

# FEWMATCH: DYNAMIC PROTOTYPE REFINEMENT FOR SEMI-SUPERVISED FEW-SHOT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

1 Semi-Supervised Few-shot Learning (SS-FSL) investigates the benefit of incorpo-  
 2 rating unlabelled data in few-shot settings. Recent work has relied on the popular  
 3 Semi-Supervised Learning (SSL) concept of iterative pseudo-labelling, yet often  
 4 yield models that are susceptible to error propagation and are sensitive to initial-  
 5 isation. Alternative work utilises the concept of consistency regularisation (CR),  
 6 a popular SSL state of the art technique where a student model is trained to con-  
 7 sistently agree with teacher predictions under different input perturbations, with-  
 8 out pseudo-label requirements. However, applications of CR to the SS-FSL set-  
 9 up struggle to outperform pseudo-labelling approaches; limited available training  
 10 data yields unreliable early stage predictions and requires fast convergence that is  
 11 not amenable for, typically slower to converge, CR approaches.

12 In this paper, we introduce a prototype-based approach for SS-FSL that exploits  
 13 model consistency in a robust manner. Our Dynamic Prototype Refinement (DPR)  
 14 approach is a novel training paradigm for few-shot model adaptation to new un-  
 15 seen classes, combining concepts from metric and meta-gradient based FSL meth-  
 16 ods. New class prototypes are alternatively refined 1) explicitly, using labelled  
 17 and unlabelled data with high confidence class predictions and 2) implicitly, by  
 18 model fine-tuning using a data selective CR loss. DPR affords CR convergence,  
 19 with the explicit refinement providing an increasingly stronger initialisation. We  
 20 demonstrate method efficacy and report extensive experiments on two competitive  
 21 benchmarks; *miniImageNet* and *tieredImageNet*. The ability to effectively utilise  
 22 and combine information from both labelled base-class and auxiliary unlabelled  
 23 novel-class data results in significant accuracy improvements.

## 24 1 INTRODUCTION

25 Few-Shot Learning (FSL) has recently made steady progress in the directions of both metric learn-  
 26 ing (Vinyals et al., 2016; Snell et al., 2017; Sung et al., 2018; Qiao et al., 2018), where class rep-  
 27 resentative features are learned to optimise intra- and inter-class distances, and meta-gradient ap-  
 28 proaches (Finn et al., 2017; Antoniou et al., 2018; Rajeswaran et al., 2019; Rusu et al., 2018) that  
 29 focus on optimising model convergence with very few training examples. Despite recent progress,  
 30 FSL performance remains limited by the small available data from which to learn from. One promis-  
 31 ing direction for progress involves introducing unlabelled training examples, allowing for expansion  
 32 of training set variability without increasing data labelling costs. Recent work has shown that this  
 33 strategy, referred to as semi-supervised few-shot learning (SS-FSL), can substantially boost FSL  
 34 performance in classification settings (Ren et al., 2018; Li et al., 2019b). These works take advan-  
 35 tage of semi-supervised learning (SSL) techniques, that historically focus on large data regimes, to  
 36 leverage information from additional unlabelled samples in combination with state of the art FSL  
 37 approaches. State of the art SS-FSL (Liu et al., 2018; Li et al., 2019b) relies on popular SSL tech-  
 38 niques of label propagation (Iscen et al., 2019), propagating label predictions to unlabelled data,  
 39 and self-training (Lee, 2013) that repeatedly labels unlabelled data, based on confidence scores, and  
 40 then retrains with this additional pseudo-annotated data. An important drawback of such strategies  
 41 is their reliance on *iteratively extending the training set using pseudo-label predictions*. Building  
 42 on pseudo-label decisions can propagate and amplify errors during training, yielding brittle methods  
 43 sensitive to model initialisation and noisy data. This problem is exacerbated in few-shot scenarios,  
 44 where available labelled data is highly limited and pseudo labels therefore have larger influence.

45 In light of these limitations, alternative work has explored the use of self-supervision techniques (Gi-  
 46 daris et al., 2019; Yu et al., 2020) to leverage information from unlabelled data. This involves the  
 47 introduction of auxiliary tasks and artificial labels (e.g. image rotation prediction, jigsaw puzzles) or  
 48 training process regularisation via a low density assumption (regularisation of consistency). These  
 49 techniques are able to exploit unlabelled data without introducing reliance on pseudo-labels. No-  
 50 tably, Consistency Regularisation (CR) (Tavainen & Valpola, 2017; Laine & Aila, 2016; Berthelot  
 51 et al., 2019b;a) regularises models to output consistent predictions under varying input perturbations.  
 52 This constitutes a state of the art SSL strategy, typically outperforming pseudo-label approaches  
 53 in large data regimes. In SS-FSL settings, however, self-supervision methods struggle to outper-  
 54 form pseudo-labelling approaches and fail to fully exploit the benefits, especially in the lowest data  
 55 regimes. This commonly results in more modest improvements from the use of unlabelled data.

56 In this work we propose a strategy that enables harnessing of the aforementioned strong performance  
 57 of CR in standard SSL, for the SS-FSL setting. We hypothesise that CR currently fails in the SS-  
 58 FSL scenario due to 1) slow convergence of CR techniques (Berthelot et al., 2019a), which is in  
 59 conflict with FSL fast convergence requirements to alleviate overfitting risks and 2) poor reliability  
 60 of model predictions in early stages, when training with limited data. We introduce a novel method  
 61 specifically designed to address these issues and demonstrate empirically that our strategy allows  
 62 successful exploitation of CR in the SS-FSL setting, outperforming state of the art techniques.

63 Our formulation exploits the popular concept of prototypes (Snell et al., 2017), commonly used in  
 64 metric-learning based FSL. Prototypes  $\mathcal{P}=\{p_1, p_2, \dots, p_{c_b}\}$  are learned global feature representa-  
 65 tions, each describing a particular class to recognise. Class prototypes are typically defined as the  
 66 average feature representation of the labelled set. They are learned using a set of base classes such  
 67 that the distances between input samples, of a given class, and the respective class prototype is min-  
 68 imised (else maximised). Our approach builds on the imprinted weights model (Qi et al., 2018), a  
 69 variant of prototypical networks, that use a simple normalisation trick to learn prototypes as clas-  
 70 sifier weights in an end-to-end manner (c.f. commonly used episode training Snell et al. (2017)).  
 71 Our proposed two-stage approach comprises pre-training on base classes, followed by our key in-  
 72 novation, a Dynamic Prototype Refinement (DPR) on novel classes. Using the imprinted weights  
 73 (IW) model we are able to seamlessly introduce an auxiliary CR loss in our base training process.  
 74 This allows to leverage unlabelled data from base classes and learn a robust initialisation for our  
 75 DPR stage. Our novel DPR method exploits unlabelled samples *from novel classes* towards learning  
 76 prototypes of higher quality. Our approach alternates between explicit updating of prototypes using  
 77 selected unlabelled samples yielding the most confident predictions (i.e. nearest to their assigned  
 78 class prototype), and implicit fine-tuning of the model with CR on a second selection of unlabelled  
 79 samples. We will show that alternating between typically smaller, more conservative updates (im-  
 80 plicit refinement) and larger, often times more disruptive feature averaging based updates (explicit  
 81 refinement), results in faster convergence for CR and often large performance gains, whilst at the  
 82 same time affording robustness to pseudo-labelling errors. We highlight that in contrast to pseudo-  
 83 labelling based approaches (Liu et al., 2018; Li et al., 2019b); *estimated labels are not propagated*  
 84 and are used exclusively to strengthen prototype initialisation, prior to fine-tuning. It is this property  
 85 that enables recovery from potential erroneous labels and the prevention of gradual drift.

