
Infinite Mixture Prototypes for Few-Shot Learning

Kelsey R. Allen¹ Evan Shelhamer^{*2} Hanul Shin^{*1} Joshua B. Tenenbaum¹

Abstract

We propose infinite mixture prototypes to adaptively represent both simple and complex data distributions for few-shot learning. Our infinite mixture prototypes represent each class by a set of clusters, unlike existing prototypical methods that represent each class by a single cluster. By inferring the number of clusters, infinite mixture prototypes interpolate between nearest neighbor and prototypical representations, which improves accuracy and robustness in the few-shot regime. We show the importance of adaptive capacity for capturing complex data distributions such as alphabets, with 25% absolute accuracy improvements over prototypical networks, while still maintaining or improving accuracy on the standard Omniglot and mini-ImageNet benchmarks. In clustering labeled and unlabeled data by the same clustering rule, infinite mixture prototypes achieves state-of-the-art semi-supervised accuracy. As a further capability, we show that infinite mixture prototypes can perform purely unsupervised clustering, unlike existing prototypical methods.

1. Introduction

Few-shot classification is the problem of learning to recognize new classes from only a few examples of each class (Lake et al., 2015; Fei-Fei et al., 2006; Miller et al., 2000). This requires careful attention to generalization, since overfitting or underfitting to the sparsely available data is more likely. Nonparametric methods are well suited to this task, as they can model decision boundaries that more closely reflect the data distribution by using the data itself.

Two popular classes of nonparametric methods are nearest neighbor methods and prototypical methods. Nearest neighbor methods represent a class by storing all of its examples, and are high-capacity models that can capture complex distributions. Prototypical methods, such as Gaussian mixture

models, represent a class by the mean of its examples, and are low-capacity models that can robustly fit simple distributions. Neighbors and prototypes are thus two ends of a spectrum from complex to simple decision boundaries, and the choice of which to apply generally requires knowledge about the complexity of the distribution.

Adaptively modulating model capacity is thus an important problem, especially in few-shot learning where the complexity of individual tasks can differ. Several approaches exist to tackle this, such as choosing k for k -nearest neighbours, selecting the number of mixture components for Gaussian mixture models, or adjusting the bandwidth (Jones et al., 1996) for kernel density estimation (Parzen, 1962).

Infinite mixture modeling (Hjort et al., 2010) represents one way of unifying these approaches for adaptively setting capacity. By inferring the number of mixture components for a given class from the data, it is possible to span the spectrum from nearest neighbors to prototypical representations.

This is particularly important in few-shot learning, where both underfitting and overfitting are common problems, because current models are fixed in their capacity.

To give an example, consider the problems of character and alphabet recognition. Recognizing characters is fairly straightforward: each character looks alike, and can be represented as a single prototype (a uni-modal Gaussian distribution). Recognizing alphabets is more complex: the uni-modal distribution assumption could be violated, and a multi-modal approach could better capture the complexity of the distribution. Figure 1 shows a prototypical network embedding for alphabets with this very issue. Even though the embedding was optimized for uni-modality, the uni-modal assumption is not guaranteed on held-out data.

We therefore propose Infinite Mixture Prototypes (IMP) to represent a class as a set of clusters, with the number of clusters determined directly from the data. IMP learns a deep embedding while also adapting the model capacity based on the complexity of the embedded data. As a further benefit, the infinite mixture modeling approach can naturally incorporate unlabeled data. We accordingly extend IMP to semi-supervised few-shot learning, and even to fully-unsupervised clustering inference.

An alternative approach to IMP would be to learn a para-

^{*}Equal contribution ¹Brain and Cognitive Sciences, MIT, Cambridge, MA ²Computer Science, UC Berkeley, Berkeley, CA. Correspondence to: Kelsey R. Allen <krallen@mit.edu>.

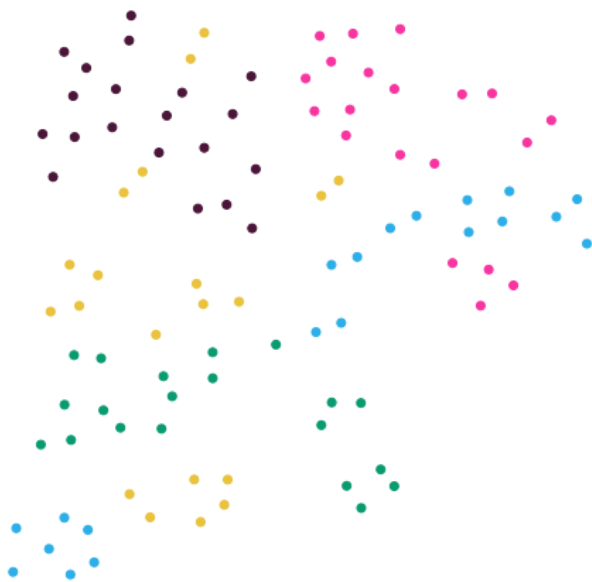


Figure 1. t-SNE visualization of the embedding from a prototypical network trained for alphabet recognition on Omniglot. Each point is a character colored by its alphabet label. The data distribution of each class is clearly not uni-modal, in violation of the modeling assumption for existing prototypical methods, causing errors. Our infinite mixture prototypes represent each class by a *set* of clusters, and infer their number, to better fit such distributions.

