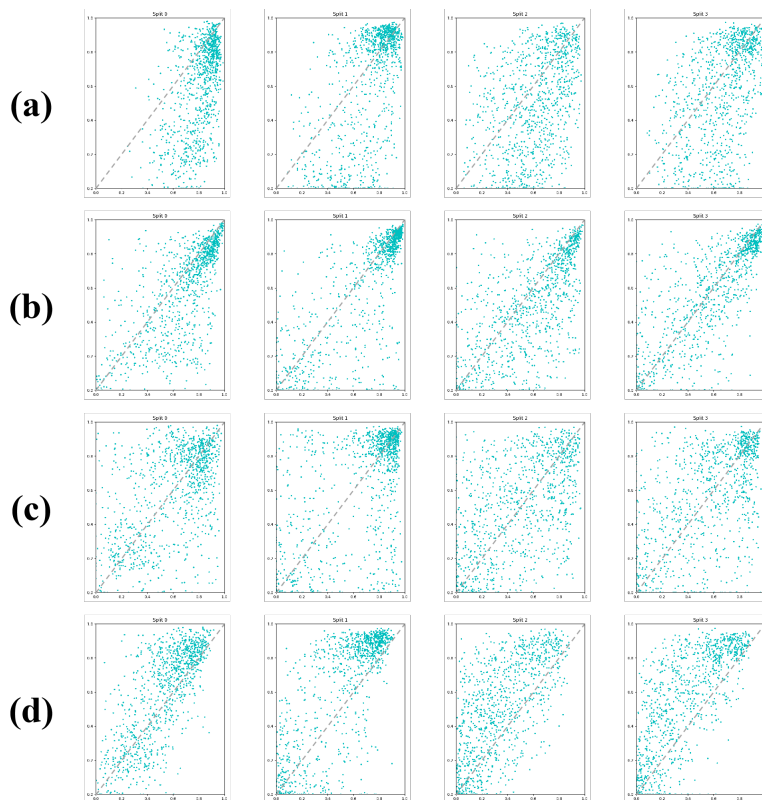


Figure 7: Scatter graphs of each term in score R . The y-axis indicates the mIoU score based on ground truth. The x-axis indicates the values of: (a) E_{sc} , (b) E_{imc} , (c) E_{cyc} , and (d) R . Each row shows the scatter graphs on the 4 folds of PASCAL-5ⁱ.

576 In Table 10, the “baseline” performance is 66.5% mIoU score and 78.2% FB-
 577 IoU score over 4 folds on the PASCAL-5ⁱ dataset. Then, with the increasing
 578 number of unlabelled examples, the performance scores of F4S also gradually
 579 improve. Finally, when with “+ 29 examples”, the proposed F4S achieves
 580 7.3% of mIoU improvements and 5.5% of FB-IoU improvements over the
 581 “baseline”. In Table 11, when with “+ 29 examples” on the COCO-20ⁱ
 582 dataset, the proposed F4S also outperforms “baseline” with a sizable margin
 583 as well, achieving 9.9% of mIoU improvements and 4.0% of FB-IoU improve-
 584 ments. Furthermore, we observed that with “+ 29 examples”, the perfor-
 585 mance eventually plateaus in both PASCAL-5ⁱ and COCO-20ⁱ datasets. This
 586 outcome is attributed to the increased number of pseudo-labeled examples
 587 with lower scores E .

Table 10: F4S performance change with different numbers of unlabelled examples on PASCAL-5ⁱ.

setting		Fold-0	Fold-1	Fold-2	Fold-3	mean	FB-IoU
baseline	1-shot	67.8	72.2	62.4	63.4	66.5 (± 0.2)	78.2 (± 0.2)
F4S	+ 4 examples	72.3	75.4	71.1	70.6	72.3 (± 0.1)	82.3 (± 0.1)
	+ 9 examples	73.0	76.0	72.2	71.6	73.2 (± 0.1)	83.4 (± 0.1)
	+ 19 examples	73.4	76.4	72.6	72.2	73.6 (± 0.1)	83.5 (± 0.2)
	+ 29 examples	73.5	76.5	72.8	72.6	73.8 (± 0.1)	83.7 (± 0.1)

