
Supplemental Materials: Model-Agnostic Graph Regularization for Few-Shot Learning

Anonymous Author(s)

Affiliation

Address

email

1 Appendix A Problem Statement and Related Work

2 **Episodic Training** A common approach is to match the training and evaluation conditions by
3 learning on C_{train} in an episodic manner, called *learning episodes* [21]. Note that training on support
4 set examples during episode evaluation is distinct from training on C_{train} . Many metric learners and
5 optimization-based learners use this training method, including Matching Networks [22], Prototypical
6 Networks [17], Relation Networks [18], and MAML [5].

7 **Non-episodic Baselines** Inspired by the transfer learning paradigm of pre-training and fine-tuning,
8 a natural non-episodic approach is to train a classifier on all examples in C_{train} at once. After training,
9 the final classification layer is removed, and this neural network is used as an embedding function f
10 that maps images \mathbf{x}_i to $x_i \in \mathbb{R}$ feature representations, including those from novel classes. It then
11 fine-tunes the final classifier layer using support set examples from the novel classes. The models
12 are a function of the parameters of a softmax layer, $\theta \in \mathbb{R}^d$. The softmax layer is formulated as the
13 similarity between image feature embeddings and the classifier parameters where θ_j is the parameters
14 for the j^{th} class, sim is the cosine similarity function.

$$p(y_i|x_i;\theta) = \frac{\exp(sim(x_i, \theta_{y_i}))}{\sum_{y' \in \mathcal{Y}} \exp(sim(x_i, \theta_{y'}))} \quad (1)$$

15 A.1 Related work

16 **Few-Shot Learning** Canonical approaches to few-shot learning include memory-based [7, 8, 13],
17 metric learning [15, 17, 18, 22], and optimization-based methods [5, 16]. However, recent studies
18 have shown that simple baseline learning techniques (i.e. simply training a backbone, then fine-tuning
19 the output layer on a few labeled examples) outperform or match performance of many meta-learning
20 methods [2, 4], prompting a closer look at the tasks [21] and contexts in which meta-learning is
21 helpful for few-shot learning [14, 20].

22 **Few-Shot Learning with Graphs** Beyond the canonical few-shot literature, studies have explored
23 learning GNNs over episodes as partially observed graphical models [6] and using GCNs to transfer
24 knowledge of semantic labels and categorical relationships to unseen classes in zero-shot learning
25 [23]. Recently, Chen et al. presented a knowledge graph transfer network (KGTN), which uses
26 a Gated Graph Neural Network (GGNN) to propagate information from base categories to novel
27 categories for few-shot learning [1]. Other works use domain knowledge graphs to provide task
28 specific customization [19], and propagate prototypes [10, 11]. However, these models have highly
29 complex architectures and consist of multiple sub-modules that all seem to impact performance.

30 Appendix B Experimental Setup

31 B.1 Mini-ImageNet

32 **Dataset** The Mini-ImageNet dataset is a subset of ILSVRC-2012 [3]. The classes are randomly
33 split into 64, 16 and 20 classes for meta-training, meta-validation, and meta-testing respectively. Each
34 class contains 600 images. We use the commonly-used split proposed in [22].

35 **Training details** We pre-train the feature extractor on $\mathcal{C}_{\text{train}}$ using the method proposed by [12].
36 Activations in the penultimate layer are pre-computed and saved as feature embeddings of 640
37 dimensions to simplify the fine-tuning process. For an N -way K -shot problem, we sample N novel
38 classes per episode, sample K support examples from those classes, and sample 15 query examples.
39 During pre-training and meta-training stages, input images are normalized using the mean and
40 standard-deviation computed on ILSVRC-2012. We apply standard data augmentation including
41 random crop, left-right flip, and color jitter in both the training or meta-training stage. We use
42 ResNet-18, ResNet-50 [9], and WRN-28-10 [24] for our backbone architectures. For pre-training
43 WRN-28-10, we follow the original hyperparameters and training procedures for S2M2_R [12]. For
44 meta-training ResNet-18, we follow the hyperparameters from [2]. At evaluation time, we choose
45 hyperparameters based on performance on the meta-validation set. Some implementation details are
46 adjusted for each method. Specifically, for ProtoNet and LEO, we include base examples during an
47 additional adaptation step per class. We show that these alterations have a minimal contribution to
48 performance in Appendix C.

49 B.2 ImageNet-FS

50 **Dataset** In the ImageNet-FS benchmark task, the 1000 ILSVRC-2012 categories are split into 389
51 base categories and 611 novel categories. From these, 193 of the base categories and 300 of the novel
52 categories are used during cross-validation and the remaining 196 base categories and 311 novel
53 categories are used for the final evaluation. Each base category has around 1,280 training images and
54 50 test images.

