

# Iterative label cleaning for transductive and semi-supervised few-shot learning

Michalis Lazarou<sup>1</sup> Yannis Avrithis<sup>2</sup> Tania Stathaki<sup>1</sup>  
<sup>1</sup>Imperial College London  
<sup>2</sup>Inria, Univ Rennes, CNRS, IRISA

## Abstract

*Few-shot learning amounts to learning representations and acquiring knowledge such that novel tasks may be solved with both supervision and data being limited. Improved performance is possible by transductive inference, where the entire test set is available concurrently, and semi-supervised learning, where more unlabeled data is available. These problems are closely related because there is little or no adaptation of the representation in novel tasks.*

*Focusing on these two settings, we introduce a new algorithm that leverages the manifold structure of the labeled and unlabeled data distribution to predict pseudo-labels, while balancing over classes and using the loss value distribution of a limited-capacity classifier to select the cleanest labels, iteratively improving the quality of pseudo-labels. Our solution sets new state of the art on four benchmark datasets, namely miniImageNet, tieredImageNet, CUB and CIFAR-FS, while being robust over feature space pre-processing and the quantity of available data.*

## 1. Introduction

*Few-shot learning* [60, 55] is challenging the deep learning paradigm in that, not only supervision is limited, but data is limited too. Despite the initial promise of *meta-learning* [38, 12], *transfer learning* [10, 58] is becoming increasingly successful in decoupling representation learning from learning novel tasks on limited data. *Semi-supervised learning* [29, 5] is one of the dominant ways of dealing with limited supervision and indeed, its few-shot learning counterparts [49, 65] are miniature versions where both labeled and unlabeled are limited proportionally, while representation learning may be decoupled. In this sense, these methods are closer to *transductive inference* [35, 50], which was a pillar of semi-supervised learning before deep learning [8].

Predicting *pseudo-labels* on unlabeled data [29] is one of the oldest ideas in semi-supervised learning [53]. Graph-based methods, in particular *label propagation* [67, 66], are prominent in transductive inference and translate to inductive inference in deep learning exactly by predicting pseudo-

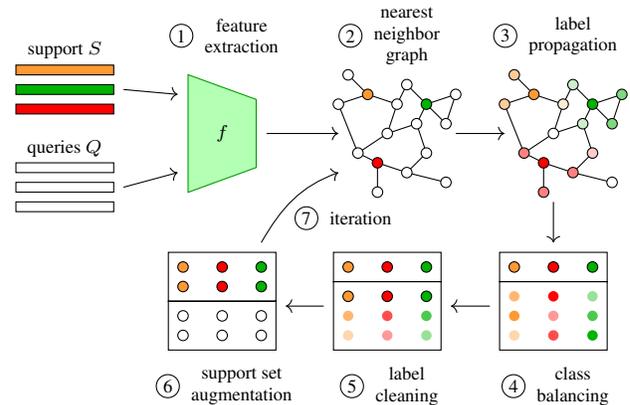


Figure 1. Overview of the proposed method.

labels [21]. However, with the representation being fixed, the quality of pseudo-labels is critical in few-shot learning [62, 28]. At the same time, in *learning with noisy labels* [3, 20, 56], it is common to clean labels based on the loss value statistics of a small-capacity classifier.

In this work, we leverage these ideas to improve transductive and semi-supervised few-shot learning. As shown in **Figure 1**, focusing on transduction, a set of labeled *support* examples  $S$  and unlabeled *queries*  $Q$  are given, represented in a *feature space* by mapping  $f$ . By *label propagation* [66], we obtain a matrix that associates examples to classes. The submatrix corresponding to unlabeled examples,  $P$ , is *normalized* over examples and classes using the Sinkhorn-Knopp algorithm [23], assuming a uniform distribution over classes. We extract pseudo-labels from  $P$ , which we clean following O2U-Net [20], keeping only one example per class. Finally, inspired by [25], we move these examples from  $Q$  to  $S$  and iterate until  $Q$  is empty.

## 2. Related work and contributions

### 2.1. Few-shot learning

**Meta-learning** This is a popular paradigm, where the training set is partitioned in episodes resembling in struc-

ture the novel tasks [12, 24, 55, 38, 60]. *Model-based* methods rely on the properties of specific model architectures, such as recurrent and memory-augmented networks [38, 52, 39]. *Optimization-based* methods attempt to learn model parameters that are able to adapt fast in novel tasks [12, 47, 68, 48, 40, 30, 6]. *Metric-based* methods attempt to learn representations that are appropriate for comparisons [55, 24, 57, 60]. Of course, *metric learning* is a research area on its own [61, 22] and modern ideas are commonly effective in few-shot learning [33].

**Predicting weights, data augmentation** Also based on meta-learning, it is possible to predict new parameters or even data. For instance, it is common to learn to *predict data-dependent network parameters* in the last layer (classifier) [13, 44, 46] or even in intermediate convolutional layers [7]. Alternatively, one can learn to *generate novel-task data* in the feature space [14] or in the input (image) space [63, 2]. The quality can be improved by translating images, similar to style transfer [34]. Such *learned data augmentation* is complementary to other ideas.

**Transfer learning** More recently, it is recognized that learning a powerful representation on the entire training set is more effective than sampling few-shot training episodes that resemble novel tasks [13, 10, 58, 37]. In doing so, one may use standard loss functions [13, 10], knowledge distillation [58] or other common self-supervision and regularization methods [37]. We follow this *transfer learning* approach, which allows us to decouple representation learning from the core few-shot learning idea and provide clearer comparisons with the competition.

## 2.2. Using unlabeled data

Leveraging unlabelled data is of interest due to the ease of obtaining such data. Two common settings are transductive inference and semi-supervised learning.

**Transductive inference** In this setting, all novel-class unlabeled query examples are assumed available at the same time at inference [35, 32, 45, 18, 19, 50]. These examples give additional information on the distribution of novel classes on top of labeled support examples.

Common transductive inference solutions are adapted for few-shot classification, notably *label propagation* [35] and *embedding propagation* [50], which smooths embeddings as in image segmentation [4]. Both operations are also used at representation learning, as in meta-learning. Using dimensionality reduction, TAFSSL [32] learns a task-specific feature subspace that is highly discriminant for novel tasks. *Meta-confidence transduction* (MCT) [28] meta-learns a data-dependent scaling function to normalize every example and iteratively updates class centers. PT+MAP [19] uses a similar iterative process but also balances over classes. *Cross-attention* [16], apart from aligning feature maps by

correlation, leverages query examples by iteratively making predictions and using the most confident ones to update the class representation.

**Semi-supervised learning** In this case, labeled novel-class support examples and additional unlabeled data are given. A classifier may be learned on both to make predictions on novel-class queries [49, 62, 65].

One of the first contributions uses unlabelled examples to adapt prototypical networks [55], while discriminating from distractor classes [49]. Common semi-supervised solutions are also adapted to few-shot classification, for instance *learning to self-train* [31], which adapts *pseudo-label* [29] and TransMatch [65], which is an adaptation of MixMatch [5]. *Instance credibility inference* [62] predicts pseudo-labels iteratively, using a linear classifier to select the most likely to be correct and then augmenting the support set. *Adaptive subspaces* [54] are learned from labeled and unlabeled data, yielding a discriminative subspace classifier that maximizes the margin between subspaces.

## 2.3. Contributions

In this work, focusing on the transfer learning paradigm to learn novel tasks given a fixed representation [58, 37], we make the following contributions:

1. We combine the power of *predicting pseudo-labels* in semi-supervised learning [29, 31] with *label cleaning* in learning from noisy labels [20].
2. According to manifold assumption, we use *label propagation* [66, 35] to infer pseudo-labels, while *balancing* over classes [23, 19] and iteratively *re-using pseudo-labels* in the propagation process [25, 62].
3. We achieve new state of the art in both *transductive* and *semi-supervised* few-shot learning.

## 3. Method

### 3.1. Problem formulation

At *representation learning*, we assume access to a labelled dataset  $D_{\text{base}}$  with each example having a label in one of the classes in  $C_{\text{base}}$ . This dataset is used to learn a mapping  $f : \mathcal{X} \rightarrow \mathbb{R}^d$  from an input space  $\mathcal{X}$  to a  $d$ -dimensional *feature* or *embedding* space.