Table 11: F4S performance change with different numbers of unlabelled examples on COCO-20ⁱ.

setting		Fold-0	Fold-1	Fold-2	Fold-3	mean	FB-IoU
baseline	1-shot	38.4	47.8	43.2	41.8	42.8 (± 0.2)	69.8 (± 0.2)
F4S	+ 4 examples	46.6	56.7	51.5	50.7	51.4 (± 0.2)	73.3 (± 0.3)
	+ 9 examples	47.5	56.6	52.1	50.6	51.7 (± 0.6)	73.6 (± 0.5)
	+ 19 examples	47.2	57.9	52.7	50.5	52.1 (± 0.6)	73.7 (± 0.6)
	+ 29 examples	48.2	58.9	52.8	50.8	52.7 (± 0.8)	73.8 (± 0.7)

588 6.6. Discussion

589 In this section, we introduce the task settings of few-shot learning and
 590 semi-supervised learning, and summarize the similarities and differences be-
 591 tween them.

592 **Setting of Few-shot Learning.** Few-shot learning (FSL) has a few
 593 available samples per class as the support set and aims to recognize the
 594 objects in the query set. In fact, FSL does not classify the data specifically,
 595 but makes a cluster to learn the similarity metric function [10]. Increasing
 596 the number of support images is a direct way to improve the performance
 597 of FSL models. However, it requires manual annotation and selection of
 598 high-quality intra-class data as new support images, which is a time- and
 599 labor-consuming process.

600 **Setting of Semi-Supervised Learning.** Semi-supervised learning (SSL)
 601 concerns with using labeled as well as unlabeled data to perform certain
 602 learning tasks. It permits harnessing the large amounts of unlabeled data
 603 available in many use cases in combination with typically smaller sets of
 604 labeled data [56]. Existing SSL methods based on deep neural networks
 605 can be categorized into: deep generative methods, consistency regularization
 606 methods, graph-based methods, pseudo-labeling methods, and hybrid meth-
 607 ods [57]. Our proposed method falls within the category of pseudo-labeling
 608 methods.

609 **Similarities.** Both few-shot learning and semi-supervised learning face
 610 the challenge of data scarcity. In the FSL, there are typically very few samples

611 available for training each category, while in the SSL, there is a small portion
612 of labeled training data and the rest is unlabeled. Besides, both FSL and
613 SSL place great demand on the model’s generalization capability. The FSL
614 and SSL models need to make accurate predictions on new data under data
615 scarcity.

616 **Differences.** Few-shot learning and Semi-supervised learning differ in
617 their primary objectives and approaches. FSL emphasizes how to effectively
618 recognize novel classes with very few labeled samples. Therefore, existing
619 FSL methods focus on the designing of network architectures, loss functions,
620 and optimizers to improve FSL performance. However, SSL concerns with
621 the utilization of unlabeled data to enhance supervised learning tasks. Taking
622 pseudo-labeling methods as an illustration, this type of method concentrates
623 on the generation of pseudo labels and the reduction of noise in order to
624 enhance the diversity of classes within the dataset, consequently facilitating
625 the supervised training of models.

626 7. Conclusion

627 We have presented a novel semi-supervised few-shot segmentation frame-
628 work named F4S, where noisy and unlabeled support images, e.g., from other
629 available datasets, are utilized to benefit both the training and test of few-
630 shot segmentation networks via generating pseudo labels. Due to the feature-
631 biased problem caused by noisy intra- and inter-class samples and resulting
632 in FSS performance degradation, we propose a ranking algorithm in F4S to
633 identify and eliminate the noisy samples via calculating and ranking con-
634 fidence scores of noisy support images. Specifically, the ranking algorithm
635 consists of an intra-class confidence score R to identify noisy intra-class sam-
636 ples based on their prediction confidence, and an inter-class confidence score
637 T to identify noisy inter-class samples based on channel-wise feature simi-
638 larity. Additionally, we have theoretically explained the effectiveness of the
639 proposed method based on a Structural Causal Model (SCM) from the view
640 of causal inference. We have conducted extensive experiments on PASCAL-5ⁱ
641 and COCO-20ⁱ datasets to validate the proposed method. Compared with re-
642 cent inductive and transductive FSS methods, the proposed method achieves
643 superior performance under 1-shot and 5-shot settings. Besides, the ablation
644 studies prove the effectiveness of each component in the score R and score
645 T .