86 In summary, **our contributions** are three-fold: **(a)** We present “Fewmatch”; a novel semi-supervised  
 87 few-shot learning approach that robustly exploits the concept of consistency regularisation, allevi-  
 88 ating the requirement of iterative pseudo-labelling and consistently outperforming approaches that  
 89 alternatively do possess such a requirement. **(b)** We introduce a dynamic prototype refinement pro-  
 90 cess, a novel training paradigm designed to harness the power of CR in few-shot regimes through the  
 91 use of both implicit and explicit prototype refinement steps. **(c)** Extensive experiments demonstrate  
 92 that we achieve state of the art performance on two standard benchmarks, outperforming prior CR  
 93 and self-supervised methods with significant accuracy gains. Further to this, we additionally explore  
 94 more realistic few-shot test conditions in terms of inequalities relating to unlabelled data availability.

## 95 2 RELATED WORK

96 **Semi-supervised learning** Existing SSL methods generally fall into two categories: (1) Pseudo-  
 97 labelling and (2) Consistency Regularisation. Techniques in the former category iteratively assign  
 98 pseudo labels to the unlabelled samples such that they can then be used with a supervised loss.

99 These include directly using the network class prediction (Lee, 2013) and graph-based label prop-  
 100 agation (Isken et al., 2019). A number of SSL works build on the second category of Consistency  
 101 Regularisation (Sajjadi et al., 2016; Laine & Aila, 2016; Tarvainen & Valpola, 2017), and have  
 102 achieved impressive results. The crux of the idea of CR is to encourage invariant (stable) predic-  
 103 tions for a given sample under different perturbations towards improving class decision boundaries.  
 104 CR ideas were first explored in (Sajjadi et al., 2016; Laine & Aila, 2016) and extended in (Tar-  
 105 vainen & Valpola, 2017) where the authors propose a mean teacher framework to perform CR be-  
 106 tween a student and teacher model in a learning paradigm involving models that share the same  
 107 architecture and teacher parameters are updated as an exponential moving average of the student  
 108 weights. Several works such as ICT (Verma et al., 2019), Mixmatch (Berthelot et al., 2019b) and  
 109 Remixmatch (Berthelot et al., 2019a) have then enabled sample perturbations by creating variants of  
 110 mixup samples (Zhang et al., 2017) that can then be further perturbed. Encouraged by the benefits  
 111 that result from representing class information using prototypes (Snell et al., 2017; Qi et al., 2018),  
 112 we take an alternative approach to CR in the context of SS-FSL and influence model prediction by  
 113 considering a measure of distance between unlabelled data and class prototypes.

114 **Few-Shot Learning** Existing FSL approaches can be broadly divided into two categories (1) Metric  
 115 based (Snell et al., 2017; Vinyals et al., 2016; Qi et al., 2018) and (2) Gradient based (Finn et al.,  
 116 2017; Antoniou et al., 2018; Rajeswaran et al., 2019; Rusu et al., 2018). Metric based methods aim  
 117 to learn global class feature representations (i.e. prototypes) whose distance is minimal to samples  
 118 of the same class. In this paper, we take advantage of one such method; Imprinted weights (Qi  
 119 et al., 2018) in order to provide per class prototypes. One of the main advantages of this approach  
 120 is that it does not require the standard, restrictive episode training strategy. Episode training is  
 121 framed as a sequence of artificially designed FSL tasks with fixed category and labelled sample  
 122 counts and also imposes an identical test time set-up. This in theory affords us greater flexibility  
 123 with the learning problem definition, allows for consideration of more practical problem setups, and  
 124 for easier combination with techniques from other fields such as integration of auxiliary losses.

125 **Semi-Supervised Few-Shot Learning (SS-FSL)** Existing SS-FSL approaches are based on the  
 126 pseudo-labelling strategy that was discussed in the context of SSL. Ren et al. (2018) propose mask  
 127 soft K-means, based on the metric learning approach, ProtoNets (Snell et al., 2017). The authors  
 128 use a soft K-means and iteratively assign pseudo labels to tune prototypes. More recently (Liu et al.,  
 129 2018) propose a Transductive Propagation Network (TPN) that propagates labels from unlabelled  
 130 data through a graph of samples and meta-learns key hyperparameters. Li et al. (2019b) proposed a  
 131 Learning to Self-Train (LST) approach that is based on self-training and meta-learns a soft weighting  
 132 network to control the influence of pseudo labelled samples and reduces label-noise during training.  
 133 Another set of approaches explore the use of self-supervision to leverage unlabelled data. Gidaris  
 134 et al. (2019) introduce auxiliary tasks, exploiting image rotations and jigsaw puzzles to learn better  
 135 feature representations. More aligned with SSL approaches and closer to our work, Yu et al. (2020)  
 136 pre-trains a classification model on base classes (in the standard FSL setting) using the imprinted  
 137 weights model and fine-tunes (without prototypes) on novel classes using the CR based mixmatch  
 138 algorithm (Berthelot et al., 2019b). While these approaches alleviate the error propagation problem  
 139 that is common when pseudo labelling is employed (Laine & Aila, 2016), their performance gains  
 140 remain limited; the techniques are not specifically adapted to the few-shot setting. Conversely,  
 141 we propose a unique training scheme that iteratively refines prototypes using both explicit average  
 142 feature representation and implicit CR refinement. We will show that this enables more flexible  
 143 feature adaptation to novel tasks and obtains more accurate class prototypes.

### 144 3 METHODOLOGY

145 We consider a base training dataset  $D_{base} = \{X_b^l, X_b^u\}$  comprising Labelled Data (LD)  
 146  $X_b^l = \{\mathbf{x}_1^l, \dots, \mathbf{x}_n^l\}$  with labels  $Y_b = \{y_1, \dots, y_n\}$ , as well as an additional set of Unlabelled Data  
 147 (UD)  $X_b^u = \{\mathbf{x}_1^u, \dots, \mathbf{x}_m^u\}$ . All examples in  $D_{base}$  belong to one of  $C_b$  base categories. Our novel  
 148 dataset  $D_{novel}$  contains  $C_n$  disjoint novel classes each with only a handful of labelled samples (e.g.  
 149  $\leq 5$ ) as well as a further limited set of unlabelled samples per class (e.g.  $\leq 100$ ) with which to  
 150 fine-tune the model.  $D_{novel}$  further comprises UD used for evaluation. Our objective, similarly to  
 151 standard few-shot settings, is to learn a classifier capable of accurately recognising novel classes,  
 152 despite having only a limited amount of available LD. However in contrast to standard FSL, we  
 153 possess additional UD for both base and novel classes, which we aim to leverage to maximise per-

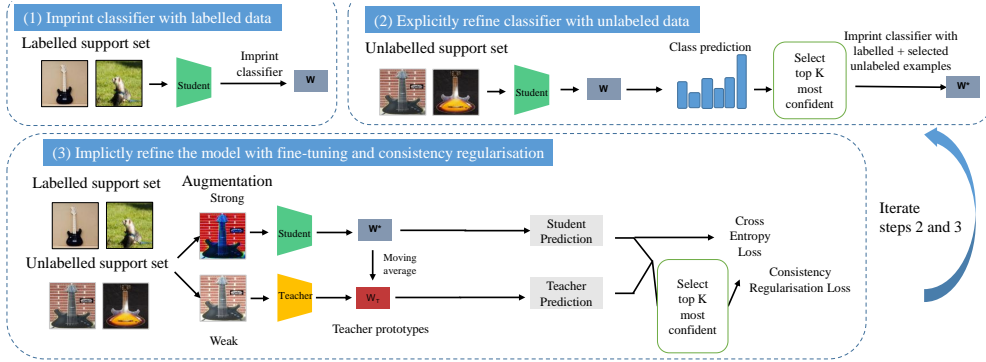


Figure 1: Overview of the Dynamic Prototype Refinement process. See main text for details.

