1   **R2's Q1: In vanilla SSL task.** Sorry that we don't have this experiment. Our thought is that if taking pseudo-labels
2   as noisy labels, perhaps one may refer to the supervised re-weighting methods [A][B]. **Q2: One issue of results[37],**
3   **52.8<53.8.** This is because the training splits of dataset are different between SSFSC and FSC (see details in the Sec.
4   4.4 of [37]). On the FSC dataset, a big proportion of labeled data are used as unlabeled for sampling SSFSC tasks.
5   Therefore, the total SSFSC training tasks contain less supervision (than FSC). **Q3: State of the art SSL methods.**
6   Actually, we began our project by trying Virtual Adversarial Training (VAT) [15] which has been shown top-performing
7   in most settings of vanilla SSL [18]. We found that VAT brings limited improvement, e.g. less than $1\%$ on miniImageNet
8   1-shot, and it works slightly better for 5-shot. We think this is because of the high-variance of FSC classifiers trained
9   with very limited supervision. In contrast, our method can greatly increase this supervision by carefully choosing
10  and weighting high-confidence pseudo labels, and thus can make a visible improvement, e.g., $9\%$ over the supervised
11  baseline on miniImageNet 1-shot. We agree that distracting class is a challenge (kindly refer to **R3**'s **Q5**). One of our
12  future works is to find out an effective way of deploying regularization-based SSL methods to tackle this.

13  **R3's Q1: "by us".** It means we implement the open-sourced MTL code on the
14  tieredImageNet. **Q2: Comparing with FSC.** We agree this is unfair in terms of (1)
15  the additional unlabeled data in each single SSFSC task or (2) the muted labels on
16  the whole dataset (kindly refer to **R2**'s **Q2**). For paper revision, we will preserve only
17  the comparison to baseline supervised methods and remove others. We will add more
18  results related to SSFSC, e.g., Figure B1. **Q4: Vary module selections: "+recursive"**
19  **or "+mixing".** Sorry for the confusion. "+recursive" and "+mixing" are actually
20  the same method (LST) with different hyperparameters, not different modules. E.g.
21  in 5-way, 1-shot setting, "+recursive" has 6 recursive stages, and every stage it uses
22  $5 \times 30$ unlabeled samples. While, "+mixing" has only one stage, using $6 \times 5 \times 30$
23  samples for once. **Q5: Quantitative analysis for the number of distracting classes.**
24  Our experiment results in Figure B1 show that both our LST and related methods,
25  Soft $k$-Means [22] and TPN [37], are obviously affected by distracting classes. Other
26  observations are that (1) LST achieves top performances, especially more than $2\%$
27  higher than TPN [37] at $classNum = 7$; (2) LST with less re-training steps, i.e., a
28  smaller $m$ value, works better for reducing the effect from a larger number of distracting
29  classes. **Q6: Require a large number of unlabeled samples.** In the supplementary
30  Table S1, we provided the results of using $5$ unlabeled samples (LINE 22-30) for both
31  our LST (w/o *recursive*) and related methods [22][37], validating our superiority in the
32  low-data settings. Note that in the Table 2 of the main paper, we reported the results of
33  LST (*recursive,hard,soft*) and related works using the same number of unlabeled data.



miniImageNet

tieredImageNet

Figure B1: The 5-way, 1-shot
results using different num-
bers of distracting classes.

34  **Q7: Distracting classes in Formula 5. SWN for distracting classes.** (1) Samples from distracting classes are mixed
35  with other unlabeled data without distinction, thus have no special role in Formula 5. (2) SWN does reduce the effect
36  of distracting classes. When comparing "*recursive,hard*" to "*recursive,hard,soft*" in Table 2 (w/$\mathcal{D}$), we can see the
37  improvements ($2.3\% \sim 5.2\%$) (*soft* = using SWN). **Q8: Accuracy of pseudo label.** Taking the miniImageNet 1-shot as
38  an example, during meta-training episodes, we can see the accuracy growing from $47.0\%$ (iter=0) to $71.5\%$ (iter=15$k$).
39  There are 6 recursive stages during meta-test. From stage-1 to stage-6, the average accuracy (of 600 meta-test episodes)
40  increases from $63.6\%$ ($62.2\%$ w/o *soft* weighting) to $68.8\%$ ($66.1\%$ w/o *soft* weighting). Detailed numbers will be
41  reported in our paper. **Q9: Insufficient aspect.** LST has some discrete hyperparameters (e.g., the numbers of hard
42  selected samples and recursive stages) that are manually set. Our future work is to make them optimizable.

43  **R4's Q2 (R3's Q3): MTL helps SSFSC.** MTL transfers the superior pre-trained DNN weights for efficient feature
44  extraction in unseen classes. It is independent from the learning method, either supervised or semi-supervised, for
45  base classifiers. Our implementations of MTL in three methods, [22][37] and ours, validate its efficiency for SSFSC.
46  **Q1: Without finetuning.** Please kindly refer to Figure 3(a). "$m = 40$" means the number of re-training steps is equal
47  to the number of total steps (40), i.e., without finetuning step. Its corresponding curve clearly drops after the 18-th
48  iteration. **Q3: Using other backbones/FSC approaches.** We incorporate the 4CONV arch. of MAML [3] and the
49  recent FSC method LEO [25] into our LST, respectively. For example, on tieredImageNet 1-shot, LST-MAML-4CONV
50  outperforms TPN-4CONV[37] by $2.9\%$ and $2.0\%$(w/$\mathcal{D}$). LST-LEO-ResNet12 outperforms TPN-ResNet12 by $3.8\%$
51  and $2.8\%$(w/$\mathcal{D}$). Other results will be reported in the final paper. **Q3: Sample statistics.** For example, in 5-way,
52  1-shot case, we use 1 labeled and 20 unlabeled samples *per class* to meta-train SWN. In each meta-test task, we have 6
53  recursive (base-)training stages. At each stage, we select 100 samples (globally ranked by pseudo-labeling confidences)
54  out of $5 \times 30$ unlabeled inputs, and then weight them by the SWN. If using distracting classes, we simply add 30
55  samples *per distracting class* to the input, without distinction. Please kindly refer to LINE 193-199 for more details.

56  [A] Ren et al. Learning to Reweight Examples for Robust Deep Learning. *ICLR'18*. [B] Jiang et al. MentorNet: Learning
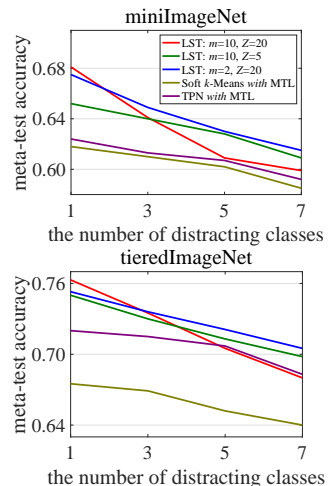57  Data-Driven Curriculum for Very Deep Neural Networks on Corrupted Labels. *ICML'18*.