646 The proposed work still has some primary limitations: (1) the computa-
647 tional complexity in the stage II of the proposed method is costly. How to
648 optimize the selection of pseudo labels to reduce the computational complex-
649 ity is a crucial concern in the future. (2) The underlying characteristics of
650 noisy samples need further investigation for designing the confidence score
651 E and making the selection of pseudo labels more reliable. We hope our
652 work may inspire the study of exploring the combination of semi-supervised
653 learning with few-shot segmentation task.

654 8. Acknowledgement

655 This work was supported in part by the National Key R&D Program
656 of China under Grant 2021ZD0112001, the National Natural Science Founda-
657 tion of China under Grant 62271119, the Natural Science Foundation of
658 Sichuan Province under Grant 2023NSFSC1972, the A*STAR AME Pro-
659 grammatic Funding A18A2b0046, the RobotHTPO Seed Fund under Project
660 C211518008, the EDB Space Technology Development Grant under Project
661 S22-19016-STDP, and the Natural Science Foundation of Jiangsu Province
662 under Grant BZ2021013.

663 References

- 664 [1] A. Shaban, S. Bansal, Z. Liu, I. Essa, B. Boots, One-shot learning for
665 semantic segmentation, British Machine Vision Conference (2017).
- 666 [2] K. Wang, J. H. Liew, Y. Zou, D. Zhou, J. Feng, Panet: Few-shot image
667 semantic segmentation with prototype alignment, in: Proceedings of
668 the IEEE/CVF International Conference on Computer Vision, 2019,
669 pp. 9197–9206.
- 670 [3] Y. Liu, X. Zhang, S. Zhang, X. He, Part-aware prototype network for
671 few-shot semantic segmentation, in: Computer Vision – ECCV 2020,
672 Springer International Publishing, Cham, 2020, pp. 142–158.
- 673 [4] B. Yang, C. Liu, B. Li, J. Jiao, Q. Ye, Prototype mixture models for
674 few-shot semantic segmentation, in: European Conference on Computer
675 Vision, Springer, 2020, pp. 763–778.

- 676 [5] L. Yang, W. Zhuo, L. Qi, Y. Shi, Y. Gao, Mining latent
677 classes for few-shot segmentation, in: 2021 IEEE/CVF Interna-
678 tional Conference on Computer Vision (ICCV), 2021, pp. 8701–8710.
679 doi:10.1109/ICCV48922.2021.00860.
- 680 [6] H. Sun, X. Lu, H. Wang, Y. Yin, X. Zhen, C. G.
681 Snoek, L. Shao, Attentional prototype inference for few-
682 shot segmentation, *Pattern Recognition* 142 (2023) 109726.
683 doi:https://doi.org/10.1016/j.patcog.2023.109726.
- 684 [7] C. Zhang, G. Lin, F. Liu, R. Yao, C. Shen, Canet: Class-agnostic seg-
685 mentation networks with iterative refinement and attentive few-shot
686 learning, in: *Proceedings of the IEEE/CVF Conference on Computer
687 Vision and Pattern Recognition*, 2019, pp. 5217–5226.
- 688 [8] Z. Tian, H. Zhao, M. Shu, Z. Yang, R. Li, J. Jia, Prior guided fea-
689 ture enrichment network for few-shot segmentation, *IEEE Transactions
690 on Pattern Analysis and Machine Intelligence* 44 (2) (2022) 1050–1065.
691 doi:10.1109/TPAMI.2020.3013717.
- 692 [9] H. Min, Y. Zhang, Y. Zhao, W. Jia, Y. Lei, C. Fan,
693 Hybrid feature enhancement network for few-shot seman-
694 tic segmentation, *Pattern Recognition* 137 (2023) 109291.
695 doi:https://doi.org/10.1016/j.patcog.2022.109291.
- 696 [10] Y. Song, T. Wang, P. Cai, S. K. Mondal, J. P. Sahoo, A comprehensive
697 survey of few-shot learning: Evolution, applications, challenges, and
698 opportunities, *ACM Computing Surveys* (2023).
- 699 [11] K. Huang, J. Geng, W. Jiang, X. Deng, Z. Xu, Pseudo-loss confidence
700 metric for semi-supervised few-shot learning, in: *Proceedings of the
701 IEEE/CVF International Conference on Computer Vision*, 2021, pp.
702 8671–8680.
- 703 [12] Y. Chen, C. Wei, D. Wang, C. Ji, B. Li, Semi-supervised contrastive
704 learning for few-shot segmentation of remote sensing images, *Remote
705 Sensing* 14 (17) (2022) 4254.
- 706 [13] Y. Tang, Z. Cao, Y. Yang, J. Liu, J. Yu, Semi-supervised few-shot object
707 detection via adaptive pseudo labeling, *IEEE Transactions on Circuits
708 and Systems for Video Technology* (2023).