metric model. The decision boundary would then be linear in the embedding, which is more complex than uni-modal prototypes, but less complex than nearest neighbors. However, it may not be possible to find an embedding that yields a linear decision boundary. In practice, either a parametric method or uni-modal mixture model is sensitive to the choice of model capacity, and may not successfully learn complex classes such as Omniglot (Lake et al., 2015) alphabets. Instead, a higher-capacity nonparametric method like nearest neighbors can work far better. For simpler classes such as characters, a parametric model from a meta-learned initialization (Finn et al., 2017) or a prototypical network that assumes uni-modal data (Snell et al., 2017) suffice. Infinite mixture prototypes span these extremes, learning to adapt to both simple and complex classes.

In this paper, we extend prototypical networks from uni-modal to multi-modal clustering through infinite mixture modeling to give 25% improvement in accuracy for alphabet recognition (complex classes) while preserving accuracy on character recognition (simple classes) on Omniglot. In the semi-supervised setting infinite mixture prototypes are more accurate than semi-supervised prototypical networks. Infinite mixture modeling also allows for fully unsupervised clustering unlike existing prototypical methods. We demonstrate that the DP-means algorithm is suitable for

instantiating new clusters and that our novel extensions are necessary for best results in the few-shot regime. By end-to-end learning with infinite mixture modeling, IMP adapts its model capacity to simple or complex data distributions, shown by equal or better accuracy compared to neighbors and uni-modal prototypes in all experiments.

2. Background

For nonparametric representation learning methods, the model parameters are for the embedding function $h_\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$ that map an input point x into a feature. The embedding of point x is the M -dimensional feature vector from the embedding function. In deep models the parameters ϕ are the weights of a deep network, and the embedding is the output of the last layer of this network. (Such methods are still nonparametric because they represent decisions by the embedding of the data, and not parameters alone.)

2.1. Few-shot Classification

In few-shot classification we are given a *support* set $S = \{(x_1, y_1), \dots, (x_K, y_K)\}$ of K labeled points and a *query* set $Q = \{(x'_1, y'_1), \dots, (x'_{K'}, y'_{K'})\}$ of K' labeled points where each $x_i, x'_i \in \mathbb{R}^D$ is a D -dimensional feature vector and $y_i, y'_i \in \{1, \dots, N\}$ is the corresponding label. In the semi-supervised setting, y_i may not be provided for every point x_i . The support set is for learning while the query set is for inference: the few-shot classification problem is to recognize the class of the queries given the labeled supports.

Few-shot classification is commonly learned by constructing few-shot tasks from a large dataset and optimizing the model parameters on these tasks. Each task, comprised of the support and query sets, is called an *episode*. Episodes are drawn from a dataset by randomly sampling a subset of classes, sampling points from these classes, and then partitioning the points into supports and queries. The number of classes in the support is referred to as the “way” of the episode, and the number of examples of each class is referred to as the “shot” of the episode. Episodic optimization (Vinyals et al., 2016) iterates by making one episode and taking one update at a time. The update to the model parameters is defined by the task loss, which for classification could be the softmax cross-entropy loss.

2.2. Neighbors & Prototypes

Neighbors Nearest neighbors classification (Cover & Hart, 1967) assigns each query the label of the closest support. Neighbor methods are extremely simple but remarkably effective, because the classification is local and so they can fit complex data distributions. This generality comes at a computational cost, as the entire training set has to be stored and searched for inference. More fundamentally, there is a