The knowledge acquired at representation learning is used to solve *novel tasks*, assuming access to a dataset  $D_{\text{novel}}$  with each example being associated with one of the classes  $C_{\text{novel}}$ , where  $C_{\text{novel}}$  is disjoint from  $C_{\text{base}}$ . Examples in  $D_{\text{novel}}$  may be labeled or not.

In *few-shot classification* [60], a novel task is defined by sampling a *support set*  $S$  from  $D_{\text{novel}}$ , consisting of  $N$  classes with  $K$  labeled examples per class, for a total of  $L := NK$  examples. Given the mapping  $f$  and the support set  $S$ , the problem is to learn an  $N$ -way classifier that makes

predictions on unlabeled *queries* also sampled from  $D_{\text{novel}}$ . Queries are treated independently of each other. This is referred to as *inductive inference*.

In *transductive inference*, a *query set*  $Q$  consisting of  $M$  unlabeled examples is also sampled from  $D_{\text{novel}}$ . Given the mapping  $f$ ,  $S$  and  $Q$ , the problem is to make predictions on  $Q$ , without necessarily learning a classifier. In doing so, one may leverage the distribution of examples in  $Q$ , which is important because  $M$  is assumed greater than  $L$ .

In *semi-supervised* few-shot classification, an unlabelled set  $U$  of  $M$  unlabeled examples is also sampled from  $D_{\text{novel}}$ . Given  $f$ ,  $S$  and  $U$ , the problem is to learn to make predictions on new queries from  $D_{\text{novel}}$ , as in inductive inference. Again,  $M > L$  and we may leverage the distribution of  $U$ .

In this work, we focus on transductive inference and semi-supervised classification, given  $f$ . The performance of  $f$  on inductive inference is one of our baselines. We develop our solution for transductive inference. In the semi-supervised case, we follow the same solution with  $Q$  replaced by  $U$ . Using the predictions on  $U$ , we then proceed as in the inductive case, with  $S$  replaced by  $S \cup U$ .

### 3.2. Nearest-neighbor graph

We are given the mapping  $f$ , the labeled support set  $S := \{(x_i, y_i)\}_{i=1}^L$  and the query set  $Q := \{x_{L+i}\}_{i=1}^M$ , where  $y_i \in [N] := \{1, \dots, N\}$ . We embed all examples from  $S$  and  $Q$  into  $V = \{\mathbf{v}_1, \dots, \mathbf{v}_T\} \subset \mathbb{R}^d$ , where  $T := L + M$  and  $\mathbf{v}_i := f(x_i)$  for  $i \in [T]$ . Following [21], we construct a  $k$ -nearest neighbour graph of the features in  $V$ , represented by a sparse  $T \times T$  nonnegative *affinity matrix*  $A$ , with

$$A_{ij} := \begin{cases} [\mathbf{v}_i^\top \mathbf{v}_j]^\gamma, & \text{if } i \neq j \wedge \mathbf{v}_i \in \text{NN}_k(\mathbf{v}_j) \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

for  $i \in [T]$ ,  $j \in [N]$ , where  $\text{NN}_k(\mathbf{v})$  are the  $k$ -nearest neighbors of  $\mathbf{v}$  in  $V$  and  $\gamma > 1$  is a hyperparameter. Finally, we obtain the symmetric  $T \times T$  *adjacency matrix*  $W := \frac{1}{2}(A + A^\top)$  and we symmetrically normalize it as

$$\mathcal{W} := D^{-1/2} W D^{-1/2}, \quad (2)$$

where  $D = \text{diag}(W \mathbf{1}_T)$  is the  $T \times T$  degree matrix of  $W$ .

### 3.3. Label propagation

Following [66], we define the  $T \times N$  *label matrix*  $Y$  as

$$Y_{ij} := \begin{cases} 1, & \text{if } i \leq L \wedge y_i = j \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

for  $i \in [T]$ ,  $j \in [N]$ . Matrix  $Y$  has one column per class and one row per example, which is an one-hot encoded label for  $S$  and a zero vector for  $Q$ . Label propagation amounts to solving  $N$  linear systems

$$Z := (I - \alpha \mathcal{W})^{-1} Y, \quad (4)$$

where  $\alpha \in [0, 1)$  is a hyperparameter. The resulting  $T \times N$  matrix  $Z$  can be used to make predictions by taking the maximum element per row [66]. However, before making predictions, we balance over classes.

### 3.4. Class balancing

We focus on the  $M \times N$  submatrix

$$P := Z_{L+1:T,:} \quad (5)$$

(the last  $M$  rows) of  $Z$  that corresponds to unlabeled queries. We first perform an element-wise *power transform*

$$P_{ij} \leftarrow P_{ij}^\tau \quad (6)$$

for  $i \in [M]$ ,  $j \in [N]$ , where  $\tau > 1$ , encouraging hard predictions. Parameter  $\tau$  is analogous to the *scale* (or *inverse temperature*) of logits in softmax-based classifiers [13, 44, 41], only here the elements of  $P$  are proportional to class probabilities rather than logits.

Inspired by [19], we normalize  $P$  to a given row-wise sum  $\mathbf{p} \in \mathbb{R}^M$  and column-wise sum  $\mathbf{q} \in \mathbb{R}^N$ . Each element  $p_i \in [0, 1]$  of  $\mathbf{p}$  represents a confidence of example  $x_{L+i}$  for  $i \in [M]$ ; it can be a function of the  $i$ -th row of  $P$  or set to 1. Each element  $q_j \geq 0$  of  $\mathbf{q}$  represents a weight of class  $j$  for  $j \in [N]$ . In the absence of such information, we set

$$\mathbf{q} := \frac{1}{N} (\mathbf{p}^\top \mathbf{1}_M) \mathbf{1}_N, \quad (7)$$

assuming a uniform distribution of queries over classes.

The normalization itself is a projection of  $P$  onto the set  $\mathbb{S}(\mathbf{p}, \mathbf{q})$  of nonnegative  $M \times N$  matrices having row-wise sum  $\mathbf{p}$  and column-wise sum  $\mathbf{q}$ ,

$$\mathbb{S}(\mathbf{p}, \mathbf{q}) := \{X \in \mathbb{R}^{M \times N} : X \mathbf{1}_N = \mathbf{p}, X^\top \mathbf{1}_M = \mathbf{q}\}. \quad (8)$$

We use the *Sinkhorn-Knopp* algorithm [23] for this projection, which alternates between rescaling the rows of  $P$  to sum to  $\mathbf{p}$  and its columns to sum to  $\mathbf{q}$ ,

$$P \leftarrow \text{diag}(\mathbf{p}) \text{diag}(P \mathbf{1}_N)^{-1} P \quad (9)$$

$$P \leftarrow P \text{diag}(P^\top \mathbf{1}_M)^{-1} \text{diag}(\mathbf{q}), \quad (10)$$

until convergence. Finally, for each query  $x_{L+i}$ ,  $i \in [M]$ , we predict the *pseudo-label*

$$\hat{y}_{L+i} := \arg \max_{j \in [N]} P_{ij} \quad (11)$$

that corresponds to the maximum element of the  $i$ -th row of the resulting matrix  $P$ , for  $i \in [M]$ .

### 3.5. Label cleaning

The predicted pseudo-labels are not necessarily correct, yet a classifier can be robust to such noise. This is the case

when enough data is available to adapt the representation [29, 21], such that the quality of pseudo-labels improves with training. Since data is limited here, we would like to select pseudo-labeled queries in  $Q$  that are most likely to be correct, treat them as truly labeled and add them to the support set  $S$ . Iterating this process is an alternative way of improving the quality of pseudo-labels.

We interpret this problem as *learning with noisy labels*, leveraging recent advances in label cleaning [3, 20, 56]. Assuming that the classifier does not overfit the data, *e.g.* with small capacity, high learning rate or few iterations, the principle is that examples with clean labels exhibit less loss than examples with noisy labels.

In particular, given the labeled support set  $S := \{(x_i, y_i)\}_{i=1}^L$  and the pseudo-labeled query set  $\hat{Q} := \{(x_{L+i}, \hat{y}_{L+i})\}_{i=1}^M$ , we train an  $N$ -way classifier  $g$  using a weighted *cross-entropy* loss

$$\ell := - \sum_{i=1}^L \log g(x_i)_{y_i} - \sum_{i=1}^M p_i \log g(x_{L+i})_{\hat{y}_{L+i}}, \quad (12)$$

where  $p_i$  is the confidence weight of example  $x_{L+i}$ . Here, the classifier  $g$  is assumed to yield a vector of probabilities over classes using softmax and  $g(x)_y$  refers to element  $y \in [N]$  of  $g(x)$ . In practice, it is obtained by a linear classifier on top of embedding  $f$ , optionally allowing the adaptation of the last layers of the network implementing  $f$ .