154 formance. To formalise our setting, we consider that  $D_{novel}$  comprises of a fixed *support set* of  $K_n^l$   
 155 labelled and  $K_n^u$  unlabelled examples per class, and refer to the remaining unlabelled test images as  
 156 the query set  $Q_n$ . This  $C_n$ -way  $K_n^l$ -shot classification problem defines a standard SS-FSL setting.

157 Our proposed ‘‘FewMatch’’ method first trains a classification model on  $D_{base}$  by exploiting the con-  
 158 cept of imprinted weights (IW) (Sec 3.1). IW allow end-to-end model training, while at the same time  
 159 learning global class feature representations (commonly referred to as prototypes (Snell et al.,  
 160 2017)) utilised as classifier weights. This is achieved by computing predictions as the cosine similar-  
 161 ity between input features and classifier. End-to-end training allows seamless introduction of a CR  
 162 loss, effectively leveraging UD to train a strong feature extractor and learn high quality prototypes.  
 163 The second stage involves model fine-tuning on  $D_{novel}$  in order to leverage UD for novel classes.  
 164 We introduce a novel training scheme: our Dynamic Prototype Refinement (DPR) process (Sec 3.2).  
 165 Our iterative strategy alternates between explicit prototype refinement using feature averaging and  
 166 implicit parameter updates using fine-tuning and CR. The strong initialisation provided by the ex-  
 167 plicit averaging and longer training times afforded by the iterative scheme enable to successfully  
 168 harness the power of CR in even the lowest data regimes. An overview of the proposed DPR is  
 169 provided in Figure 1 with an analogous overview of the base training process in Appendix A.6.

170 3.1 BASE TRAINING: PROTOTYPE DRIVEN CONSISTENCY REGULARISATION

171 In this section we firstly introduce our imprinted weight formulation and then describe the integra-  
 172 tion of this within a teacher-student framework, enabling the introduction of our CR loss.

173 **Imprinted weights formulation.** Our classification model uses a standard architecture, compris-  
 174 ing a feature extraction network  $\theta_f$ , and a classifier defined by a fully connected layer without bias  
 175  $W \in F \times C_b$ , where  $F$  is the output dimension of  $\theta_f$ . The main idea of imprinted weights is to  
 176 train the model such that, for a given class  $c$ , the cosine similarity between the embedding vector  
 177  $\theta_f(\mathbf{x})$  of input image  $\mathbf{x}$  and the corresponding column  $w_c$  of  $W$  is maximised. By normalising  
 178 the classifier and embedding vectors, the model can be trained end-to-end using a standard cross  
 179 entropy loss. In this setting,  $w_c$  is regarded as the *prototype* representation of class  $c$  and can be  
 180 learned implicitly without the, typically required, episode training strategy and support set aver-  
 181 aging. More formally, for input sample  $\mathbf{x}$ , the set of classification scores output by the model is  
 182  $f(x) = \{f^1(x), f^2(x), \dots, f^c(x), \dots, f^{C_b}(x)\}$  and the score for a given class  $c$  is computed as:

$$f^c(\mathbf{x}) = \frac{\exp(\gamma(\mathbf{w}_c^T, \theta_f(\mathbf{x})))}{\sum_{i=1}^{C_b} \exp(\gamma(\mathbf{w}_i^T, \theta_f(\mathbf{x})))} \tag{1}$$

183 where  $w_i$  is the  $i^{th}$  column of weight matrix  $W$  and the prototype  $p_i$  of class  $i$ . The scaled cosine  
 184 similarity is then given by  $\gamma(\mathbf{w}_i^T, \theta_f(\mathbf{x})) = s \cdot \mathbf{w}_i^T(\theta_f(\mathbf{x}))$ .  $w_i$  and  $\theta_f(\mathbf{x})$  are normalized using the  $L_2$   
 185 norm, and  $s$  is a trainable scalar, as introduced by (Qi et al., 2018) to avoid the risk that the cosine  
 186 distance yields distributions lacking in discriminative power. Finally, the classification loss can be

187 calculated as:  $\mathcal{L}_{ce}(\mathbf{x}) = -\sum_{c=1}^{C_b} \delta_{c,y} \log f^c(\mathbf{x})$  where  $\delta_{c,y}$  is the Dirac delta function. Defining  
 188 class prototypes as learnable model weights affords end-to-end training and enables introduction of  
 189 CR to our model in a natural fashion. These decisions allow us to leverage UD and implicitly refine  
 190 prototypes without explicit pseudo-labelling. Furthermore, this approach optimises the base class  
 191 learning process by allowing full exploitation of the available LD without the typical requirement  
 192 that necessitates simulation of the few-shot set-up (episode training) (Finn et al., 2017).

193 **Consistency Regularisation.** We highlight that the described training strategy does not yet lever-  
 194 age UD, available in the considered SS-FSL problem setting. Towards taking advantage of UD,  
 195 we introduce a CR loss (Tarvainen & Valpola, 2017) that is driven by the learned prototypes. The  
 196 idea underlying CR is to regularise predictions such that they become invariant to small input per-  
 197 turbations that do not affect class semantics. This strategy has been used successfully for a variety  
 198 of problems and is particularly appealing in the semi-supervised context as it leverages UD with-  
 199 out explicit pseudo-labelling. A key difference in our setting, with respect to conventional SSL,  
 200 is that our CR loss directly depends on prototype instantiations, as predictions are based on the  
 201 distance between input and each class prototype. This strategy drives our approach to learn more  
 202 discriminative and robust prototypes towards maintaining classification accuracy under different in-  
 203 put perturbations. Following strategies adopted in the recent SSL state of the art (Berthelot et al.,  
 204 2019a), we embed our IW model within a teacher-student framework (Tarvainen & Valpola, 2017)  
 205 where we seek to impose consistency between teacher and student predictions. Both teacher and  
 206 student networks share the same architecture, however only student weights are optimised by back-  
 207 propagation. Teacher weights  $\theta_T$  are computed as an Exponential Moving Average (EMA) of the  
 208 student weights  $\theta$ ,  $\theta_T = (1 - \alpha)\theta_T + \alpha\theta$ . Such temporal averaging strategies have been shown to  
 209 yield more robust and accurate models and are therefore desirable in often noisy few-shot settings.