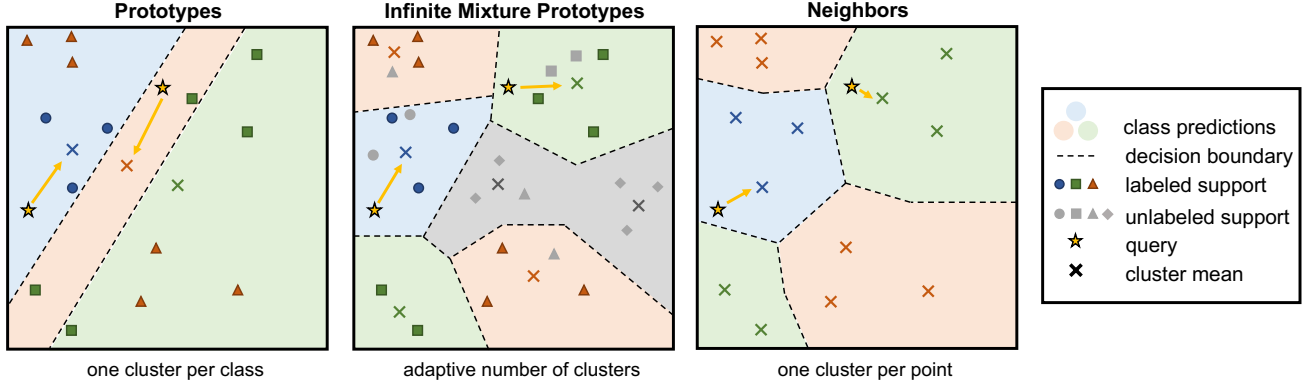


Figure 2. Our infinite mixture prototypes (IMP) method represents each class by a set of clusters, and infers the number of clusters from the data to adjust its modeling capacity. IMP is optimized end-to-end to cluster labeled and unlabeled data into *multi-modal* prototypes.

modeling issue: how should the distance metric to determine the “nearest” neighbor be defined?

Neighborhood component analysis (Goldberger et al., 2004) learns the distance metric by defining stochastic neighbors to make the classification decision differentiable. The metric is parameterized as a linear embedding A , and the probability of a point x_i having neighbor x_j is given by the softmax over Euclidean distances in the embedding:

$$p_{ij} = \frac{\exp(\|Ax_i - Ax_j\|^2)}{\sum_{k \neq j} \exp(\|Ax_i - Ax_k\|^2)}. \quad (1)$$

The probability that a point x_i is in class n is given by the sum of probabilities of neighbors in the class:

$$p_A(y = n | x_i) = \sum_{j: y_j = n} p_{ij}. \quad (2)$$

Stochastic neighbors naturally extend to a non-linear embedding trained by episodic optimization. Deep nearest neighbors classification therefore serves as a high-capacity nonparametric method for few-shot learning.

Prototypes Prototypical networks (Snell et al., 2017) form *prototypes* as the mean of the embedded support points in each class:

$$\mu_n = \frac{1}{|S_n|} \sum_{(x_i, y_i) \in S_n} h_\phi(x_i), \quad (3)$$

with S_n denoting the set of support points in class n . Paired with a distance $d(x_i, x_j)$, the prototypes classify a query point x' by the softmax over distances to the prototypes:

$$p_\phi(y' = n | x') = \frac{\exp(-d(h_\phi(x'), \mu_n))}{\sum_{n'} \exp(-d(h_\phi(x'), \mu_{n'}))}. \quad (4)$$

For the standard choice of the Euclidean distance function, the prototypes are equivalent to a Gaussian mixture model in the embedding with an identity covariance matrix.

ϕ is optimized by minimizing the negative log-probability of the true class of each query point by stochastic gradient descent over episodes. Prototypical networks therefore learn to create *uni-modal* class distributions for *fully-labeled* supports by representing each class by one cluster.

2.3. Infinite Mixture Modeling

Infinite mixture models (Hjort et al., 2010) do not require the number of mixture components to be known and finite. Instead, the number of components is inferred from data through Bayesian nonparametric methods (West et al., 1994; Rasmussen, 2000). In this way infinite mixture models adapt their capacity to steer between overfitting with high capacity and underfitting with low capacity.

The advantage of adaptivity is countered by the implementation and computational difficulties of Gibbs sampling and variational inference for infinite mixtures. To counter these issues, DP-means (Kulis & Jordan, 2012) is a deterministic, hard clustering algorithm derived via Bayesian nonparametrics for the Dirichlet process. DP-means iterates over the data points, computing each point’s minimum distance to all existing cluster means. If this distance is greater than a threshold λ , a new cluster is created with mean equal to the point. It optimizes a k -means-like objective for reconstruction error plus a penalty for making clusters.

λ , the distance threshold for creating a new cluster, is the sole hyperparameter for the algorithm. In deriving DP-means, Kulis & Jordan (2012) relate α , the concentration parameter for the Chinese restaurant process (CRP) (Aldous, 1985), to λ :

$$\lambda = 2\sigma \log\left(\frac{\alpha}{(1 + \frac{\rho}{\sigma})^{d/2}}\right) \quad (5)$$

where ρ is a measure of the standard deviation for the base distribution from which clusters are assumed to be drawn in the CRP. They then derive DP-means by connection to a Gibbs sampling procedure in the limit as σ approaches 0.