The loss term  $\ell_i := -p_i \log g(x_{L+i})_{\hat{y}_{L+i}}$ , corresponding to the pseudo-labeled query  $x_{L+i}$ , is used for selection. Following O2U-Net [20], we use a *cyclical schedule* of learning rate and collect the average loss  $\bar{\ell}_i$  over all epochs, for  $i \in [M]$ . In learning with noisy labels, it is common to detect noisy labels based on statistics of this loss value for clean and noisy labels [3, 56]. However, this does not work well with predicted pseudo-labels [1], hence we select queries having the *least average loss* [20, 1]. Since we shall iterate the process, we take the extreme approach of selecting *one query example per class*:

$$\mathcal{I} := \left\{ \arg \min_{\hat{y}_{L+i}=j} \bar{\ell}_i : j \in [N] \right\}. \quad (13)$$

Finally, we augment the support set  $S$  with the selected queries and their pseudo-labels, while at the same time removing the selected queries from  $Q$ :

$$S \leftarrow S \cup \{(x_{L+i}, \hat{y}_{L+i})\}_{i \in \mathcal{I}} \quad (14)$$

$$Q \leftarrow Q \setminus \{x_{L+i}\}_{i \in \mathcal{I}} \quad (15)$$

### 3.6. Iterative inference

Although label propagation and class balancing make predictions on the entire unlabeled query set  $Q$ , we apply cleaning to only keep one pseudo-labeled query per class, which we move from  $Q$  to the support set  $S$ . We iterate

the entire process, selecting one pseudo-labeled query per class at a time, until  $Q$  is empty and  $S$  is augmented with all pseudo-labeled queries. Assuming that the selections are correct, the idea is that treating them as truly labeled in  $S$  improves the quality of the pseudo-labels.

---

#### Algorithm 1: Iterative label propagation and cleaning (iLPC).

---

**input** : embedding  $f$   
**input** : labeled support set  $S$  with  $|S| = L$   
**input** : unlabeled query set  $Q$  with  $|Q| = M$   
**output** : augmented support set  $S$  with  $|S| = L + M$

```

1 repeat
2    $\mathcal{W} \leftarrow \text{GRAPH}(f, S, Q; \gamma, k)$   $\triangleright$  adjacency matrix (1),(2)
3    $Y \leftarrow \text{LABEL}(S)$   $\triangleright$  label matrix (3)
4    $Z \leftarrow \text{LP}(\mathcal{W}, Y; \alpha)$   $\triangleright$  label propagation (4)
5    $P \leftarrow Z_{L+1:L+M,:}$   $\triangleright$  unlabeled submatrix (5)
6    $P \leftarrow \text{POWER}(P; \tau)$   $\triangleright$  power transform (6)
7    $(\mathbf{p}, \mathbf{q}) \leftarrow \text{BALANCE}(P)$   $\triangleright$  class balance (7)
8    $P \leftarrow \text{SINKHORN}(P; \mathbf{p}, \mathbf{q})$   $\triangleright$  Sinkhorn-Knopp (9),(10)
9    $\hat{Y} \leftarrow \text{PREDICT}(P)$   $\triangleright$  pseudo-labels (11)
10   $\mathcal{I} \leftarrow \text{CLEAN}(f, S, Q, \hat{Y}, \mathbf{p})$   $\triangleright$  label cleaning (12),(13)
11   $(S, Q) \leftarrow \text{AUGMENT}(S, Q, \mathcal{I})$   $\triangleright$  augment support (14),(15)
12 until  $Q = \emptyset$   $\triangleright$  all queries are predicted
```

---

Algorithm 1 summarizes this process, called *iterative label propagation and cleaning* (iLPC). Given  $S$ ,  $Q$  and the embedding  $f$ , we construct the nearest neighbor graph represented by the normalized adjacency matrix  $\mathcal{W}$  (1),(2) and we perform label propagation on the current label matrix  $Y$  (4). Focusing on the unlabeled submatrix  $P$  of the resulting matrix  $Z$ , we perform power transform (6) and row/column normalization to balance over classes (9),(10). We predict pseudo-labels  $\hat{Y}$  from the normalized  $P$  (11), which we use along with  $S$  and  $Q$  to train a linear classifier on top of  $f$  with cross entropy loss (12) and a cyclical learning rate schedule [20]. We select one query per class with the least average loss over all epochs (13), which we move from  $Q$  to  $S$  as labeled (14),(15). With  $Q$ ,  $S$  redefined, we repeat the process until  $Q$  is empty.

At termination, all data is labeled in  $S$ . The predicted labels over the original queries are the output in the case of *transductive* inference. In *semi-supervised* classification, we use  $S$  to learn a new classifier and make predictions on new queries, as in *inductive* inference.

## 4. Experiments

### 4.1. Datasets

**miniImageNet** This is a widely used few-shot image classification dataset [60, 48]. It contains 100 randomly sampled classes from ImageNet [27]. These 100 classes are split into 64 training (base) classes, 16 validation (novel) classes and 20 test (novel) classes. Each class contains 600 examples (images). We follow the commonly used split proposed

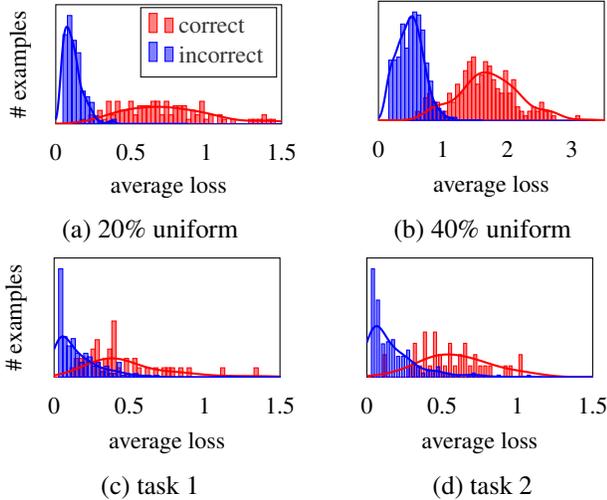


Figure 2. (a,b) Distributions of loss values (12) for correctly and incorrectly labeled examples, normalized independently. Uniform label noise: (a) 20%, (b) 40%. Pseudo-labels predicted by (11) for two different 1-shot transductive *miniImageNet* tasks (c,d).

in [48]. All images are resized to  $84 \times 84$ .

**tieredImageNet** This is also sampled from ImageNet [27] but has a hierarchical structure. Classes are partitioned into 34 categories, organized into 20 training, 6 validation and 8 test categories, containing 351, 97 and 160 classes, respectively. This ensures that training classes are semantically distinct from test classes, which is more realistic. We follow the common split of [9]. Again, all images are  $84 \times 84$ .

**CUB** This is a fine-grained classification dataset consisting of 200 classes, each corresponding to a bird species. We follow the split defined by [10, 15], with 100 training, 50 validation and 50 test classes. To compare fairly with competitors, we *use* bounding boxes on ResNet features following [59] to compare with [62] but we *do not use* bounding boxes on WRN [51] features to compare with [19].

**CIFAR-FS** This dataset is derived from CIFAR-100 [26], consisting of 100 classes with 600 examples per class. We follow the split provided by [10], with 64 training, 16 validation and 20 test classes. To compare fairly with competitors, we use the original image resolution of  $32 \times 32$  on WRN features to compare with [19] but we resize images to  $84 \times 84$  on ResNet features to compare with [62].