210 Considering an unlabelled sample  $u_b$  we realise sample perturbations, as suggested in (Xie et al.,  
 211 2019; Berthelot et al., 2019a), by generating  $\tilde{u}_b$  and  $\hat{u}_b$  using *weak* and *strong* augmentations re-  
 212 spectively. The weak augmentation sample  $\tilde{u}_b$  has the goal of improving prediction stability in the  
 213 teacher network. This strategy helps to constrain the strong augmentation sample prediction. The  
 214 consistency loss is then computed as:  $\mathcal{L}_{cons}(u_b) = \|\text{Sharp}(f_t(\hat{u}_b), \mathcal{T}) - f_s(\tilde{u}_b)\|^2$ ; where  $f_s$  and  $f_t$   
 215 are predictions computed by the student and teacher networks respectively; and  $\text{Sharp}(\cdot)$  is a sharp-  
 216 ening function, parametrised by temperature  $\mathcal{T}$ , introduced in (Berthelot et al., 2019b) to reduce the  
 217 entropy of the label distribution. In summary, the model is trained on the base classes using global  
 218 loss  $\mathcal{L}_{base} = \mathcal{L}_{ce} + \lambda\mathcal{L}_{cons}$ , where hyperparameter  $\lambda$  balances the relative influence of the terms.

### 219 3.2 DYNAMIC PROTOTYPE REFINEMENT

220 Our training stage, considering  $D_{base}$ , yields a model capable of estimating reliable class prototypes  
 221 on novel, unseen categories. In a standard few-shot setting (i.e. without available UD), prototypes  
 222 are often estimated directly from the support set and reliable performance can be achieved without  
 223 further training. In our problem setting, we set the objective of exploiting the additionally available  
 224 UD in order to obtain strong prototype initialisations that then lend themselves to further refine-  
 225 ment. Towards this goal, the main component of FewMatch constitutes our Dynamic Prototype  
 226 Refinement (DPR) strategy, taking advantage of the UD available from  $D_{novel}$ , with the aim of im-  
 227 proving model adaption to novel categories. By design our approach is able to improve performance  
 228 on novel categories despite the presence of limited data regimes. DPR comprises three stages: **(1)**  
 229 Prototype Initial Inference (PII), via the introduced IW procedure **(2) Explicit** prototype refinement  
 230 using top-K selection and **(3) Implicit** prototype refinement using CR. Prototypes are initially es-  
 231 timated during the first step and then dynamically updated using iterations of steps two and three,  
 232 such that prototype quality is iteratively improved. The remainder of this section provides further  
 233 detail on steps **(1)**-**(3)** and the iterative process.

234 **Prototype Initial Inference.** Given new category  $j$  from  $D_{novel}$  with support set  
 235  $S_j = \{\mathbf{x}_1^s, y_1^s, \dots, \mathbf{x}_n^s, y_n^s\} \cup \{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ , compute an initial prototype using the labelled sup-  
 236 port set as:

$$p_j^* = P(S_j) = \frac{1}{|S_j|} \sum_{\mathbf{x}_i^s \in S_j} \theta_f(\mathbf{x}_i^s), \quad (2)$$

237 The estimated prototype is then imprinted in classifier  $W$  as  $w_j = p_j$  and the process is repeated  
 238 for each new category (see Figure 1). This allows for recognition of new classes without model  
 239 retraining and provides high quality initialisation for our dynamic refinement stage.

240 **Explicit Prototype Refinement.** We highlight that initial prototypes, computed using Eq. 2, do not  
 241 make use of the additional UD available for novel classes. Exploiting UD can be considered crucial  
 242 for novel classes due to the availability of only limited labelled data. Towards reducing prototype  
 243 biases, we expand the support set using pseudo-labelled UD, where labels are assigned according  
 244 to respective prediction scores. The prediction scores  $f_s(u)$  are again obtained with Eq. 1 using up-  
 245 dated prototype estimates and current model parameters. We mitigate the varying quality of pseudo  
 246 labels by selecting the top- $K$  samples with the most confident predictions per class which, by defi-  
 247 nition, consist of the  $K$  unlabelled samples that are closest to their assigned class’ prototypes. This  
 248 augmentation results in an extended annotated support set defined for each class  $j$  as  $S_j^* = S_j \cup U_j$ ,  
 249 where  $U_j = \text{top-}K(f_s^j(u))$  is the set of unlabelled samples selected for class  $j$ . The prototype is then  
 250 refined using Eq. 2 by replacing  $S$  with  $S^*$ . Crucially, we emphasise that per stage pseudo-labels are  
 251 used *uniquely* to update prototypes and that samples, pseudo-labelled at this stage, are considered  
 252 unlabelled again at the next iteration. Importantly *pseudo-labels are therefore not propagated*, al-  
 253 lowing for recovery from potentially erroneous predictions during the subsequent fine-tuning stage.

254 **Implicit Refinement using Consistency Regularisation.** Our implicit refinement stage inherits  
 255 ideas from gradient-based FSL, which typically adapts the entire model to novel classes via a fine-  
 256 tuning stage. This stage is generally missing from prototype-based methods, which explicitly repre-  
 257 sent prototypes as an average feature representation, and thus lose the flexibility afforded by learning  
 258 implicit network parameters. This fine-tuning stage is particularly desirable in our setting, where we  
 259 seek to maximally leverage the available UD and our prototypes are defined as model weights. It is a  
 260 natural choice to consider deploying Consistency Regularisation to fine-tune the model, noting that  
 261 the refined prototypes obtained at this stage afford high quality teacher predictions. We implement  
 262 the strategy described in Sec. 3.1 to fine-tune the model on novel classes with CR. To further im-  
 263 prove robustness to noisy teacher predictions and difficult examples, we adopt a selective prototype  
 264 CR strategy. By calculating teacher prediction scores  $f_t(u)$  according to their prototype distance,  
 265 we can select the top- $K$  unlabelled examples with the least ambiguous label predictions to compute  
 266 the CR loss. Note that this second top- $K$  selection set  $V$  will differ from top- $K$  set  $U$  computed  
 267 during the explicit stage, as 1) prototypes were updated 2) they are computed on the teacher model  
 268 subject to weak input augmentation. The model is fine-tuned for R gradient updates by minimising  
 269  $\mathcal{L}(\mathbf{x}, v_b^u) = \mathcal{L}_{ce}(\mathbf{x}) + \lambda_{ft} \mathcal{L}_{cons}(v^u)$ , where  $\mathcal{L}_{ce}$  and  $\mathcal{L}_{cons}$  are computed as described in Sec. 3.1,  
 270 where labelled sample  $\mathbf{x}$  is from  $\mathcal{D}_{novel}$  and  $v^u \in V$ .