3. Infinite Mixture Prototypes (IMP)

Our infinite mixture prototypes (IMP) method pursues two approaches for adapting capacity: learning cluster variance to scale assignments, and multi-modal clustering to interpolate between neighbor and prototypical representations. This capability lets our model adapt its capacity to avoid underfitting, unlike existing prototypical models with fixed capacity. Figure 2 gives a schematic view of our multi-modal representation and how it differs from existing prototype and neighbor representations. Algorithm 1 expresses infinite mixture prototypes inference in pseudocode.

Within an episode, we initially cluster the support into class-wise means. Inference proceeds by iterating through all support points and computing their minimum distance to all existing clusters. If this distance exceeds a threshold λ , a new cluster is made with mean equal to that point. IMP then updates soft cluster assignments $z_{i,c}$ as the normalized Gaussian density for cluster membership. Finally, cluster means μ_c are computed by the weighted mean of their members. Since each class can have multiple clusters, we classify a query point x' by the softmax over distances to the closest cluster in each class n :

$$p_\phi(y' = n | x') = \frac{\exp(-d(h_\phi(x'), \mu_{c_n^*}))}{\sum_{n'} \exp(-d(h_\phi(x'), \mu_{c_{n'}^*}))} \quad (6)$$

with $c_n^* = \arg \min_{c:l_c=n} d(h_\phi(x'), \mu_c)$ indexing the clusters, where each cluster c has label l_c .

IMP optimizes the embedding parameters ϕ and cluster variances σ by stochastic gradient descent across episodes.

3.1. Adapting capacity by learning cluster variance σ

We learn the cluster variance σ to scale the assignment of support points to clusters. When σ is small, the effective distance is large and the closest points dominate, and when σ is large, the effective distance is small so farther points are more included. σ is differentiable, and therefore learned jointly with the embedding parameters ϕ . In practice, learning σ can improve the accuracy of prototypical networks, which we demonstrate by ablation in Table 1. For IMP, σ has a further role in creating new clusters.

3.2. Adapting capacity by multi-modal clustering

To create multi-modal prototypes, we extend the clustering algorithm DP-means (Kulis & Jordan, 2012) for compatibility with classification and end-to-end optimization. For classification, we distinguish labeled and unlabeled clusters, and incorporate labels into the point-cluster distance calculation. For end-to-end optimization, we soften cluster assignment, propose a scheme to select λ , and mask the classification loss to encourage multi-modality.

Indirect optimization of λ While λ is non-differentiable,

Algorithm 1 IMP: support prototypes and query inference

Require: supports $(x_1, y_1), \dots, (x_K, y_K)$ and queries $x'_1, \dots, x'_{K'}$
Return: clusters (μ_c, l_c, σ_c) and query classifications $p(y' | x')$

1. Init. each cluster μ_c with label l_c and $\sigma_c = \sigma_l$ as class-wise means of the supports, and C as the number of classes
 2. Estimate λ as in Equation 5
 3. Infer the number of clusters
 - for** each point x_i **do**
 - for** c in $\{1, \dots, C\}$ **do**
 - $d_{i,c} = \begin{cases} \|h_\phi(x_i) - \mu_c\|^2 & \text{if } (x_i \text{ is labeled and } l_c = y_i) \\ & \text{or } x_i \text{ is unlabeled} \\ +\infty & \text{otherwise} \end{cases}$
 - end for**
 - If $\min_c d_{i,c} > \lambda$: set $C = C + 1$, $\mu_C = h_\phi(x_i)$, $l_C = y_i$, $\sigma_C = \{\sigma_l \text{ if } x_i \text{ labeled, } \sigma_u \text{ otherwise}\}$.
 - end for**
 4. Assign supports to clusters by $z_{i,c} = \frac{\mathcal{N}(h_\phi(x_i); \mu_c, \sigma_c)}{\sum_c \mathcal{N}(h_\phi(x_i); \mu_c, \sigma_c)}$
 5. For each cluster c , compute mean $\mu_c = \frac{\sum_i z_{i,c} h_\phi(x_i)}{\sum_i z_{i,c}}$
 6. Classify queries by Equation 6
-

we propose an indirect optimization of the effective threshold for creating a new cluster. Based on Equation 5, λ depends on the concentration hyperparameter α , a measure of standard deviation in the prior ρ , and the cluster variance σ . α remains a hyperparameter, but with lessened effect. We estimate ρ as the variance between prototypes in each episode. As noted, σ is differentiable, so we learn it.