## 4.2. Setup

**Tasks** We consider  $N$ -way,  $K$ -shot classification tasks with  $N = 5$  randomly sampled novel classes and  $K \in \{1, 5\}$  randomly selected examples per class for support set  $S$ , that is,  $L = 5K$  examples in total. For the query set  $Q$ , we randomly sample 15 additional examples per class, that is,  $M = 75$  examples in total, which is the most common

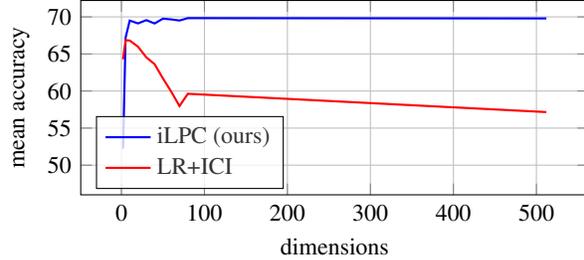


Figure 3. 1-shot transductive inference on *miniImageNet*, ablation over LR+ICI [62] pre-processing: dimension reduction by PCA.

choice in the literature [35, 31, 65].

In the semi-supervised setting, the unlabeled set  $U$  contains an additional number of randomly sampled examples per novel class. This number depends on  $K$ . We use two settings, namely 30/50 and 100/100, where the first number (30 or 100) refers to 1-shot and the second (50 or 100) to 5-shot. Again, these are the two most common choices in semi-supervised few shot learning [31, 62, 28, 49, 65].

Unless otherwise stated, we use 1000 tasks and report mean accuracy and 95% confidence interval on the test set.

**Competitors** As discussed in Appendix C, there are several flaws in experimental evaluation in the literature, like the use of different networks, training, versions of datasets, dimensionality and feature pre-processing. Fair comparison is impossible, unless one uses public code to reproduce results under exactly the same setup.

In this work, we do provide completely fair comparisons with such reproduced results of three state-of-the-art methods: LR+ICI [62], PT+MAP [19] and MCT [28]. Only [62] is published, while the other two are pre-prints.

**Networks** We use publicly available pre-trained backbone convolutional neural networks that are trained on the base-class training set. We experiment with two popular networks, namely, the residual network ResNet-12 [41] and the wide residual network WRN-28-10 [51].

In particular, to compare with [62], we use pre-trained weights of the ResNet-12 provided by [62], which we call ResNet-12A, as well as official public code<sup>1</sup> for testing. To compare with [19], we use pre-trained weights of a WRN-28-10 provided by [37]<sup>2</sup>, which are the same used by [19], as well as official public code<sup>3</sup> for testing. To compare with [28], we use official public code<sup>4</sup> to train from scratch another version of ResNet-12 used by [28], which we call ResNet-12B, as well as the same code for testing.

**Feature pre-processing** Each method uses its own feature pre-processing. LR+ICI [62] uses  $\ell_2$ -normalization and PCA

<sup>1</sup><https://github.com/Yikai-Wang/ICI-FSL>

<sup>2</sup>[https://github.com/nupurkmr9/S2M2\\_fewshot](https://github.com/nupurkmr9/S2M2_fewshot)

<sup>3</sup><https://github.com/yhu01/PT-MAP>

<sup>4</sup><https://github.com/seongmin-kye/MCT>

INFERENCE	COMPONENTS				RESNET-12		WRN-28-10		
	LP	Balance	iLC	iProb	Class	1-shot	5-shot	1-shot	5-shot
Inductive						56.30 $\pm$ 0.62	75.59 $\pm$ 0.47	60.04 $\pm$ 0.64	81.35 $\pm$ 0.43
Transductive	✓					61.09 $\pm$ 0.70	75.32 $\pm$ 0.50	74.35 $\pm$ 0.68	84.76 $\pm$ 0.42
Transductive	✓	✓				65.04 $\pm$ 0.75	76.82 $\pm$ 0.50	79.21 $\pm$ 0.72	88.00 $\pm$ 0.41
Transductive†	✓	✓	✓			<b>69.79</b> $\pm$ 0.99	79.82 $\pm$ 0.55	<b>83.43</b> $\pm$ 0.79	<b>89.10</b> $\pm$ 0.42
Transductive	✓	✓		✓		58.27 $\pm$ 0.91	74.11 $\pm$ 0.56	72.28 $\pm$ 0.90	85.88 $\pm$ 0.46
Transductive		✓	✓		✓	68.79 $\pm$ 0.96	<b>79.93</b> $\pm$ 0.56	82.01 $\pm$ 0.85	89.04 $\pm$ 0.43

Table 1. Ablation study of algorithmic components of our method iLPC on *miniImageNet*. Inductive: baseline using only support examples. LP: label propagation. Balance: class balancing (7). iLC: iterative label cleaning, without which we just output predictions (11). iProb: iterative selection of top examples per class directly as column-wise maxima of  $P$  (5) instead of iLC. Class: linear classifier used for prediction instead of LP, as in [62], with balancing still applied on output probabilities. †: default setting of iLPC.

METHOD	NETWORK	<i>miniImageNet</i>		<i>tieredImageNet</i>		CIFAR-FS		CUB	
		1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
LR+ICI [62]	ResNet-12A	66.80	79.26	80.79	87.92	73.97	84.13	88.06	92.53
LR+ICI [62]*	ResNet-12A	66.85 $\pm$ 0.92	78.89 $\pm$ 0.55	82.40 $\pm$ 0.84	88.80 $\pm$ 0.50	75.36 $\pm$ 0.97	84.57 $\pm$ 0.57	86.53 $\pm$ 0.79	92.11 $\pm$ 0.35
<b>iLPC (ours)</b>	ResNet-12A	<b>69.79</b> $\pm$ 0.99	<b>79.82</b> $\pm$ 0.55	<b>83.49</b> $\pm$ 0.88	<b>89.49</b> $\pm$ 0.38	<b>77.14</b> $\pm$ 0.95	<b>85.23</b> $\pm$ 0.55	<b>89.00</b> $\pm$ 0.70	<b>92.74</b> $\pm$ 0.35
PT+MAP [19]	WRN-28-10	82.92 $\pm$ 0.26	88.82 $\pm$ 0.13	-	-	87.69 $\pm$ 0.23	90.68 $\pm$ 0.15	91.55 $\pm$ 0.19	93.99 $\pm$ 0.10
PT+MAP [19]*	WRN-28-10	82.88 $\pm$ 0.73	88.78 $\pm$ 0.40	88.15 $\pm$ 0.71	92.32 $\pm$ 0.40	86.91 $\pm$ 0.72	90.50 $\pm$ 0.49	91.37 $\pm$ 0.61	93.93 $\pm$ 0.32
LR+ICI [62]*	WRN-28-10	80.61 $\pm$ 0.80	87.93 $\pm$ 0.44	86.79 $\pm$ 0.76	91.73 $\pm$ 0.40	84.88 $\pm$ 0.79	89.75 $\pm$ 0.48	90.18 $\pm$ 0.65	93.35 $\pm$ 0.30
<b>iLPC (ours)</b>	WRN-28-10	<b>83.43</b> $\pm$ 0.79	<b>89.10</b> $\pm$ 0.42	<b>88.50</b> $\pm$ 0.77	<b>92.99</b> $\pm$ 0.42	<b>87.57</b> $\pm$ 0.78	<b>91.26</b> $\pm$ 0.46	<b>91.53</b> $\pm$ 0.59	<b>94.27</b> $\pm$ 0.30

Table 2. Transductive inference, comparison with LR+ICI [62] and PT+MAP [19]. \*: our reproduction with official code on our datasets.

to reduce ResNet-12A to 5 dimensions. PR+MAP [19] uses element-wise power transform,  $\ell_2$ -normalization and centering of WRN-28-10 features. MCT [28] uses flattening of the output tensor of ResNet-12B rather than spatial pooling. By default, we use the same choices as [19, 28] for WRN-28-10 and ResNet-12B. For ResNet-12A however, we use  $\ell_2$ -normalization only on transductive inference and we do not use any dimensionality reduction.

**Implementation details** We use PyTorch [42] and scikit-learn [43]. Label cleaning is based on a linear classifier on top of  $f$ , initialized by imprinting the average of support features per class and then trained using (12). We use SGD with momentum 0.9 and weight decay 0.0005. Following [20], we use a cyclical schedule of 10 cycles with 20 iterations per cycle for the learning rate, which linearly decreases from a maximum  $\eta_{\max}$  to zero. For inductive (resp. semi-supervised) learning, we use logistic regression on support (resp. also pseudo-labeled) examples, learned using scikit-learn [62]. The row-wise sum  $\mathbf{p}$  (9) is fixed to 1. Appendix B includes more choices.