271 **Dynamic Prototype Refinement.** Our implicit and explicit refinement steps allow iterative pro-  
 272 totype refinement towards further performance improvement. We alternate between explicit and  
 273 implicit steps for  $M$  iterations, reinitialising estimated pseudo-label at each iteration. Top- $K$  selec-  
 274 tion, for the first explicit stage, relies on student predictions since teachers are randomly initialised.  
 275 Teacher predictions, presumed to be more accurate and stable, are used in subsequent iterations. Im-  
 276 portantly, we note that teacher parameters are reinitialised before each implicit stage (after explicit  
 277 selection) thus introducing stochasticity, increasing robustness to pseudo-label errors and aiding loss  
 278 optimization. Algorithm details for dynamic prototype refinement are provided in Appendix A.5

## 279 4 EXPERIMENTS

280 **Experimental set-up.** We evaluated Fewmatch on two standard SS-FSL benchmarks:  
 281 *miniImageNet* (Vinyals et al., 2016) and *tieredImageNet* (Ren et al., 2018), both subsets of the  
 282 ImageNet dataset (Russakovsky et al., 2015) designed specifically for FSL. *MiniImageNet* cons-  
 283 sists of 100 classes with 600 image samples per class. We use the standard 64/16/20 classes split  
 284 for train/val/test sets (Vinyals et al., 2016) and use 40%/60% of the data for labelled/unlabelled  
 285 splits following previous works (Ren et al., 2018; Li et al., 2019b). *TieredImageNet* contains 608  
 286 classes from 34 super-level categories. These are divided into 20/6/8 coarse super-level categories  
 287 for train/val/test splits and contain 351, 97 and 160 classes, respectively. We follow the standard  
 288 semi-supervised split (Ren et al., 2018; Li et al., 2019b), with 10% of the images of each class  
 289 forming the labelled split and the remaining 90% being the UD. We consider  $K_n^l = 5$  way  $N=1, 5$   
 290 shot classification problems and follow the strategy adopted in (Ren et al., 2018; Li et al., 2019b) to

Table 1: Mean classification accuracies of the 5-way 1/5-shot tasks. (**Bold**: Best results per set-up). SL+U setting uses all available training LD (SL setting) with additional UD vs SSL using 10% (*tieredImageNet*) or 40% LD (*miniImageNet*). Grey rows: methods using self-supervision.

Setting	Model	Backbone	<i>miniImageNet</i>		<i>tieredImageNet</i>	
			1-shot	5-shot	1-shot	5-shot
SL	MTL (Sun et al., 2019)	<i>ResNet-12</i>	61.20 ±1.80	75.50 ±0.80	-	-
	CTM (Li et al., 2019a)	<i>ResNet-18</i>	62.05 ±0.55	78.63 ±0.06	64.78 ±0.11	81.05 ±0.52
	CC+rot (Gidaris et al., 2019)	<i>WRN-28-10</i>	62.93 ±0.45	79.87 ±0.33	70.53 ±0.51	84.98 ±0.36
SL + U	CC+rot+unlabelled	<i>WRN-28-10</i>	64.03 ±0.46	80.68 ±0.33	-	-
	TransMatch (Yu et al., 2020)	<i>WRN-28-10</i>	63.02±1.07	81.19±0.59	-	-
SSL	MS k-Means (Ren et al., 2018)	4Conv	50.4	64.4	52.4	69.9
	MS k-Means with MTL	ResNet-12	62.1	73.6	68.6	81.0
	TPN (Liu et al., 2018)	4Conv	52.8	66.4	55.7	71.0
	TPN with MTL	ResNet-12	62.7	74.2	72.1	83.3
	LST (Li et al., 2019b)	ResNet-12	70.1 ±1.9	78.7 ±0.8	77.7 ±1.6	85.2 ±0.8
	Ours	ResNet-12	<b>75.66±0.95</b>	<b>82.93±0.62</b>	<b>78.70±0.93</b>	<b>85.40±0.58</b>
Distractor Setting						
SL + U	TransMatch	<i>WRN-28-10</i>	59.32±1.10	79.29±0.62	-	-
SSL	MS k-Means with MTL	ResNet-12	61.0	72.0	66.9	80.2
	TPN with MTL	ResNet-12	61.3	72.4	71.5	82.7
	LST	ResNet-12	64.1	77.4	73.4	83.4
	Ours	ResNet-12	<b>70.35±0.98</b>	<b>80.23±0.66</b>	<b>74.24±0.95</b>	<b>83.64±0.63</b>

291 generate test episodes: we randomly sample  $K_n^l$  classes from the test set,  $N$  labelled images from  
 292 each class, 100 unlabelled images as support images and 15 query images.

293 The previous protocol can be regarded as a standard set-up that we follow for fair comparisons.  
 294 Towards exploring more realistic few-shot testing scenarios, we consider two additional directions.  
 295 Firstly, the *distractor* setting (Li et al., 2019b) introduces UD from irrelevant classes, providing a  
 296 more challenging test environment. Testing involves randomly selecting 100 unlabelled images from  
 297 three task-irrelevant classes to serve as *distractors* and adding these to the unlabelled set. Table 1  
 298 (lower), reports mean accuracy for 600 randomly generated test episodes in comparison to the state-  
 299 of-the-art for this challenging setting. Secondly, the absence of an episode-based training require-  
 300 ment affords FewMatch additional flexibility and enables more realistic SS-FSL testing schemes,  
 301 e.g. investigating model adaptation capabilities under varying amounts of UD per class. We pro-  
 302 vide classification accuracies for settings with unbalanced class sampling: (1) randomly selecting  
 303 between 70-130 US per class; (2) 80-120 US per class. As Table 2 shows, FewMatch performance  
 304 retains stability in unbalanced settings, c.f. the balanced default (exactly 100 US per class).

305 The method was implemented with PyTorch (Paszke et al., 2017) using the same ResNet-12 back-  
 306 bone as (Li et al., 2019b). For base category training, we follow parameters used in (Gidaris &  
 307 Komodakis, 2018): our model is optimised using SGD with momentum 0.9, weight decay 0.0005,  
 308 mini-batchsize 256 (128 LD and 128 UD) for 30 epochs. All input images were resized to  $84 \times 84$ .  
 309 The learning rate was initialised to 0.1, and updated to 0.01 at epoch 20. Following SSL practice (Tar-  
 310 vainen & Valpola, 2017), weighting parameter  $\lambda$  is defined as a linear ramp-up function increasing  
 311 from 0 to 300 in the first 15 epochs. We set the total number of DPR iterations as  $M = 3$  and each  
 312 implicit refinement step fine-tunes the model for 20 steps with 0.01 learning rate. Each mini-batch  
 313 comprises all LD and 40 randomly sampled UD per-category. We linearly increase weighting pa-  
 314 rameter  $\lambda_{ft}$  from 0 to 10 in the first 10 steps. The number of unlabelled samples selected is set to  
 315  $K = 25$ . We set EMA rate  $\alpha = 0.5$ , and  $\mathcal{T} = 0.5$ . Strong augmentations for the student network  
 316 are computed using RandAugment (e.g. color, shear) (Cubuk et al., 2019), applying three random  
 317 operations with magnitude set to 9. Teacher weak augmentations use random cropping and flipping.