We model separate variances for labeled and unlabeled clusters, σ_l and σ_u respectively, in order to capture differences in uncertainty between labeled and unlabeled data. In the fully-supervised setting, λ is estimated from σ_l , and in the semi-supervised setting λ is estimated from the mean of σ_l and σ_u . In summary, learning the cluster variances σ affects IMP by scaling the distances between points and clusters, and through its role in our episodic estimation of λ .

Multi-modal loss We optimize all models with the cross-entropy loss. For the multi-modal methods (nearest neighbors and IMP), we mask the loss to only include the closest neighbor/cluster for each class, in the same manner as inference. That is, for a class n , we find the most likely cluster $c_n^* \leftarrow \arg \max_{c:l_c=n} \log p(h_\phi(x) | \mu_c, \sigma_c)$ and then take the loss over the queries in the class (Q_n):

$$J = \frac{1}{|Q_n|} \sum_{x \in Q_n} \left[-\log p(h_\phi(x) | \mu_{c_n^*}, \sigma_{c_n^*}) + \log \sum_{n' \neq n} p(h_\phi(x) | \mu_{c_{n'}^*}, \sigma_{c_{n'}^*}) \right].$$

Taking the loss for the closest clusters avoids over-penalizing multi-modality in the embedding, as taking the loss over all the clusters would. We found that masking improves the few-shot accuracy of these methods over other losses that incorporate all clusters.

Table 1. Multi-modal clustering and learning cluster variances on fully-supervised 10-way, 10-shot Omniglot alphabet recognition and 5-way, 5-shot mini-ImageNet. Scaling distances with the learned variance gives a small improvement and multi-modal clustering gives a further improvement.

METHOD	σ	MULTI-MODAL	ALPH. ACC.	MINI. ACC.
PROTOTYPES	-	-	65.2 \pm 0.6	66.1 \pm 0.6
PROTOTYPES	✓	-	65.2 \pm 0.6	67.2 \pm 0.5
IMP (OURS)	✓	✓	92.0 \pm 0.1	68.1 \pm 0.8

Table 2. Learning labeled cluster variance σ_l and unlabeled cluster variance σ_u on semi-supervised 5-way, 1-shot Omniglot and mini-ImageNet with 5 unlabeled points per class and 5 distractors (see Section 4). Learning σ_l, σ_u is better than learning a tied σ for labeled and unlabeled clusters.

METHOD	σ	OMNI. ACC.	MINI. ACC.
TIED	σ	93.5 \pm 0.3	48.6 \pm 0.4
IMP (OURS)	σ_l, σ_u	98.9 \pm 0.1	49.6 \pm 0.8

3.3. Ablations and Alternatives

We ablate our episodic and end-to-end extensions of DP-means to validate their importance for few-shot learning. Learning and performing inference with IMP is more robust to different choices of λ than simply using DP-means during inference (Figure 3). Multi-modality and learned variance make their own contributions to accuracy (Table 1). Learning separate σ_l, σ_u , for labeled and unlabeled clusters respectively, is more accurate than learning a shared σ for all clusters (Table 2). For full details of the datasets and settings in these ablations, refer to Section 4.

In principle, IMP’s clustering can be iterated multiple times during training and inference. However, we found that one clustering iteration suffices. Two iterations during training had no effect on accuracy, and even 100 iterations during inference still had no effect on accuracy, showing that the clustering is stable.

Alternative Algorithms DP-means was derived through the limit of a Gibbs sampler as the variance approaches 0, and so it does hard assignment of points to clusters. With hard assignment, it is still possible to learn the embedding parameters ϕ end-to-end by differentiating through the softmax over distances between query points and support clusters as in Equation 4. However, hard assignment of labeled and unlabeled data is harmful in our experiments, especially early on in training (see supplement).

When reintroducing variance into multi-modal clustering as we do, a natural approach would be to reconsider the Gibbs sampler for the CRP (West et al., 1994; Neal, 2000) from

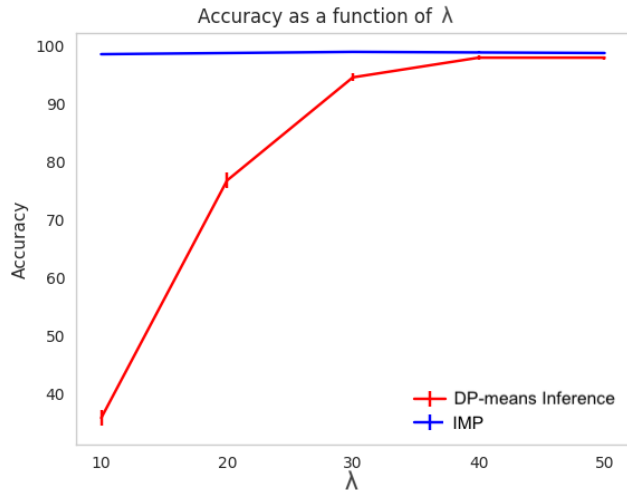


Figure 3. Learning and inference with IMP is more accurate and robust than DP-means inference on a prototypical network embedding alone. This plot shows the accuracy for the standard benchmark of semi-supervised 5-way, 1-shot Omniglot for different choices of the distance threshold λ for creating a new cluster.