- 709 [14] J. Li, R. Socher, S. C. Hoi, Dividemix: Learning with noisy labels as
710 semi-supervised learning, in: International Conference on Learning Rep-
711 resentations, 2019.
- 712 [15] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. Tsang, M. Sugiyama,
713 Co-teaching: Robust training of deep neural networks with extremely
714 noisy labels, Neural Information Processing Systems(NeurIPS) (2018).
- 715 [16] J.-Y. Zhu, T. Park, P. Isola, A. A. Efros, Unpaired image-to-image trans-
716 lation using cycle-consistent adversarial networks, in: Proceedings of the
717 IEEE international conference on computer vision, 2017, pp. 2223–2232.
- 718 [17] T. Zhou, P. Krahenbuhl, M. Aubry, Q. Huang, A. A. Efros, Learning
719 dense correspondence via 3d-guided cycle consistency, in: Proceedings
720 of the IEEE Conference on Computer Vision and Pattern Recognition,
721 2016, pp. 117–126.
- 722 [18] K. Nguyen, S. Todorovic, Feature weighting and boosting for few-shot
723 segmentation, in: Proceedings of the IEEE/CVF International Confer-
724 ence on Computer Vision, 2019, pp. 622–631.
- 725 [19] R. Zhang, H. Zhu, H. Zhang, C. Gong, J. T. Zhou, F. Meng, Semi-
726 supervised few-shot segmentation with noisy support images, in: 2023
727 IEEE International Conference on Image Processing (ICIP), IEEE, 2023,
728 pp. 1550–1554.
- 729 [20] N. Dong, E. Xing, Few-shot semantic segmentation with prototype
730 learning, in: British Machine Vision Conference, 2018.
- 731 [21] G. Li, V. Jampani, L. Sevilla-Lara, D. Sun, J. Kim, J. Kim, Adaptive
732 prototype learning and allocation for few-shot segmentation, in: 2021
733 IEEE/CVF Conference on Computer Vision and Pattern Recognition
734 (CVPR), 2021, pp. 8330–8339. doi:10.1109/CVPR46437.2021.00823.
- 735 [22] C. Lang, G. Cheng, B. Tu, J. Han, Learning what not to segment:
736 A new perspective on few-shot segmentation, in: Proceedings of the
737 IEEE/CVF conference on computer vision and pattern recognition,
738 2022, pp. 8057–8067.