55 **Training details** We follow the procedure by [8] to pre-train the ResNet-50 feature extractor, and
56 adopt the Square Gradient Magnitude loss to regularize representation learning, which we scale by
57 0.005. The model is trained using the SGD algorithm with a batch size of 256, momentum of 0.9
58 and weight decay of 0.0005. The learning rate is initialized as 0.1 and is divided by 10 for every 30
59 epochs. During fine-tuning, we train for 10,000 iterations using the SGD algorithm with a batch size
60 of 256, momentum of 0.9, weight decay of 0.005, and learning rate of 0.01.

61 B.3 Label Graph

62 **WordNet ontology** ImageNet comprises of 82,115 ‘synsets’, which are based on the WordNet
63 ontology. For both the Mini-ImageNet and ImageNet-FS experiments, we first choose the synsets
64 corresponding to the output classes of each task – 100 for Mini-ImageNet and 1000 for ImageNet-FS.
65 ImageNet provides IS-A relationships over the synsets, defining a DAG over the classes. We only
66 consider the sub-graph consisting of the chosen classes and their ancestors. The classes are all leaves
67 of the DAG.

68 **Training details** The hyperparameter settings used for the node2vec-based graph regularization
69 objective are in line with typical values. For all experiments, we set $p = 1$, $q = 1$ and temperature
70 $T = 2$. We set the batch size to 128 for Mini-ImageNet and 256 for ImageNet-FS. Empirically, we
71 find that setting the regularization λ scaling higher for lower shots results in better performance, and
72 set $\lambda = 5, 3, 1$ for 1-, 2-, and 5-shot tasks respectively.

73 Appendix C Ablations

74 C.1 Mini-ImageNet Ablations

75 C.1.1 Model re-implementations with adaptation

76 For episodically-evaluated few-shot models, it is common practice to disregard base classes during
 77 evaluation. To implement graph regularization, we include both base and novel classes during test
 78 time and perform a further adaptation step per task. We show that the boost in performance is not due
 79 to these modifications.

Table 1: Validation of baseline model modifications.

Model	Backbone	1-shot	5-shot
ProtoNet	ResNet-18	54.16 ± 0.82	73.68 ± 0.65
ProtoNet (adaptation) [†]	ResNet-18	54.86 ± 0.73	74.14 ± 0.50
ProtoNet (adaptation) + Graph (Ours)	ResNet-18	55.47 ± 0.73	74.56 ± 0.49
LEO [†]	WRN 28-10	58.22 ± 0.09	74.46 ± 0.19
LEO (adaptation)	WRN 28-10	57.85 ± 0.20	74.25 ± 0.17
LEO (adaptation) + Graph (Ours)	WRN 28-10	60.93 ± 0.19	76.33 ± 0.17

80 C.1.2 Finding good parameter initializations for novel classes

81 Recent works have shown that good parameter initialization is important for few-shot adaptations
 82 [14]. For example, Dhillon et al. [4] showed that initializing novel classifiers with the mean of the
 83 support set improves few-shot performance.

84 Here, we explore various methods of incorporating graph relations to improve parameter initialization
 85 for novel classes. We compare our proposed method with simpler methods to show that the our graph
 86 regularization method is boosting performance in a non-trivial manner. For each method, we keep the
 87 adaptation procedure the same, namely, the fine-tuning procedure described by Baseline++ [2].

88 We then vary parameter initialization using the following methods: (A) random initialization, (B)
 89 initializing novel classes with the weights of the closest training class in graph distance in the
 90 knowledge graph, (C) our method.

Table 2: Mini-Imagenet with different parameter initialization methods (in % measured over 600 evaluation iterations).

Model	Backbone	1-shot	5-shot
S2M2 _R + Init A [12]	WRN 28-10	64.93 ± 0.18	83.18 ± 0.11
S2M2 _R + Init B	WRN 28-10	65.50 ± 0.81	83.32 ± 0.57
S2M2_R + Init C	WRN 28-10	66.93 ± 0.65	83.35 ± 0.53

91 C.2 ImageNet-FS Ablations

92 Here, we justify our model design decisions by considering alternatives. We first probe the benefits
 93 of using random walk neighborhoods by defining $N(y)$ as only nodes that have direct edges with y
 94 (“child-parent loss”). We try separately learning label graph embeddings, and passing the information
 95 to the classifier layer via “soft target” classification loss (“Independent graph w/ soft targets”). Results
 96 show that computing the graph loss directly on the classifier parameters is important for performance.
 97 Finally, we show that the quality of the label graph affects performance by removing layers of internal
 98 nodes of the WordNet hierarchy, starting from the bottom-most nodes (“Remove last 5, 10 layers”).