which DP-means was derived, or other Dirichlet process inference methods such as expectation maximization (Kimura et al., 2013). These alternatives are less accurate in our experiments, mainly as a result of the CRP prior’s “rich get richer” dynamics, which prefers clusters with more assignments (leading to accuracy drops of 5–10%). This is especially problematic early in training, when unlabeled points are often incorrectly assigned. The supplement includes derivations and experiments regarding these multi-modal clustering alternatives.

4. Experiments

We experimentally show that infinite mixture prototypes are more accurate and more general than uni-modal prototypes.

We control for architecture and optimization by comparing methods with the same base architecture of Vinyals et al. (2016) and same episodic optimization settings of Snell et al. (2017). For further implementation details see Appendix A.1 of the supplement. All code for our method and baselines will be released.

We consider two widely-used datasets for few-shot learning:

Omniglot (Lake et al., 2015) is a dataset of 1,623 handwritten characters from 50 alphabets. There are 20 examples of each character, where the images are resized to 28x28 pixels and each image is rotated by multiples of 90°. This gives 6,492 classes in total, which are then split into 4,112 training classes, 688 validation classes, and 1,692 test classes.

mini-ImageNet (Vinyals et al., 2016) is a reduced ver-

Table 3. Alphabet and character recognition accuracy on Omniglot. Alphabets have more complex, multi-modal data distributions while characters have simpler, uni-modal data distributions. IMP improves accuracy for multi-modal alphabet classes, preserves accuracy for uni-modal character classes (Chars), and generalizes better from super-classes to sub-classes.

TRAINING	TESTING	PROTOTYPES	IMP	NEIGHBORS
ALPHABET	ALPHABET (10-WAY 10-SHOT)	65.6±0.4	92.0±0.1	92.4±0.2
ALPHABET	CHARS (20-WAY 1-SHOT)	82.1±0.4	95.4±0.2	95.4±0.2
CHARS	CHARS (20-WAY 1-SHOT)	94.9±0.2	95.1±0.1	95.1±0.1

sion of the ILSVRC’12 dataset (Russakovsky et al., 2015), which contains 600 84x84 images for 100 classes randomly selected from the full dataset. We use the split from Ravi & Larochelle (2017) with 64/16/20 classes for train/val/test.

4.1. Accuracy and Generality of Multi-modal Clustering by Infinite Mixture Prototypes

Our experiments on Omniglot alphabets and characters show that multi-modal prototypes are significantly more accurate than uni-modal prototypes for recognizing complex classes (alphabets) and recover uni-modal prototypes as a special case for recognizing simple classes (characters). Multi-modal prototypes generalize better for super-class to sub-class transfer learning, improving accuracy when training on alphabets but testing on characters. By unifying the clustering of labeled and unlabeled data alike, our multi-modal prototypes also address fully unsupervised clustering, unlike prior prototypical network models that are undefined without labels.

We first show the importance of multi-modality for learning representations of multi-modal classes: Omniglot alphabets. For these experiments, we train for alphabet classification, using only the super-class labels. Episodes are constructed by sampling n alphabets, and n_c characters within each alphabet. 1 image of each character is randomly sampled for the support, with 5 examples of each character for the query. We refer to these episodes as n -way n_c -shot episodes. For training, we sample 10-way, 10-shot episodes.

For character testing, we provide 1 labeled image of 20 different characters in the support, and score the correct character classification of the queries. Note that both alphabet and character testing are on held-out alphabets and characters respectively.

As seen in Table 3, IMP substantially outperforms prototypical networks for both alphabet and character recognition from alphabet training. On 20-way 1-shot character recognition, IMP achieves 95.4% from alphabet supervision alone, slightly out-performing prototypical networks trained directly on character recognition (94.9%). By clustering each super-class into multiple modes, IMP is better able to generalize to sub-classes.

For a parametric alternative, we trained MAML (Finn et al., 2017) on alphabet recognition, with the same episode composition as IMP. MAML achieves only 61.9% accuracy on 10-way 10-shot alphabet recognition. This demonstrates that a parametric classifier of this capacity, with decisions that are linear in the embedding, is not enough to solve alphabet recognition—instead, multi-modality is necessary.