- 739 [23] C. Lang, G. Cheng, B. Tu, C. Li, J. Han, Base and meta: A new
740 perspective on few-shot segmentation, *IEEE Transactions on Pattern*
741 *Analysis and Machine Intelligence* (2023).
- 742 [24] C. Lang, G. Cheng, B. Tu, J. Han, Few-shot segmentation via divide-
743 and-conquer proxies, *International Journal of Computer Vision* (2023)
744 1–23.
- 745 [25] Z. Hu, Z. Yang, X. Hu, R. Nevatia, Simple: Similar pseudo label
746 exploitation for semi-supervised classification, in: *Proceedings of the*
747 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*
748 (CVPR), 2021, pp. 15099–15108.
- 749 [26] M. Yang, J. Ling, J. Chen, M. Feng, J. Yang, Discrimina-
750 tive semi-supervised learning via deep and dictionary representa-
751 tion for image classification, *Pattern Recognition* 140 (2023) 109521.
752 doi:<https://doi.org/10.1016/j.patcog.2023.109521>.
- 753 [27] J. Li, G. Li, Y. Shi, Y. Yu, Cross-domain adaptive clustering for semi-
754 supervised domain adaptation, in: *Proceedings of the IEEE/CVF Con-*
755 *ference on Computer Vision and Pattern Recognition*, 2021, pp. 2505–
756 2514.
- 757 [28] K. Huang, J. Geng, W. Jiang, X. Deng, Z. Xu, Pseudo-loss confi-
758 dence metric for semi-supervised few-shot learning, in: *Proceedings of*
759 *the IEEE/CVF International Conference on Computer Vision (ICCV)*,
760 2021, pp. 8671–8680.
- 761 [29] M. Xu, Z. Zhang, H. Hu, J. Wang, L. Wang, F. Wei, X. Bai, Z. Liu, End-
762 to-end semi-supervised object detection with soft teacher, in: *Proceed-*
763 *ings of the IEEE/CVF International Conference on Computer Vision*
764 (ICCV), 2021, pp. 3060–3069.
- 765 [30] Y. Jin, J. Wang, D. Lin, Semi-supervised semantic segmentation via
766 gentle teaching assistant, *Advances in Neural Information Processing*
767 *Systems* 35 (2022) 2803–2816.
- 768 [31] Y. Wang, J. Zhang, M. Kan, S. Shan, Learning pseudo labels for semi-
769 and-weakly supervised semantic segmentation, *Pattern Recognition* 132
770 (2022) 108925. doi:<https://doi.org/10.1016/j.patcog.2022.108925>.

- 771 [32] P. Mazumder, P. Singh, V. P. Namboodiri, Rnnp: A robust few-shot
772 learning approach, in: Proceedings of the IEEE/CVF Winter Conference
773 on Applications of Computer Vision, 2021, pp. 2664–2673.
- 774 [33] J. Lu, S. Jin, J. Liang, C. Zhang, Robust few-shot learning for user-
775 provided data, IEEE transactions on neural networks and learning sys-
776 tems 32 (4) (2020) 1433–1447.
- 777 [34] O. B. Baran, R. G. Cinbis, Semantics-driven attentive few-shot learning
778 over clean and noisy samples, Neurocomputing 513 (2022) 59–69.
- 779 [35] K. J. Liang, S. B. Rangrej, V. Petrovic, T. Hassner, Few-shot learning
780 with noisy labels, in: Proceedings of the IEEE/CVF Conference on
781 Computer Vision and Pattern Recognition, 2022, pp. 9089–9098.
- 782 [36] Z. Chen, T. Ji, S. Zhang, F. Zhong, Noise suppression for improved few-
783 shot learning, in: ICASSP 2022-2022 IEEE International Conference on
784 Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2022, pp.
785 1900–1904.
- 786 [37] X. Luo, Z. Tian, T. Zhang, B. Yu, Y. Y. Tang, J. Jia, Pfenet++: Boost-
787 ing few-shot semantic segmentation with the noise-filtered context-aware
788 prior mask, IEEE Transactions on Pattern Analysis and Machine Intel-
789 ligence (2023).
- 790 [38] J. Pearl, M. Glymour, N. P. Jewell, Causal inference in statistics: A
791 primer. 2016, Google Scholar there is no corresponding record for this
792 reference (2016).
- 793 [39] D. B. Rubin, Essential concepts of causal inference: a remarkable history
794 and an intriguing future, Biostatistics & Epidemiology 3 (1) (2019) 140–
795 155.
- 796 [40] Z. Yue, H. Zhang, Q. Sun, X.-S. Hua, Interventional few-shot learning,
797 Advances in neural information processing systems 33 (2020) 2734–2746.
- 798 [41] D. Zhang, H. Zhang, J. Tang, X.-S. Hua, Q. Sun, Causal intervention
799 for weakly-supervised semantic segmentation, Advances in Neural Infor-
800 mation Processing Systems 33 (2020) 655–666.

