



Figure 5: Qualitative segmentation results on novel classes on PASCAL-5<sup>i</sup>. From left to right: support image with mask, query input, query ground-truth mask, query prediction of the baseline, and prediction of RiFeNet.

| Un | MP | split0      | split1      | split2      | split3      | mIoU        |
|----|----|-------------|-------------|-------------|-------------|-------------|
|    |    | 65.7        | 71.0        | 59.5        | 59.7        | 64.0        |
| ✓  |    | 67.3        | 71.8        | 66.2        | 59.2        | 66.1        |
|    | ✓  | 66.0        | 72.1        | 66.2        | <b>60.4</b> | 66.2        |
| ✓  | ✓  | <b>68.4</b> | <b>73.5</b> | <b>67.1</b> | 59.4        | <b>67.1</b> |

Table 3: Ablation studies on the key components of RiFeNet. “Un” and “MP” denote the use of the unlabeled branch and the multi-level prototypes, respectively.

| components        | split0      | split1      | split2      | split3      | mIoU        |
|-------------------|-------------|-------------|-------------|-------------|-------------|
| gp (support-only) | 67.3        | 71.8        | 66.2        | 59.2        | 66.1        |
| gp+gp             | 67.5        | 73.1        | 66.2        | 58.4        | 66.3        |
| gp+lp (w/o CA)    | 68.1        | 73.2        | 66.7        | 59.1        | 66.8        |
| gp+lp (w/ CA)     | <b>68.4</b> | <b>73.5</b> | <b>67.1</b> | <b>59.4</b> | <b>67.1</b> |

Table 4: Ablation studies on multi-level prototypes. “gp” and “lp” denote global and local prototypes, respectively. That is, “gp+gp” means extracting both query and support prototypes globally. “CA” refers to channel-wise attention.

is consistent with our analysis in Sec.3.3. On the other hand, because the unlabeled inputs come from resampling the training dataset, we double the training iterations of the baseline for a fair comparison. Increased training iterations have little effect on the baseline due to early convergence. This proves that the effectiveness of our method is not from the multiple sampling of data but from the learned discriminative and semantic features.

**Different hyper-parameters.** We first look into the effect of different numbers of unlabeled input in a single meta-training process. Tab.6 shows the results on PASCAL-5<sup>i</sup> under a 1-shot setting, with ResNet50 as its backbone. The best results are obtained when the number of unlabeled images is set to 2. Initially, the segmentation effect of the model

| components     | epoch | split0      | split1      | split2      | split3      | mIoU        |
|----------------|-------|-------------|-------------|-------------|-------------|-------------|
| w/o unlabeled  | 200   | 66.0        | 72.1        | 66.2        | <b>60.4</b> | 66.2        |
| w/o unlabeled  | 400   | 66.5        | 72.4        | 65.5        | 59.5        | 66.0        |
| un (w/o guide) | 200   | 66.9        | 72.2        | 65.9        | 58.3        | 65.8        |
| un (w/ guide)  | 200   | <b>68.4</b> | <b>73.5</b> | <b>67.1</b> | 59.4        | <b>67.1</b> |

Table 5: Ablation studies on the unlabeled branch. “w/ guide” refers to the use of query local prototypes in the unlabeled branch for guidance, while “w/o guide” means using prototypes generated from the unlabeled branch itself.

| num | split0      | split1      | split2      | split3      | mIoU        |
|-----|-------------|-------------|-------------|-------------|-------------|
| 0   | 66.0        | 72.1        | 66.2        | 60.4        | 66.2        |
| 1   | 66.8        | 72.8        | 66.9        | 59.8        | 66.6        |
| 2   | <b>68.4</b> | <b>73.5</b> | <b>67.1</b> | 59.4        | <b>67.1</b> |
| 3   | 65.9        | 72.6        | 66.9        | <b>59.8</b> | 66.0        |

Table 6: Ablation studies of different numbers of unlabeled images in the single meta-training process.

increased as the number of unlabeled images increased. When the number continues to increase, the accuracy decreases instead. We deem the reason is that when the effect of unlabeled enhancement counts much more than the query branch itself, the attention of feature mining may turn to the unlabeled branch, thus disturbing the query prediction. The segmentation accuracy decreases after the features are blurred. We also conduct detailed ablation experiments with other parameters, which are included in the Appendix.

## 5 Conclusion

In few-shot segmentation, traditional methods suffer from semantic ambiguity and inter-class similarity. Thus from the perspective of pixel-level binary classification, we propose RiFeNet, an effective model with an unlabeled branch constraining foreground semantic consistency. Without extra data, this unlabeled branch improves the in-class generalization of the foreground. Moreover, we propose to further enhance the discrimination of background and foreground by a multi-level prototype generation and interaction module. Rerum doloribus voluptatem, quibusdam aperiam mollitia culpa?Maiores ipsa quasi, hic aut laborum neque dicta maiores vel?Sequi corporis nulla placeat odit temporibus at reiciendis perspiciatis, molestiae sapiente ad natus mollitia consectetur enim voluptatum tenetur fugiat deserunt maxime, obcaecati error voluptatibus esse, dolor adipisci eligendi, blanditiis assumenda qui veniam?Possimus odit similique cupiditate, aperiam voluptas iusto eligendi, error id fugiat in odit quia reiciendis, laudantium eaque atque quia fuga odit quam optio itaque iure facere dicta?Ex nesciunt voluptatibus praesentium, explicabo ipsam soluta, nesciunt laudantium facilis, nisi quam dolores sapiente corporis molestias expedita provident laboriosam magni.Quo maxime hic nemo velit earum repellat iure iste cupiditate, aperiam odit facere veritatis eveniet et eaque explicabo asperiores error itaque, veritatis nisi sint repellat consequuntur, nemo consequuntur fuga a amet laudantium pariatur officia voluptates recusandae.Aspernatur ea accusantium vel eum modi,

voluptas sapiente magni?