Table 4. Generalization to held-out characters on 10-way, 5-shot Omniglot alphabet recognition. 40% of the characters are kept for training and 60% held out for testing. IMP maintains accuracy on held-out characters, suggesting that multi-modal clustering is more robust to new and different sub-classes from the same super-class.

METHOD	TRAINING MODES	TESTING MODES	BOTH MODES
IMP (OURS) PROTOTYPES	99.0±0.1	94.9±0.2	96.6±0.2
	92.4±0.3	77.7±0.4	82.9±0.4

To further examine generalization, we consider holding out character sub-classes during alphabet super-class training. In this experiment the training and testing alphabets are the same, but the characters within each alphabet are divided into training (40%) and testing (60%) splits. We compare alphabet recognition accuracy using training characters, testing characters, and all characters to measure generalization to held-out modes in Table 4. While prototypical networks achieve good accuracy on training modes, their accuracy drops 16% relative on testing modes, and still drops 10% relative on the combination of both modes. The reduced accuracy of prototypical networks on held-out modes indicates that uni-modality is not maintained on unseen characters even when they are from the same alphabets. IMP accuracy drops less than 5% relative from training to testing modes and both modes, showing that multi-modal clustering generalizes better to unseen data.

Fully Unsupervised Clustering IMP is able to do fully unsupervised clustering via multi-modality. Prototypical networks (Snell et al., 2017) and semi-supervised prototypical networks (Ren et al., 2018) are undefined without labeled data during testing because the number of clusters is defined by the number of classes.

For this unsupervised clustering setting, we use the models

Table 5. Unsupervised clustering of unseen Omniglot characters by IMP. Learning with IMP makes substantially purer clusters than DP-means inference on a prototypical network embedding, showing that the full method is necessary for best results.

METHOD	METRIC	10-WAY	100-WAY	200-WAY
IMP	PURITY	0.97	0.90	0.91
DP-MEANS		0.91	0.73	0.71
IMP	NMI	0.97	0.95	0.94
DP-MEANS		0.89	0.88	0.87
IMP	AMI	0.92	0.81	0.70
DP-MEANS		0.76	0.58	0.51

that were optimized for alphabet recognition. For testing, we randomly sample 5 examples of n character classes from the test set without labels.

IMP handles labeled and unlabeled data by the same clustering rule, infers the number of clusters as needed, and achieves good results under the standard clustering metrics of purity, and normalized/adjusted mutual information (NMI/AMI). We examine IMP’s clustering quality on purely unlabeled data in Table 5. IMP maintains strong performance across a large number of unlabeled clusters, without knowing the number of classes in advance, and without having seen any examples from the classes during training.

As a baseline, we evaluate multi-modal inference by DP-means (Kulis & Jordan, 2012) on the embedding from a prototypical network with the same architecture and training data as IMP. We cross-validate the cluster threshold λ on validation episodes for each setting, choosing by AMI.

4.2. Few-Shot Classification Benchmarks

We evaluate IMP on the standard few-shot classification benchmarks of Omniglot and mini-ImageNet in the fully-supervised and semi-supervised regimes.

We consider five strong fully-supervised baselines trained on 100% of the data. We compare to three parametric methods, MAML (Finn et al., 2017), Reptile (Nichol & Schulman, 2018), and few-shot graph networks (Garcia & Bruna, 2018), as well as three nonparametric methods, nearest neighbors, prototypical networks (Snell et al., 2017), and the memory-based model of Kaiser et al. (2017).

Fully-supervised results are reported in Table 6. In this setting, we evaluate IMP in the standard episodic protocol of few-shot learning: shot and way are fixed and classes are balanced within an episode. IMP learns to recover uni-modal clustering as a special case, matching or out-performing prototypical networks when the classes are uni-modal.

In the semi-supervised setting of labeled and unlabeled examples we follow Ren et al. (2018). We take only 40% of the data as labeled for both supports and queries while

the rest of the data is included as unlabeled examples. The unlabeled data is then incorporated into episodes as (1) within-support examples that allow for semi-supervised refinement of the support classes or (2) *distractors* which lie in the complement of the support classes. Semi-supervised episodes augment the fully supervised n -way, k -shot support with 5 unlabeled examples for each of the n classes and include 5 more distractor classes with 5 unlabeled instances each. The query set still contains only support classes.