318 **Comparison to State-of-the-Art (SOTA) methods.** We compared FewMatch with SOTA ap-  
 319 proaches including (a) 3 FSL and (b) 5 SS-FSL methods in Table 1. We note that several SS-FSL  
 320 approaches, including FewMatch, outperform SOTA FSL approaches, highlighting the potential  
 321 of using additional UD to learn more accurate models. We observe that FewMatch outperforms  
 322 the SS-FSL state of the art in both standard and distractor settings and that strongest performance

Table 2: Ablation study on *miniImageNet*. PCR: base training prototype Consistency Regularisation; ER: Explicit prototype refinement; IR: Implicit refinement using Selective Consistency Regularisation; DR: Dynamic Refinement

Model Components				<i>miniImageNet</i>	
PCR	ER	IR	DR	1-shot	5-shot
Remixmatch				53.52	66.50
Imprinted-weights (IW)				59.09	75.59
IW + Remixmatch (no mixup)				62.20	76.31
✓				61.59	77.90
✓	✓			71.35	81.75
✓	✓	✓		72.52	82.25
✓	✓	✓	✓	75.66	82.93
Unbalanced Number of Unlabelled Samples					
min/max US 70/130				74.24	82.51
min/max US 80/120				75.14	82.82

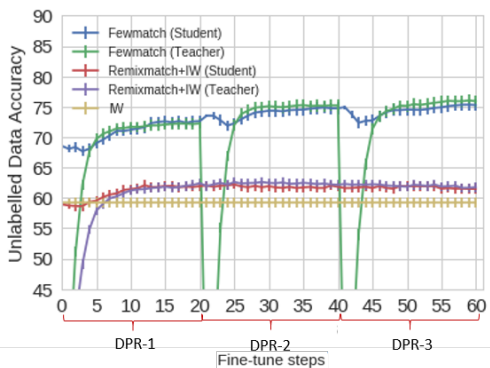


Figure 2: Accuracy on training unlabelled data with  $M = 3$  iterations of the DPR stage.

gains are observed in the 1-shot setting. We further highlight that 1) we significantly outperform self-supervision methods that use a more powerful backbone encoder and were trained in a more favourable setting (SL+U: using all base LD with additional UD, vs SSL setting using a fraction of LD only) and 2) the closest SOTA method LST, requires, in contrast to FewMatch, complex episode training, requiring a fixed number of LD and UD at both training and test time.

**Ablation experiments.** We evaluate the influence of each model component using *miniImageNet* under 5 way 1/5 shot settings. Specifically, we evaluate the influence of using CR in the base training stage (PCR), Explicit Prototype refinement (ER), Implicit Refinement (IR) and Dynamic Refinement (DR) which iterates between ER and IR. We additionally include three baselines: Imprinted Weights (Qi et al., 2018) (no use of unlabelled data), SOTA CR based SSL method Remixmatch (Berthelot et al., 2019a) (no accounting for the few-shot setup), and Imprinted Weights combined with Remixmatch. We highlight that the latter baseline is highly similar to the method of Yu et al. (2020) and provides context towards the performance expected in the SSL setting. We note that methods using Remixmatch use CR during *both base and novel* training stages and that the latter method is implemented without mixup (used in the Remixmatch method) as the label mixing strategy is not compatible with the prototype approach and would require the definition of infinitely many prototypes. Results are reported in Table 2 and show that each component makes a clear contribution to the performance gain; with ER (providing a strong initialisation) and DR (addressing slow CR convergence rates) yielding the strongest performance gains.

**Analysis of the DPR process.** Figure 2 evaluates the improved reliability of teacher predictions throughout our DPR process ( $M=3$ ). We report accuracy on training UD during the DPR stage, compared to baseline imprinted weights + remixmatch (IWR) which uses CR without addressing the underlying challenges. We observe that our iterative process continuously improves performance, successfully exploiting CR towards reaching higher quality predictions. Conversely, the IWR model fails to exploit UD, obtaining a minimal performance gain with respect to baseline FSL method IW.

## 5 CONCLUSION

We introduced a novel prototype-driven approach named FewMatch, designed specifically to exploit the power of consistency regularisation in limited data regimes. In contrast with pre-existing state of the art methods, we alleviate requirements for iterative pseudo-labelling, preventing propagation of errors induced by inaccurate model predictions. We go beyond the introduction of self-supervised auxiliary losses and propose a novel training strategy: a dynamic prototype refinement that alternates between explicit pseudo label based updates and implicit model fine-tuning. Our extensive experiments demonstrate that this iterative strategy allows successful exploitation of unlabelled data within a consistency regularisation framework, yielding large performance gains.



## 357 REFERENCES

- 358 Antreas Antoniou, Harrison Edwards, and Amos Storkey. How to train your maml. *arXiv*, 2018.
- 359 David Berthelot, Nicholas Carlini, Ekin D Cubuk, Alex Kurakin, Kihyuk Sohn, Han Zhang, and  
360 Colin Raffel. Remixmatch: Semi-supervised learning with distribution alignment and augmenta-  
361 tion anchoring. *arXiv*, 2019a.
- 362 David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A  
363 Raffel. Mixmatch: A holistic approach to semi-supervised learning. In *NeurIPS*, 2019b.
- 364 Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated  
365 data augmentation with a reduced search space. *arXiv*, 2019.
- 366 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation  
367 of deep networks. In *ICML*, 2017.
- 368 Spyros Gidaris and Nikos Komodakis. Dynamic few-shot visual learning without forgetting. In  
369 *CVPR*, 2018.
- 370 Spyros Gidaris, Andrei Bursuc, Nikos Komodakis, Patrick Pérez, and Matthieu Cord. Boosting  
371 few-shot visual learning with self-supervision. *arXiv*, 2019.
- 372 Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. Label propagation for deep semi-  
373 supervised learning. In *CVPR*, 2019.
- 374 Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. *arXiv*, 2016.
- 375 Dong-Hyun Lee. Pseudo-label: The simple and efficient semi-supervised learning method for deep  
376 neural networks. In *ICMLW*, 2013.
- 377 Hongyang Li, David Eigen, Samuel Dodge, Matthew Zeiler, and Xiaogang Wang. Finding task-  
378 relevant features for few-shot learning by category traversal. In *CVPR*, 2019a.
- 379 Xinzhe Li, Qianru Sun, Yaoyao Liu, Qin Zhou, Shibao Zheng, Tat-Seng Chua, and Bernt Schiele.  
380 Learning to self-train for semi-supervised few-shot classification. In *NeurIPS*, 2019b.
- 381 Yanbin Liu, Juho Lee, Minseop Park, Saehoon Kim, Eunho Yang, Sung Ju Hwang, and Yi Yang.  
382 Learning to propagate labels: Transductive propagation network for few-shot learning. *arXiv*,  
383 2018.
- 384 Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito,  
385 Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in  
386 pytorch. 2017.
- 387 Hang Qi, Matthew Brown, and David G Lowe. Low-shot learning with imprinted weights. In *CVPR*,  
388 2018.
- 389 Siyuan Qiao, Chenxi Liu, Wei Shen, and Alan L Yuille. Few-shot image recognition by predicting  
390 parameters from activations. In *CVPR*, 2018.
- 391 Aravind Rajeswaran, Chelsea Finn, Sham Kakade, and Sergey Levine. Meta-learning with implicit  
392 gradients. *arXiv*, 2019.
- 393 Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B Tenenbaum,  
394 Hugo Larochelle, and Richard S Zemel. Meta-learning for semi-supervised few-shot classifica-  
395 tion. *arXiv*, 2018.
- 396 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng  
397 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual  
398 recognition challenge. *IJCV*, 2015.
- 399 Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero,  
400 and Raia Hadsell. Meta-learning with latent embedding optimization. *arXiv*, 2018.