- 801 [42] R. Wang, M. Yi, Z. Chen, S. Zhu, Out-of-distribution generaliza-
802 tion with causal invariant transformations, in: Proceedings of the
803 IEEE/CVF Conference on Computer Vision and Pattern Recognition,
804 2022, pp. 375–385.
- 805 [43] Y. Wang, X. Li, Z. Qi, J. Li, X. Li, X. Meng, L. Meng, Meta-causal fea-
806 ture learning for out-of-distribution generalization, in: European Con-
807 ference on Computer Vision, Springer, 2022, pp. 530–545.
- 808 [44] T. Zhang, H.-R. Shan, M. A. Little, Causal graphsage: A robust graph
809 method for classification based on causal sampling, *Pattern Recognition*
810 128 (2022) 108696. doi:<https://doi.org/10.1016/j.patcog.2022.108696>.
- 811 [45] L. G. Neuberg, *Causality: Models, reasoning, and inference*, by judea
812 pearl, cambridge university press, 2000, *Econometric Theory* 19 (4)
813 (2003) 675–685. doi:[10.1017/S0266466603004109](https://doi.org/10.1017/S0266466603004109).
- 814 [46] B. Zhu, Y. Niu, X.-S. Hua, H. Zhang, Cross-domain empirical risk min-
815 imization for unbiased long-tailed classification, in: Proceedings of the
816 AAAI Conference on Artificial Intelligence, Vol. 36, 2022, pp. 3589–3597.
- 817 [47] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan,
818 P. Dollár, C. L. Zitnick, Microsoft coco: Common objects in context,
819 2014, pp. 740–755.
- 820 [48] M. Everingham, L. V. Gool, C. K. I. Williams, J. Winn, A. Zisserman,
821 The pascal visual object classes (voc) challenge, *International Journal*
822 *of Computer Vision* 88 (2) (2010) 303–338.
- 823 [49] M. Boudiaf, H. Kervadec, Z. I. Masud, P. Piantanida, I. B. Ayed, J. Dolz,
824 Few-shot segmentation without meta-learning: A good transductive
825 inference is all you need?, in: 2021 IEEE/CVF Conference on Com-
826 puter Vision and Pattern Recognition (CVPR), 2021, pp. 13974–13983.
827 doi:[10.1109/CVPR46437.2021.01376](https://doi.org/10.1109/CVPR46437.2021.01376).
- 828 [50] J. Min, D. Kang, M. Cho, Hypercorrelation squeeze for few-shot segme-
829 nation, in: 2021 IEEE/CVF International Conference on Computer Vi-
830 sion (ICCV), 2021, pp. 6921–6932. doi:[10.1109/ICCV48922.2021.00686](https://doi.org/10.1109/ICCV48922.2021.00686).

- 831 [51] G. Cheng, C. Lang, J. Han, Holistic prototype activation for few-shot
832 segmentation, *IEEE Transactions on Pattern Analysis and Machine In-*
833 *telligence* 45 (4) (2023) 4650–4666. doi:10.1109/TPAMI.2022.3193587.
- 834 [52] Y. Lu, X. Wu, Z. Wu, S. Wang, Cross-domain few-shot segmentation
835 with transductive fine-tuning, arXiv preprint arXiv:2211.14745 (2022).
- 836 [53] X. Shi, D. Wei, Y. Zhang, D. Lu, M. Ning, J. Chen, K. Ma, Y. Zheng,
837 Dense cross-query-and-support attention weighted mask aggregation
838 for few-shot segmentation, in: *Computer Vision – ECCV 2022*, Springer
839 Nature Switzerland, Cham, 2022, pp. 151–168.
- 840 [54] D. Kang, M. Cho, Integrative few-shot learning for classification and seg-
841 mentation, in: *Proceedings of the IEEE/CVF Conference on Computer*
842 *Vision and Pattern Recognition*, 2022, pp. 9979–9990.
- 843 [55] S. Kim, P. Chikontwe, S. An, S. H. Park, Uncertainty-aware semi-
844 supervised few shot segmentation, *Pattern Recognition* 137 (2023)
845 109292.
- 846 [56] J. E. Van Engelen, H. H. Hoos, A survey on semi-supervised learning,
847 *Machine learning* 109 (2) (2020) 373–440.
- 848 [57] X. Yang, Z. Song, I. King, Z. Xu, A survey on deep semi-supervised
849 learning, *IEEE Transactions on Knowledge and Data Engineering*
850 (2022).