### 4.3. Ablation study

**Hyperparameters** Our hyperparameters include  $\gamma$  and  $k$  used in the nearest neighbor graph (1),  $\alpha$  in label propagation (4),  $\tau$  in balancing (7) and the maximum learning rate  $\eta_{\max}$  of label cleaning. We optimize them on the val-

idation set of every dataset. Common choices for  $k$  and  $\alpha$  are in [15, 20] and in [0.5, 0.8], respectively. We set  $\gamma = 3$ ,  $\tau = 3$  and  $\eta_{\max} = 0.1$ . More details and precise choices per dataset are given in Appendix A.

**Algorithmic components** Table 1 ablates our method in the presence or not of individual components, as well as using alternatives to our components. The use of queries with label propagation immediately gives a gain of transductive over inductive inference, up to 14% in 1-shot. Balancing and iterative label cleaning each bring another gain of 4-5%, especially in 1-shot. Iterative label cleaning is very important because selecting examples based on  $P$  instead is even worse than baseline transductive inference. Label propagation is also superior to using a linear classifier in most cases.

### 4.4. Label cleaning: loss distribution

To illustrate our label cleaning, we conduct two experiments, showing the distribution of the loss value (12). In the first, shown in Figure 2(a,b), we inject label noise uniformly at random to the 20% (a) and 40% (b) of 500 labeled examples. The correctly and incorrectly labeled examples have very different loss distributions. Importantly, while previous work on noisy labels [3, 20, 56] attempts to detect clean examples by an optimal threshold on the loss value, we only need few clean examples per iteration. Examples with minimal loss value are clean.

METHOD	NETWORK	<i>mini</i> IMAGENET		<i>tiered</i> IMAGENET		CIFAR-FS	
		1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
PT+MAP [19]*	WRN-28-10	83.79 $\pm$ 0.71	88.94 $\pm$ 0.33	88.87 $\pm$ 0.64	92.01 $\pm$ 0.36	87.63 $\pm$ 0.66	90.15 $\pm$ 0.46
<b>iLPC (ours)</b>	WRN-28-10	<b>86.30</b> $\pm$ 0.71	<b>89.99</b> $\pm$ 0.32	<b>90.60</b> $\pm$ 0.64	<b>92.97</b> $\pm$ 0.38	<b>88.51</b> $\pm$ 0.69	<b>91.23</b> $\pm$ 0.44

Table 3. Transductive inference, 50 queries per class. \*: our reproduction with official code on our datasets.

METHOD	<i>m</i> IN	<i>t</i> IN	CIFAR-FS	CUB
PT+MAP [19]*	89.97 $\pm$ 0.34	93.33 $\pm$ 0.34	91.30 $\pm$ 0.45	94.24 $\pm$ 0.28
<b>iLPC (ours)</b>	<b>90.51</b> $\pm$ 0.35	<b>93.61</b> $\pm$ 0.38	<b>91.59</b> $\pm$ 0.44	<b>94.75</b> $\pm$ 0.26

Table 4. Transductive 10-shot inference using WRN-28-10. *m*IN: *mini*ImageNet. *t*IN: *tiered*ImageNet. \*: our reproduction with official code on our datasets.

METHOD	PRE	<i>mini</i> ImageNet		<i>tiered</i> ImageNet	
		1-shot	5-shot	1-shot	5-shot
PT+MAP [19]*		48.57 $\pm$ 0.81	75.67 $\pm$ 0.82	49.67 $\pm$ 0.77	88.32 $\pm$ 0.50
<b>iLPC (ours)</b>		78.89 $\pm$ 0.90	86.80 $\pm$ 0.46	86.52 $\pm$ 0.47	91.07 $\pm$ 0.47
PT+MAP [19]*	✓	82.88 $\pm$ 0.73	88.78 $\pm$ 0.40	88.15 $\pm$ 0.71	92.32 $\pm$ 0.40
<b>iLPC (ours)</b>	✓	<b>83.43</b> $\pm$ 0.79	<b>89.10</b> $\pm$ 0.42	<b>88.50</b> $\pm$ 0.77	<b>92.99</b> $\pm$ 0.42

Table 5. Transductive inference, ablation over PT+MAP [19] pre-processing. PRE: power transform, normalization, centering. \*: our reproduction with official code on our datasets.

METHOD	NETOWRK	<i>mini</i> IMAGENET	
		1-shot	5-shot
MCT (instance,flip) [28]	ResNet-12B	78.55 $\pm$ 0.86	86.03 $\pm$ 0.42
MCT (no scale) [28]*	ResNet-12B	67.26 $\pm$ 0.60	<b>81.90</b> $\pm$ 0.43
<b>iLPC (ours)</b>	ResNet-12B	<b>75.58</b> $\pm$ 1.16	81.58 $\pm$ 0.50
iLPC (ours)	ResNet-12A	69.79 $\pm$ 0.99	79.82 $\pm$ 0.55

Table 6. Transductive inference, comparison with MCT [28] using ResNet-12B. \*: our reproduction with official code on our datasets, without augmentation and without scaling.

The second experiment is on two novel 1-shot transductive tasks, shown in Figure 2(c,d). We use 50 unlabeled queries per class and we predict pseudo-labels according to (11). Label cleaning is more challenging here because the two distributions are more overlapping. This is natural because predictions are more informed than uniform, even if incorrect. Still, a large proportion of clean examples have less loss than the minimal value of noisy ones.

#### 4.5. Transductive inference

Table 2 compares our iLPC with LR+ICI [62] and PT+MAP [19]. The truly fair comparison is with our reproductions, indicated by \*. Apart from the default networks, we also use WRN-28-10 with LR+ICI [62], since it is more powerful. Our iLPC is superior or on par with competitors.

The gain over [62] up to 3% on *mini*ImageNet 1-shot. The gain over PT+MAP [19] is around 0.5%.

We also experiment with 50 unlabeled queries per class, or  $M = 250$  in total. As shown in Table 3, the gain over PT+MAP [19] increases up to 3% on *mini*ImageNet 1-shot. This can be attributed to the fact that PT+MAP [19] operates on Euclidean space, while we capture the manifold structure, which manifests itself in the presence of more data. A 10-shot experiment, using more labeled support examples, is shown in Table 4. The gain is around 0.5%.

Table 5 shows that PT+MAP [19] is very sensitive to feature pre-processing, losing up to 40% without it, while our iLPC more robust, losing only up to 5%. Similarly, Figure 3 shows that LR+ICI [62] is sensitive to dimension reduction, working best at only 5 dimensions. By contrast, our iLPC is very stable and only fails at 2 dimensions.

Table 6 compares our iLPC with MCT [28]. We reproduce MCT results by training from scratch ResNet-12B using the official code and we test both methods without data augmentation (horizontal flipping) and without meta-learned scaling function. The objective is to compare the two transductive methods under the same backbone network and the same training process, which is clearly superior to ResNet-12A. Under these settings, MCT is slightly better in 5-shot but iLPC outperforms it by a large margin in 1-shot.

#### 4.6. Semi-supervised learning

As shown in Table 7, iLPC is superior to LR+ICI [62] in all settings by an even larger margin than in transductive inference, *e.g.* by nearly 5% in *mini*ImageNet 1-shot. This can be attributed to capturing the manifold structure of the data, since there is more unlabeled data in this case. Because PT+MAP [19] does not experiment with semi-supervised learning, we adapt it in the same way as ours, using the default WRN-28-10, again outperforming it.