- 401 Mehdi Sajjadi, Mehran Javanmardi, and Tolga Tasdizen. Regularization with stochastic transforma-  
402 tions and perturbations for deep semi-supervised learning. In *Neurips*, 2016.
- 403 Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In  
404 *NeurIPS*, 2017.
- 405 Qianru Sun, Yaoyao Liu, Tat-Seng Chua, and Bernt Schiele. Meta-transfer learning for few-shot  
406 learning. In *CVPR*, 2019.
- 407 Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales.  
408 Learning to compare: Relation network for few-shot learning. In *CVPR*, 2018.
- 409 Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consis-  
410 tency targets improve semi-supervised deep learning results. In *NeurIPS*, 2017.
- 411 Vikas Verma, Alex Lamb, Juho Kannala, Yoshua Bengio, and David Lopez-Paz. Interpolation con-  
412 sistency training for semi-supervised learning. *arXiv*, 2019.
- 413 Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one  
414 shot learning. In *NeurIPS*, 2016.
- 415 Qizhe Xie, Eduard Hovy, Minh-Thang Luong, and Quoc V Le. Self-training with noisy student  
416 improves imagenet classification. *arXiv*, 2019.
- 417 Zhongjie Yu, Lin Chen, Zhongwei Cheng, and Jiebo Luo. Transmatch: A transfer-learning scheme  
418 for semi-supervised few-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer  
419 Vision and Pattern Recognition*, pp. 12856–12864, 2020.
- 420 Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical  
421 risk minimization. *arXiv*, 2017.

## 422 A APPENDIX

423 We provide additional material to supplement our work. Section A.1 evaluates the influence of the  
424 number of unlabelled samples on FewMatch’s performance, and demonstrates the method’s ability to  
425 leverage unlabelled examples. In Section A.2, we report a comparison between the Semi-Supervised  
426 Learning (SSL), Few Shot Learning (FSL) and Semi-Supervised Few Shot Learning (SS-FSL) set-  
427 tings, highlighting the challenges associated with SS-FSL. In Section A.3, we report an additional  
428 experiment, studying the influence of Dynamic Prototype Refinement (DPR) iterations  $M$  on our  
429 model performance. Please refer to our main paper for further method details. In Section A.4, we  
430 further synthesize the comparison between *FewMatch* and existing SS-FSL approaches, explicitly  
431 providing additional details to highlight the main differences between the considered methods. In  
432 Section A.5, we provide pseudocode description of our Dynamic Prototype Refinement process. Fi-  
433 nally, Section A.6 provides detailed pseudocode for the first stage of our method (prototype-driven  
434 consistency regularisation as described in Section 3.1 of the main paper).

### 435 A.1 INFLUENCE OF THE NUMBER OF UNLABELLED SAMPLES

436 We test the impact of using variable amounts of US per class on classification accuracy in the 5-  
437 way 1-shot setting on *miniImageNet*. Results are shown in Figure 3, showing a large increase in  
438 performance when including 50 US and a more modest yet consistent improvement as the number  
439 of US increases. This highlights the advantage provided by the use of US to complement the few-  
440 shot labelled examples, as well as FewMatch’s ability to leverage unlabelled examples.

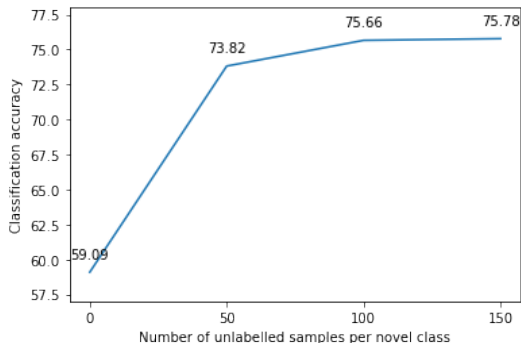


Figure 3: Mean classification accuracy on 5-way 1-shot on *miniImageNet* with varying amounts of unlabelled samples.

#### 441 A.2 COMPARISON BETWEEN SSL, FSL, AND SS-FSL SETTINGS

442 In Table 3, we report training sample counts (labelled and unlabelled) per category used in FSL,  
 443 SSL and SS-FSL settings. The stated values follow the convention in Ren et al. (2018) (5-way  
 444 1-shot) on Mini-ImageNet and, for SSL, we report the setting comprising the minimal LS with  
 445 respect to recent state of art methods Berthelot et al. (2019a) on the common benchmark, CIFAR-  
 446 10. Compared to FSL, this table highlights that 1) fewer labelled data is available during the base  
 447 training stage, increasing the difficulty of obtaining a strong initialisation and 2) a substantial amount  
 448 of additional unlabelled data is available for novel classes. Compared to SSL, the amount of labelled  
 449 and unlabelled samples is significantly reduced in the SS-FSL setting (in particular; the unlabelled  
 450 samples), highlighting the challenges associated with adapting SSL methods to the SS-FSL scenario.

Table 3: Comparison of available per category training Labelled Samples (LS) and Unlabelled Samples (US) between FSL, SS-FSL, SSL)

Data Split	FSL	SS-FSL	SSL
Base classes	600 LS	240 LS + 360 US	-
Novel classes	1 LS	1LS + 100 US	25 LS + 4750 US

#### 451 A.3 PARAMETER STUDY: DYNAMIC PROTOTYPE REFINEMENT (DPR) ITERATIONS

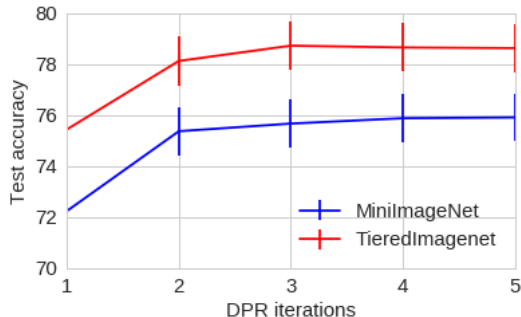


Figure 4: Dynamic Prototype Refinement (DPR) performance with respect to iterations  $M$

452 We evaluate the influence of DPR iteration count  $M$  with respect to model performance in the 5-  
 453 way 1-shot setting and report respective test accuracies in Figure 4. We observe similar behaviour  
 454 for both datasets considered (*miniImageNet* and *tieredImageNet*), with performance improving and

455 then stabilising for  $M \geq 3$ . Our model requires only three iterations to reach optimal performance  
 456 in the investigated settings.

457 A.4 COMPARISON OF FEWMATCH AND EXISTING SS-FSL APPROACHES

458 In Table 4, we provide an additional detailed comparison of FewMatch with state of the art SS-FSL  
 459 approaches, including Masked Soft  $k$ -Means (Ren et al., 2018), Transmatch (Yu et al., 2020) and  
 460 LST (Li et al., 2019b). We compare six different characteristics of the methods: Base dataset split,  
 461 Training Strategy, Prototype estimation, Classifier learning approach, backbone encoder adaptation  
 462 strategy (to novel task) and SSL approach used. Table 2 illustrates that 1) FewMatch provides a more  
 463 flexible training strategy as it does not require episodic training. This allows consideration of differ-  
 464 ent set-ups at test time in contrast to episodic training that typically enforcing a fixed  $K$ -way- $N$ -shot  
 465 setting. 2) Compared to Masked Soft  $k$ -Means, the only other method using prototypes, FewMatch  
 466 adopts a more flexible prototype learning process by combining feature averaging with fine-tuning.  
 467 This is enabled by the fact that prototypes are defined as classifier weights, allowing learning of high  
 468 quality prototype representations. Furthermore, FewMatch adapts the feature backbone to the novel  
 469 task, reducing the influence of domain shift. 3) In contrast to LST, Fewmatch combines classifier pa-  
 470 rameter updates with the concept of prototypes, allowing a stronger initialization for the fine-tuning  
 471 stage to be obtained. 4) In contrast to TransMatch, FewMatch uses fewer labelled training examples  
 472 in the base training stage, and fine tunes the model using a combination of feature averaging and  
 473 backpropagation; affording better CR convergence.