Table 3: Imagenet-FS ablations. Experiment setups, in order from the top: our proposed method, using only child-parent edges, independently learning graph embeddings, removing 5 layers of the ImageNet hierarchy, and removing 10 layers of the ImageNet hierarchy.

Ablation	1-shot
Ours	61.09
Child-parent loss	56.78
Independent graph w/ soft targets	56.22
Remove last 5 layers	57.80
Remove last 10 layers	54.86

References

- [1] R. Chen, T. Chen, X. Hui, H. Wu, G. Li, and L. Lin. Knowledge graph transfer network for few-shot recognition. *arXiv preprint arXiv:1911.09579*, 2019.
- [2] W.-Y. Chen, Y.-C. Liu, Z. Kira, Y.-C. F. Wang, and J.-B. Huang. A closer look at few-shot classification. *arXiv preprint arXiv:1904.04232*, 2019.
- [3] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [4] G. S. Dhillon, P. Chaudhari, A. Ravichandran, and S. Soatto. A baseline for few-shot image classification. *arXiv preprint arXiv:1909.02729*, 2019.
- [5] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1126–1135. JMLR. org, 2017.
- [6] V. Garcia and J. Bruna. Few-shot learning with graph neural networks. *arXiv preprint arXiv:1711.04043*, 2017.
- [7] S. Gidaris and N. Komodakis. Dynamic few-shot visual learning without forgetting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4367–4375, 2018.
- [8] B. Hariharan and R. Girshick. Low-shot visual recognition by shrinking and hallucinating features. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3018–3027, 2017.
- [9] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [10] L. Liu, T. Zhou, G. Long, J. Jiang, L. Yao, and C. Zhang. Prototype propagation networks (ppn) for weakly-supervised few-shot learning on category graph. *arXiv preprint arXiv:1905.04042*, 2019.
- [11] L. Liu, T. Zhou, G. Long, J. Jiang, and C. Zhang. Learning to propagate for graph meta-learning. In *Advances in Neural Information Processing Systems*, pages 1039–1050, 2019.
- [12] P. Mangla, N. Kumari, A. Sinha, M. Singh, B. Krishnamurthy, and V. N. Balasubramanian. Charting the right manifold: Manifold mixup for few-shot learning. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 2218–2227, 2020.
- [13] S. Qiao, C. Liu, W. Shen, and A. L. Yuille. Few-shot image recognition by predicting parameters from activations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7229–7238, 2018.
- [14] A. Raghu, M. Raghu, S. Bengio, and O. Vinyals. Rapid learning or feature reuse? towards understanding the effectiveness of maml. *arXiv preprint arXiv:1909.09157*, 2019.

- 136 [15] S. Ravi and H. Larochelle. Optimization as a model for few-shot learning. 2016.
- 137 [16] A. A. Rusu, D. Rao, J. Sygnowski, O. Vinyals, R. Pascanu, S. Osindero, and R. Hadsell.
138 Meta-learning with latent embedding optimization. *arXiv preprint arXiv:1807.05960*, 2018.
- 139 [17] J. Snell, K. Swersky, and R. Zemel. Prototypical networks for few-shot learning. In *Advances*
140 *in neural information processing systems*, pages 4077–4087, 2017.
- 141 [18] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. H. Torr, and T. M. Hospedales. Learning to compare:
142 Relation network for few-shot learning. In *Proceedings of the IEEE Conference on Computer*
143 *Vision and Pattern Recognition*, pages 1199–1208, 2018.
- 144 [19] Q. Suo, J. Chou, W. Zhong, and A. Zhang. Tadanet: Task-adaptive network for graph-enriched
145 meta-learning. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowl-*
146 *edge Discovery & Data Mining*, pages 1789–1799, 2020.
- 147 [20] Y. Tian, Y. Wang, D. Krishnan, J. B. Tenenbaum, and P. Isola. Rethinking few-shot image
148 classification: a good embedding is all you need? *arXiv preprint arXiv:2003.11539*, 2020.
- 149 [21] E. Triantafillou, T. Zhu, V. Dumoulin, P. Lamblin, K. Xu, R. Goroshin, C. Gelada, K. Swersky,
150 P.-A. Manzagol, and H. Larochelle. Meta-dataset: A dataset of datasets for learning to learn
151 from few examples. *arXiv preprint arXiv:1903.03096*, 2019.
- 152 [22] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra, et al. Matching networks for one shot learning.
153 In *Advances in neural information processing systems*, pages 3630–3638, 2016.
- 154 [23] X. Wang, Y. Ye, and A. Gupta. Zero-shot recognition via semantic embeddings and knowledge
155 graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
156 pages 6857–6866, 2018.
- 157 [24] S. Zagoruyko and N. Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*,
158 2016.