#### 4.7. Comparison with the state of the art

Table 8 and Table 9 compare our iLPC with a larger collection of recent methods on the transductive and semi-supervised settings, respectively. Even when the network and data split appears to be the same, we acknowledge that our results are not directly comparable with any method other than our reproductions. As discussed in Appendix C, this is due to the very diverse choices made in the bibliography, *e.g.* versions of network, training settings, versions of datasets,

METHOD	NETWORK	SPLIT	miniIMAGENET		tieredIMAGENET		CIFAR-FS		CUB	
			1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
LR+ICI [62]	ResNet-12A	30/50	69.66	80.11	84.01	89.00	76.51	84.32	89.58	92.48
LR+ICI [62]*	ResNet-12A	30/50	67.64 $\pm$ 0.95	80.00 $\pm$ 0.55	83.20 $\pm$ 0.89	89.15 $\pm$ 0.48	76.57 $\pm$ 0.96	84.11 $\pm$ 0.60	88.36 $\pm$ 0.71	-
<b>iLPC (ours)</b>	ResNet-12A	30/50	<b>72.36</b> $\pm$ 0.93	<b>81.23</b> $\pm$ 0.51	<b>84.34</b> $\pm$ 0.77	<b>89.66</b> $\pm$ 0.49	78.44 $\pm$ 0.89	<b>85.69</b> $\pm$ 0.50	<b>89.44</b> $\pm$ 0.62	-
LR+ICI [62]*	WRN-28-10	30/50	81.77 $\pm$ 0.81	88.53 $\pm$ 0.44	87.96 $\pm$ 0.75	92.11 $\pm$ 0.41	86.41 $\pm$ 0.77	89.90 $\pm$ 0.53	91.29 $\pm$ 0.59	-
PT+MAP [19]†	WRN-28-10	30/50	83.47 $\pm$ 0.68	89.15 $\pm$ 0.38	88.67 $\pm$ 0.68	92.44 $\pm$ 0.38	87.71 $\pm$ 0.67	90.38 $\pm$ 0.45	92.05 $\pm$ 0.52	-
<b>iLPC (ours)</b>	WRN-28-10	30/50	<b>84.00</b> $\pm$ 0.74	<b>89.68</b> $\pm$ 0.37	<b>88.87</b> $\pm$ 0.73	<b>92.66</b> $\pm$ 0.40	<b>87.78</b> $\pm$ 0.69	<b>90.58</b> $\pm$ 0.50	<b>92.35</b> $\pm$ 0.52	-

Table 7. Semi-supervised few-shot learning, comparison with LR+ICI [62] and PT+MAP [19]. \*: our reproduction with official code on our datasets. †: our adaptation to semi-supervised, based on official code. CUB 5-shot omitted: no class has the required 70 examples.

METHOD	NETWORK	miniIMAGENET		tieredIMAGENET		CIFAR-FS		CUB	
		1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
LR+ICI [62]*	ResNet-12A	66.85 $\pm$ 0.92	78.89 $\pm$ 0.55	82.40 $\pm$ 0.84	88.80 $\pm$ 0.50	75.36 $\pm$ 0.97	84.57 $\pm$ 0.57	86.53 $\pm$ 0.79	92.11 $\pm$ 0.35
CAN+Top- $k$ [16]	ResNet-12	67.19 $\pm$ 0.55	80.64 $\pm$ 0.35	73.21 $\pm$ 0.58	84.93 $\pm$ 0.38	-	-	-	-
DPGN [64]	ResNet-12	67.77 $\pm$ 0.32	84.60 $\pm$ 0.43	72.45 $\pm$ 0.51	87.24 $\pm$ 0.39	77.90 $\pm$ 0.50	90.20 $\pm$ 0.40	75.71 $\pm$ 0.47	91.48 $\pm$ 0.33
MCT (instance) [28]	ResNet-12B	78.55 $\pm$ 0.86	86.03 $\pm$ 0.42	82.32 $\pm$ 0.81	87.36 $\pm$ 0.50	85.61 $\pm$ 0.69	90.03 $\pm$ 0.46	-	-
PT+MAP [19]*	WRN-28-10	82.88 $\pm$ 0.73	88.78 $\pm$ 0.40	88.15 $\pm$ 0.71	92.32 $\pm$ 0.40	86.91 $\pm$ 0.72	90.50 $\pm$ 0.49	91.37 $\pm$ 0.61	93.93 $\pm$ 0.32
Fine-tuning [11]	WRN-28-10	65.73 $\pm$ 0.68	78.40 $\pm$ 0.52	73.34 $\pm$ 0.71	85.50 $\pm$ 0.50	76.58 $\pm$ 0.68	85.79 $\pm$ 0.50	-	-
EP [50]	WRN-28-10	70.74 $\pm$ 0.85	84.34 $\pm$ 0.53	78.50 $\pm$ 0.91	88.36 $\pm$ 0.57	-	-	-	-
SIB [17]†	WRN-28-10	70.00 $\pm$ 0.60	79.20 $\pm$ 0.40	72.90	82.80	80.00 $\pm$ 0.60	85.3 $\pm$ 0.40	-	-
SIB+E <sup>3</sup> BM [36]	WRN-28-10	71.40 $\pm$ 0.50	81.20 $\pm$ 0.40	75.60 $\pm$ 0.6	84.30 $\pm$ 0.4	-	-	-	-
<b>iLPC (ours)</b>	WRN-28-10	<b>83.43</b> $\pm$ 0.79	<b>89.10</b> $\pm$ 0.42	<b>88.50</b> $\pm$ 0.77	<b>92.99</b> $\pm$ 0.42	<b>87.57</b> $\pm$ 0.78	<b>91.26</b> $\pm$ 0.46	<b>91.53</b> $\pm$ 0.59	<b>94.27</b> $\pm$ 0.30

Table 8. Transductive inference state of the art. \*: our reproduction with official code on our datasets. †: *tieredImageNet* as reported by [36].

METHOD	NETWORK	SPLIT	miniIMAGENET	
			1-shot	5-shot
LST [31]	ResNet-12	30/50	70.10 $\pm$ 1.90	78.70 $\pm$ 0.80
LR+ICI [62]	ResNet-12A	30/50	69.66	80.11
MCT (instance) [28]	ResNet-12B	30/50	73.80 $\pm$ 0.70	84.40 $\pm$ 0.50
$k$ -means [49]†	WRN-28-10	100/100	52.35 $\pm$ 0.89	67.67 $\pm$ 0.65
TransMatch [65]	WRN-28-10	100/100	63.02 $\pm$ 1.07	81.06 $\pm$ 0.59
<b>iLPC (ours)</b>	WRN-28-10	100/100	<b>87.62</b> $\pm$ 0.67	<b>90.51</b> $\pm$ 0.36

Table 9. Semi-supervised few-shot learning state of the art. †: as reported by [65].

or pre-processing. For instance, ResNet-12 is different than either ResNet-12A or ResNet-12B.

For this reason, we focus on the best result by each method, including ours. Necessarily, methods experimenting with WRN-28-10 have an advantage. Still, at least among those, iLPC performs best by a large margin in both settings, with the closest second best being PT+MAP [19].

## 5. Conclusion

Our solution is conceptually simple and combines in a unique way ideas that have been successful in problems related to our task at hand. *Label propagation* exploits the

manifold structure of the data, which becomes important in the presence of more data, while still being competitive otherwise. *Class balancing* provides a strong hint in correcting predictions when certain classes dominate. *Label cleaning*, originally introduced for learning with noisy labels, is also very successful in cleaning predicted pseudo-labels. *Iterative reuse* of few pseudo-labels as true labels bypasses the difficulty of single-shot detection of clean examples.

Importantly, reasonable baselines, like predicting pseudo-labels by a classifier or iteratively re-using pseudo-labels without cleaning, fail completely. When compared under fair settings, our iLPC outperforms or is on par with state-of-the-art methods on a range of networks and datasets. It is also significantly more robust against feature pre-processing and other tricks on which other methods rely. Finally, there is no data augmentation in our work. Augmentation should naturally bring complementary improvement.

## References

- [1] Paul Albert, Diego Ortego, Eric Arazo, Noel E O’Connor, and Kevin McGuinness. Relab: Reliable label bootstrapping for semi-supervised learning. *arXiv preprint arXiv:2007.11866*, 2020. 4
- [2] A. Antoniou, A. Storkey, and H. Edwards. Data augmentation generative adversarial networks. *arXiv preprint arXiv:1711.04340*, abs/1711.04340, 2017. 2