Table 4: Comparison of FewMatch to existing SS-FSL approaches

Method	Masked Soft $k$ -Means	Transmatch	LST	FewMatch
<b>Base dataset</b>	60% US+40% LS	100% LS	60% US+40% LS	60% US+40% LS
<b>Training</b>	Episodic	End to end	Episodic	End to end
<b>Prototypes</b>	Feature averaging	/	/	Iterative feature
<b>Classifier</b>	/	backpropagation	backpropagation	averaging and backprop
<b>Feature</b>	Fixed	Fixed	Adapted to novel task	Adapted to novel task
<b>Learning</b>	Pseudo label	CR	Pseudo label	CR

474 A.5 DYNAMIC PROTOTYPE REFINEMENT ALGORITHM

475 We provide an algorithmic description of our Dynamic Prototype Refinement (DPR) process in  
 476 Algorithm 1. DPR contains three steps: 1) Prototypes initial inference; 2) Explicit prototype refine-  
 477 ment; 3) Implicit refinement using CR. We alternate between explicit and implicit refinement for  $M$   
 478 epochs after the initial inference step.

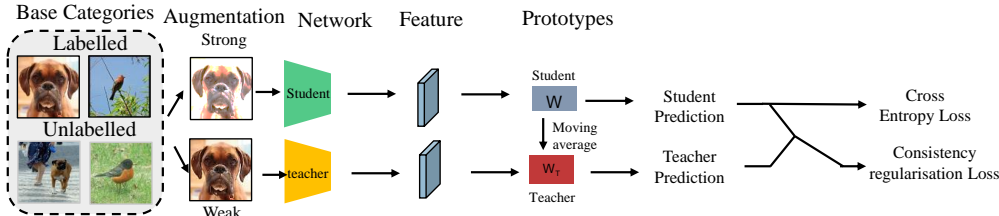


Figure 5: Base training process

**Algorithm 1** Dynamic Prototype Refinement

- 
- 1: **Input:** labelled examples  $\mathcal{S} = \{S_1, \dots, S_j, S_{C_n}\}$ , and unlabelled examples  $\mathcal{U}$ ; Number of novel categories:  $C_n$ ; number of iterations  $M$ ; number of fine-tuning steps  $R$ ; pre-trained student and teacher model parameters  $\theta, \theta_T$ ; weighting parameters  $\lambda_{ft}, \alpha$ .
  - 2: **Output:** Prototypes of novel categories  $W^{**}$ , student model parameters  $\theta$ ;
  - 3: *Prototypes initial inference:*  $W \leftarrow \{p_1^*, p_2^*, \dots, p_{C_n}^*\}$ , calculate  $p_j^* \leftarrow P(S_j)$  by Eq equation 2
  - 4: **For**  $i = 1$  **to**  $M$  :
    - 5: *Explicit prototype refinement*
    - 6:  $U_j \leftarrow \text{top-}K(f_{t, \theta_T, W_T}^j(u)), \forall j \in 1, \dots, C_n$ ,  $f_{t, \theta_T, W_T}^j$  computed by equation 1  
with parameters  $\theta_T, W_T$   $\triangleright \theta_T, W_T$  initialised to  $\theta, W$  for  $i = 1$
    - 7:  $S_j^* \leftarrow U_j \cup S_j \quad \forall j \in 1, \dots, C_n$
    - 8:  $W^* \leftarrow \{P(S_1^*), \dots, P(S_{C_n}^*)\}$
    - 9: *Implicit refinement using CR*
    - 10: Randomly re-initialise teacher parameters  $\theta_T$
    - 11: **For**  $r = 1$  **to**  $R$ :
      - 12: Sample a batch of unlabelled samples  $\mathcal{U}_s$  from  $\mathcal{U}$
      - 13:  $\bar{u} \leftarrow \text{WeakAugment}(u), \hat{u} \leftarrow \text{StrongAugment}(u), u \in \mathcal{U}_s$
      - 14:  $V_j \leftarrow \text{top-}K(f_{t, \theta_T, W^*}^j(\bar{u})) \quad \forall j \in 1, \dots, C_n$
      - 15:  $W^{**}, \theta^* \leftarrow \arg \min_{W, \theta} \mathcal{L}_{ce}(\mathbf{x}) + \lambda_{ft} \mathcal{L}_{cons}(v_b^u), \quad \mathbf{x} \in \mathcal{S}, v^u \in V = \{V_1, \dots, V_{C_n}\}$
      - 16: *Update teacher parameters*  $W_T \leftarrow (1 - \alpha)W_T + \alpha W^{**}, \quad \theta_T \leftarrow (1 - \alpha)\theta_T + \alpha \theta^*$
      - 17: **end**
- 

**Algorithm 2** Prototype Driven Consistency Regularization

- 
- 1: **Input:** Labelled examples and their one-hot labels  $\mathcal{X} = \{(x_b, y_b) : b \in 1, \dots, B\}$ , Unlabelled examples  $\mathcal{U} = \{(u_b) : b \in 1, \dots, B\}$ , weighting parameters  $\lambda, \alpha$ .
  - 2: **Output:** Optimised student model parameters  $\theta^*, W^*$
  - 3: Randomly initialise Student and Teacher model parameters and prototypes:  $\theta, \theta_T, W, W_T$
  - 4: **While** not done **do**
    - 5: Sample batch of labelled  $\mathcal{X}_b$  and unlabelled samples  $\mathcal{U}_b$  from  $\mathcal{X}, \mathcal{U}$
    - 6: **for all**  $(x_b, u_b) \in (\mathcal{X}_b, \mathcal{U}_b)$  **do**
    - 7:  $\hat{x}_b = \text{StrongAugment}(x_b)$
    - 8:  $\bar{u}_b = \text{WeakAugment}(u_b)$
    - 9:  $\hat{u}_b = \text{StrongAugment}(u_b)$
    - 10:  $q_b^l \leftarrow f_{s, \theta, W}(\hat{x}_b), f_{s, \theta, W}(\hat{x}_b)$  computed as in Eq (1) in main-manuscript with student parameters  $\theta, W$
    - 11:  $q_b^u \leftarrow f_{t, \theta_T, W_T}(\bar{u}_b), \hat{q}_b^u = f_{s, \theta, W}(\hat{u}_b)$
    - 12:  $\mathcal{L}(x_b, u_b) = \mathcal{L}_{ce}(q_b^l) + \lambda \|\text{Sharp}(q_b^l, \mathcal{T}) - \hat{q}_b^u\|^2$  as in Eq (2) in main-manuscript
    - 13:  $W^*, \theta^* \leftarrow \arg \min_{W, \theta} \sum_{\mathcal{X}_b, \mathcal{U}_b} \mathcal{L}(x_b, u_b)$
    - 14: *Update teacher parameters*  $W_T \leftarrow (1 - \alpha)W_T + \alpha W^*, \quad \theta_T \leftarrow (1 - \alpha)\theta_T + \alpha \theta^*$
    - 15: **end**
-