- [3] Eric Arazo, Diego Ortego, Paul Albert, Noel E. O’Connor, and Kevin McGuinness. Unsupervised label noise modeling and loss correction. In *ICML*, 2019. 1, 4, 6
- [4] Gedas Bertasius, Lorenzo Torresani, Stella X. Yu, and Jianbo Shi. Convolutional random walk networks for semantic image segmentation. In *CVPR*, July 2017. 2
- [5] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. In *NeurIPS*, 2019. 1, 2
- [6] Luca Bertinetto, Joao F Henriques, Philip HS Torr, and Andrea Vedaldi. Meta-learning with differentiable closed-form solvers. *arXiv preprint arXiv:1805.08136*, 2018. 2
- [7] Luca Bertinetto, João F Henriques, Jack Valmadre, Philip Torr, and Andrea Vedaldi. Learning feed-forward one-shot learners. In *NeurIPS*, 2016. 2
- [8] Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien. *Semi-Supervised Learning*. The MIT Press. 1
- [9] Da Chen, Yuefeng Chen, Yuhong Li, Feng Mao, Yuan He, and Hui Xue. Self-supervised learning for few-shot image classification. *arXiv preprint arXiv:1911.06045*, 2019. 5
- [10] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Wang, and Jia-Bin Huang. A closer look at few-shot classification. In *ICLR*, 2019. 1, 2, 5, 11, 12
- [11] Guneet S Dhillon, Pratik Chaudhari, Avinash Ravichandran, and Stefano Soatto. A baseline for few-shot image classification. *arXiv preprint arXiv:1909.02729*, 2019. 8, 11
- [12] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017. 1, 2
- [13] Spyros Gidaris and Nikos Komodakis. Dynamic few-shot visual learning without forgetting. In *CVPR*, 2018. 2, 3
- [14] Bharath Hariharan and Ross Girshick. Low-shot visual recognition by shrinking and hallucinating features. In *CVPR*, 2017. 2
- [15] Nathan Hilliard, Lawrence Phillips, Scott Howland, Artëm Yankov, Courtney D Corley, and Nathan O Hodas. Few-shot learning with metric-agnostic conditional embeddings. *arXiv preprint arXiv:1802.04376*, 2018. 5
- [16] Ruibing Hou, Hong Chang, MA Bingpeng, Shiguang Shan, and Xilin Chen. Cross attention network for few-shot classification. In *NeurIPS*, 2019. 2, 8
- [17] Shell Xu Hu, Pablo Garcia Moreno, Yang Xiao, Xi Shen, Guillaume Obozinski, Neil Lawrence, and Andreas Damianou. Empirical bayes transductive meta-learning with synthetic gradients. In *ICLR*, 2019. 8
- [18] Yuqing Hu, Vincent Gripon, and S. Pateux. Exploiting unsupervised inputs for accurate few-shot classification. *arXiv preprint arXiv:2001.09849*, 2020. 2
- [19] Yuqing Hu, Vincent Gripon, and Stéphane Pateux. Leveraging the feature distribution in transfer-based few-shot learning. *arXiv preprint arXiv:2006.03806*, 2020. 2, 3, 5, 6, 7, 8, 11
- [20] Jinchu Huang, Lie Qu, Rongfei Jia, and Binqiang Zhao. O2U-Net: A simple noisy label detection approach for deep neural networks. *ICCV*, 2019. 1, 2, 4, 6
- [21] Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. Label propagation for deep semi-supervised learning. In *CVPR*, 2019. 1, 3, 4, 11
- [22] Sungeon Kim, Dongwon Kim, Minsu Cho, and Suha Kwak. Proxy anchor loss for deep metric learning. In *CVPR*, 2020. 2
- [23] Philip A Knight. The Sinkhorn-Knopp algorithm: convergence and applications. *SIAM Journal on Matrix Analysis and Applications*, 2008. 1, 2, 3
- [24] Gregory Koch, Richard Zemel, and Ruslan Salakhutdinov. Siamese neural networks for one-shot image recognition. In *ICML workshop*, 2015. 2
- [25] Deguang Kong and Chris Ding. Maximum consistency preferential random walks. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2012. 1, 2
- [26] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 5
- [27] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep convolutional neural networks. In *NeurIPS*, 2012. 4, 5
- [28] Seong Min Kye, Hae Beom Lee, Hoirin Kim, and Sung Ju Hwang. Meta-learned confidence for few-shot learning. *arXiv preprint arXiv:2002.12017*, 2020. 1, 2, 5, 6, 7, 8
- [29] Dong-Hyun Lee. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. 2013. 1, 2, 4
- [30] Kwonjoon Lee, Subhansu Maji, Avinash Ravichandran, and Stefano Soatto. Meta-learning with differentiable convex optimization. In *CVPR*, 2019. 2, 11, 12
- [31] Xinzhe Li, Qianru Sun, Yaoyao Liu, Qin Zhou, Shibao Zheng, Tat-Seng Chua, and Bernt Schiele. Learning to self-train for semi-supervised few-shot classification. In *NeurIPS*, 2019. 2, 5, 8, 11
- [32] Moshe Lichtenstein, Prasanna Sattigeri, Rogerio Feris, Raja Giryes, and Leonid Karlinsky. TAFSSL: Task-adaptive feature sub-space learning for few-shot classification, 2020. 2
- [33] Bin Liu, Yue Cao, Yutong Lin, Qi Li, Zheng Zhang, Mingsheng Long, and Han Hu. Negative margin matters: Understanding margin in few-shot classification. In *ECCV*, 2020. 2
- [34] Ming-Yu Liu, Xun Huang, Arun Mallya, Tero Karras, Timo Aila, Jaakko Lehtinen, and Jan Kautz. Few-shot unsupervised image-to-image translation. In *CVPR*, 2019. 2
- [35] Yanbin Liu, Juho Lee, Minseop Park, Saehoon Kim, Eunho Yang, Sung Ju Hwang, and Yi Yang. Learning to propagate labels: Transductive propagation network for few-shot learning. *arXiv preprint arXiv:1805.10002*, 2018. 1, 2, 5, 11
- [36] Yaoyao Liu, Bernt Schiele, and Qianru Sun. An ensemble of epoch-wise empirical bayes for few-shot learning. In *ECCV*, 2020. 8
- [37] Puneet Mangla, Nupur Kumari, Abhishek Sinha, Mayank Singh, Balaji Krishnamurthy, and Vineeth N Balasubramanian. Charting the right manifold: Manifold mixup for few-shot learning. In *WACV*, 2020. 2, 5, 12
- [38] Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. A simple neural attentive meta-learner. *arXiv preprint arXiv:1707.03141*, 2017. 1, 2
- [39] Tsendsuren Munkhdalai and Hong Yu. Meta networks. In *ICML*, 2017. 2

- [40] Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018. [2](#)
- [41] Boris Oreshkin, Pau Rodríguez López, and Alexandre Lacoste. Tadam: Task dependent adaptive metric for improved few-shot learning. In *NeurIPS*, 2018. [3](#), [5](#)
- [42] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. [6](#)
- [43] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *JMLR*, 2011. [6](#)
- [44] Hang Qi, Matthew Brown, and David G Lowe. Low-shot learning with imprinted weights. In *CVPR*, 2018. [2](#), [3](#)
- [45] L. Qiao, Y. Shi, Jia Li, Yaowei Wang, Tiejun Huang, and Yonghong Tian. Transductive episodic-wise adaptive metric for few-shot learning. *ICCV*, 2019. [2](#)
- [46] Siyuan Qiao, Chenxi Liu, Wei Shen, and Alan L Yuille. Few-shot image recognition by predicting parameters from activations. In *CVPR*, 2018. [2](#)
- [47] Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. Meta-learning with implicit gradients. In *NeurIPS*, 2019. [2](#)
- [48] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. 2016. [2](#), [4](#), [5](#)
- [49] Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B Tenenbaum, Hugo Larochelle, and Richard S Zemel. Meta-learning for semi-supervised few-shot classification. *arXiv preprint arXiv:1803.00676*, 2018. [1](#), [2](#), [5](#), [8](#), [11](#)
- [50] Pau Rodríguez, Issam Laradji, Alexandre Drouin, and Alexandre Lacoste. Embedding propagation: Smoother manifold for few-shot classification. *ECCV*, 2020. [1](#), [2](#), [8](#), [11](#)
- [51] Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. *arXiv preprint arXiv:1807.05960*, 2018. [5](#)
- [52] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-learning with memory-augmented neural networks. In *ICML*, 2016. [2](#)
- [53] H Scudder. Probability of error of some adaptive pattern-recognition machines. *IEEE Transactions on Information Theory*, 1965. [1](#)
- [54] Christian Simon, Piotr Koniusz, Richard Nock, and Mehrtash Harandi. Adaptive subspaces for few-shot learning. In *CVPR*, 2020. [2](#)
- [55] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *NeurIPS*, 2017. [1](#), [2](#)
- [56] Jiaming Song, Lunjia Hu, Yann Dauphin, M. Auli, and Tengyu Ma. Robust and on-the-fly dataset denoising for image classification. *arXiv preprint arXiv:2003.10647*, 2020. [1](#), [4](#), [6](#)
- [57] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. Learning to compare: Relation network for few-shot learning. In *CVPR*, 2018. [2](#)
- [58] Yonglong Tian, Yue Wang, Dilip Krishnan, Joshua B Tenenbaum, and Phillip Isola. Rethinking few-shot image classification: a good embedding is all you need? *arXiv preprint arXiv:2003.11539*, 2020. [1](#), [2](#)
- [59] Eleni Triantafillou, Richard Zemel, and Raquel Urtasun. Few-shot learning through an information retrieval lens. In *NeurIPS*, 2017. [5](#)
- [60] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *NeurIPS*, 2016. [1](#), [2](#), [4](#)
- [61] Xun Wang, Xintong Han, Weilin Huang, Dengke Dong, and Matthew R Scott. Multi-similarity loss with general pair weighting for deep metric learning. In *CVPR*, 2019. [2](#)
- [62] Yikai Wang, C. Xu, Chen Liu, Liyong Zhang, and Yanwei Fu. Instance credibility inference for few-shot learning. *CVPR*, 2020. [1](#), [2](#), [5](#), [6](#), [7](#), [8](#), [11](#)
- [63] Yu-Xiong Wang, Ross Girshick, Martial Hebert, and Bharath Hariharan. Low-shot learning from imaginary data. In *CVPR*, 2018. [2](#)
- [64] Ling Yang, Liangliang Li, Zilun Zhang, Xinyu Zhou, Erjin Zhou, and Yu Liu. Dpgn: Distribution propagation graph network for few-shot learning. In *CVPR*, 2020. [8](#), [11](#)
- [65] Zhongjie Yu, L. Chen, Zhongwei Cheng, and Jiebo Luo. Transmatch: A transfer-learning scheme for semi-supervised few-shot learning. *CVPR*, 2020. [1](#), [2](#), [5](#), [8](#), [11](#)
- [66] Dengyong Zhou, Olivier Bousquet, Thomas Navin Lal, Jason Weston, and Bernhard Schölkopf. Learning with local and global consistency. In *NIPS*, 2003. [1](#), [2](#), [3](#)
- [67] Xiaojin Zhu and Zoubin Ghahramani. Learning from labeled and unlabeled data with label propagation. Technical report, 2002. [1](#)
- [68] Luisa Zintgraf, Kyriacos Shiarli, Vitaly Kurin, Katja Hofmann, and Shimon Whiteson. Fast context adaptation via meta-learning. [2](#)

## A. Hyperparameters

Table 10 shows the best hyperparameters  $k$  (1) and  $\alpha$  (4) for every dataset, network and number of support examples per class  $K \in \{1, 5\}$ . The hyperparameters are optimized on the validation set separately for each experiment. We carried out experiments in the transductive setting for  $k \in \{5, 8, 10, 15, 20, 25, 30, 40, 50, 60\}$  and  $\alpha \in \{0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9\}$  and select the combination resulting in the best mean validation accuracy. In the semi-supervised setting, we use the same optimal values.

## B. Confidence weights

The Sinkhorn-Knopp algorithm iteratively normalizes a  $M \times N$  positive matrix  $P$  to a row-wise sum  $\mathbf{p} \in \mathbb{R}^M$  and column-wise sum  $\mathbf{q} \in \mathbb{R}^N$ . We experiment with two ways of setting the value of  $\mathbf{p}$ :

1. **Uniform.** Interpreting the  $i$ -th row of  $P$  as a class probability distribution for the  $i$ -th query, it should be normalized to one, such that  $p_i = 1$  uniformly.
2. **Entropy.** Because we do not have the same confidence for each prediction, we use the entropy of the predicted class probability distribution of each example to quantify its uncertainty. Following [21], we associate to each example  $x_{L+i}$  for  $i \in [M]$  a weight

$$\omega_i := 1 - \frac{H(\hat{\mathbf{z}}_i)}{\log(N)}, \quad (16)$$

where  $N$  is the number of classes and  $\hat{\mathbf{z}}_i$  is the  $\ell_1$ -normalized  $i$ -th row of  $Z$  (4), that is,  $\hat{z}_{ij} := z_{ij} / \sum_{k=1}^N z_{ik}$ . We then set the confidence weights  $p_i = \omega_i$ . Note that  $\omega_i$  takes values in  $[0, 1]$  because  $\log(N)$  is the maximum possible entropy.

Given  $\mathbf{p}$  and assuming balanced classes,  $\mathbf{q}$  is defined by (7), that is,  $q_j = \frac{1}{N} \sum_{i=1}^M p_i$  for  $j \in [N]$ . In the special case of  $p_i = 1$ , this simplifies to  $q_j = \frac{M}{N}$ .

Table 11 compares the two approaches. Even though using non-uniform confidence weights is a reasonable choice, uniform weights are superior in all settings. This can be attributed to the fact that examples with small weight tend to be ignored in the balancing process, hence their class distribution and consequently their predictions are determined mostly by other examples with large weight. For this reason, examples with small weight may get more incorrect predictions in the case of entropy.

## C. Flaws in evaluation

Throughout our investigations we observed that comparisons are commonly published that are not under the same settings. In this section we highlight such problems.

PARAM	mIN		tIN		CFS		CUB	
$K$ (shot)	1	5	1	5	1	5	1	5
ResNet-12A								
$k$ (1)	15	25	15	60	15	15	10	8
$\alpha$ (4)	0.8	0.4	0.5	0.8	0.8	0.4	0.6	0.6
ResNet-12B								
$k$ (1)	15	15	-	-	-	-	-	-
$\alpha$ (4)	0.9	0.9	-	-	-	-	-	-
WRN-28-10								
$k$ (1)	20	30	20	20	20	25	25	25
$\alpha$ (4)	0.8	0.2	0.8	0.8	0.4	0.5	0.2	0.5

Table 10. Selected hyperparameters. mIN: *miniImageNet*. tIN: *tieredImageNet*. CFS: CIFAR-FS.

METHOD	RESNET-12		WRN-28-10	
	1-shot	5-shot	1-shot	5-shot
uniform	<b>69.79</b> $\pm 0.99$	<b>79.82</b> $\pm 0.55$	<b>83.43</b> $\pm 0.79$	<b>89.10</b> $\pm 0.42$
entropy	66.94 $\pm 1.01$	78.34 $\pm 0.58$	81.05 $\pm 0.90$	88.43 $\pm 0.44$

Table 11. Comparison between ways of setting confidence weights  $\mathbf{p}$ ; transductive inference on *miniImageNet*. Uniform:  $p_i = 1$ . Entropy:  $p_i = \omega_i$  (16).

1. In multiple works such as [50, 19, 64, 11], comparisons between state-of-the-art methods are made without explicitly differentiating between inductive and transductive methods. This is unfair since transductive methods perform better by leveraging query data.
2. Comparisons use different networks without mentioning so. For example, Table 1 of [65] does not indicate what network each method uses. [65] uses WRN-28-10, while [35] uses a 4-layer convolutional network.
3. In the semi-supervised setting, comparisons use different numbers of unlabelled data without mentioning so. In Table 4 of [50] for example, [50] uses 100 unlabelled examples while [31] uses 30 for 1-shot and 50 for 5-shot, [49] and [35] use 20 for 1-shot and 20 for 5-shot. In Table 1 of [62], the best model of [62] uses an 80/80 split for 1/5-shot while other methods such as [31] use a 30/50 split. In Table 1 of [65], [65] uses 100 or 200 unlabelled examples while [49, 35] use 20/20 split for 1/5-shot.
4. Some methods use different dataset settings when comparing with other methods without explicitly stating so. In Table 1 of [62] for instance, [62] uses the bounding box provided for CUB while other methods such as [10, 30] do not.
5. Comparisons using the same network is made but this

network has been trained using a different training regimes. Unless the novelty of the work lies in the training regime, this is unfair. As shown in [37], a better training regime can increase the performance significantly.

6. There are several different variants of the benchmark datasets, coming from different sources. The two most common variants are [10], which uses original image files, and [30], which uses pre-processed tensors stored in `pkl` files. Testing a network on a different variant than the one it was trained on may result in performance drops as large as 5%.

We believe that highlighting these evaluation flaws will help researchers avoid making such mistakes and move towards a fairer evaluation. We encourage the community to compare different methods against the same settings and if otherwise, state clearly the differences. As a contribution towards a fairer evaluation, we intend to make our code publicly available along with the pre-trained networks used in this work.