Semi-supervised results are reported in Table 7. We train and test IMP, existing prototypical methods, and nearest neighbors in this setting. Semi-supervised prototypical networks (Ren et al., 2018) incorporate unlabeled data by soft k -means clustering (of their three comparable variants, we report “Soft k -Means+Cluster” results). Prototypical networks (Snell et al., 2017) and neighbors are simply trained on the 40% of the data with labels.

Through multi-modality, IMP clusters labeled and unlabeled data by a single rule. In particular this helps with the distractor distribution, which is in fact more diffuse and multi-modal by comprising several different classes.

The results reported on these benchmarks are for models trained and tested with n -way episodes. This is to equalize comparison across methods¹.

5. Related Work

Prototypes Prototypical networks (Snell et al., 2017) and semi-supervised prototypical networks (Ren et al., 2018) are the most closely related to our work. Prototypical networks simply and efficiently represent each class by its mean in a learned embedding. They assume that the data is fully labeled and uni-modal in the embedding. Ren et al. (2018) extend prototypes to the semi-supervised setting by refining prototypes through soft k -means clustering of the unlabeled data. They assume that the data is at least partially labeled and retain the uni-modality assumption. Both Snell et al. (2017) and Ren et al. (2018) are limited to one cluster per class. Mensink et al. (2013) represent classes by the mean of their examples in a linear embedding to incorporate new classes into large-scale classifiers without re-training. They extend their approach to represent classes by multiple prototypes, but the number of prototypes per class is fixed and hand-tuned, and their approach does not incorporate unlabeled data. We define a more general and adaptive approach through infinite mixture modeling that extends prototypi-

¹ Snell et al. (2017) train at higher way than testing and report a boost in accuracy. We find that this boost is somewhat illusory, and at least partially explained away by controlling for the number of gradients per update. We show this by experiment through the use of gradient accumulation in Appendix A.2 of the supplement. (For completeness, we confirmed that our implementation of prototypical networks reproduces reported results at higher way.)

Table 6. Fully-supervised few-shot accuracy using 100% of the labeled data. IMP performs equal to or better than prototypical networks (Snell et al., 2017). Although IMP is more general, it can still recover uni-modal clustering as a special case.

Method	Omniglot				mini-ImageNet	
	5-WAY		20-WAY		5-WAY	
	1-SHOT	5-SHOT	1-SHOT	5-SHOT	1-SHOT	5-SHOT
IMP (OURS)	98.4±0.3	99.5±0.1	95.0±0.1	98.6±0.1	49.6±0.8	68.1±0.8
NEIGHBORS	98.4±0.3	99.4±0.1	95.0±0.1	98.3±0.1	49.6±0.8	59.4±1.0
SNELL ET AL. (2017)	98.2±0.3	99.6±0.1	94.9±0.2	98.6±0.1	47.0±0.8	66.1±0.7
FINN ET AL. (2017)	98.7±0.4	99.9±0.3	95.8±0.3	98.9±0.2	48.7±1.84	63.1±0.92
GARCIA & BRUNA (2018)	99.2	99.7	97.4	99	50.3	66.41
KAISER ET AL. (2017)	98.4	99.6	95	98.6	-	-

Table 7. Semi-supervised few-shot accuracy on 40% of the labeled data with 5 unlabeled examples per class and 5 distractor classes. The distractor classes are drawn from the complement of the support classes and are only included unlabeled. IMP achieves equal or better accuracy than semi-supervised prototypical networks (Ren et al., 2018).

Method	Omniglot				mini-ImageNet	
	5-WAY		20-WAY		5-WAY	
	1-SHOT	5-SHOT	1-SHOT	5-SHOT	1-SHOT	5-SHOT
IMP (OURS)	98.9 ± 0.1	99.4 ± 0.1	96.9 ± 0.2	98.3 ± 0.1	49.2 ± 0.7	64.7 ± 0.7
REN ET AL. (2018)	98.0 ± 0.1	99.3 ± 0.1	96.2 ± 0.1	98.2 ± 0.1	48.6 ± 0.6	63.0 ± 0.8
NEIGHBORS	97.9 ± 0.2	99.1 ± 0.1	93.8 ± 0.2	97.5 ± 0.1	47.9 ± 0.7	57.3 ± 0.8
SNELL ET AL. (2017)	97.8 ± 0.1	99.2 ± 0.1	93.4 ± 0.1	98.1 ± 0.1	45.1 ± 1.0	62.5 ± 0.5

cal networks to multi-modal clustering, with one or many clusters per class, of labeled and unlabeled data alike.

Metric Learning Learning a metric to measure a given notion of distance/similarity addresses recognition by retrieval: given an unlabeled example, find the closest labeled example. Kulis (2013) gives a general survey. The contrastive loss and siamese network architecture (Chopra et al., 2005; Hadsell et al., 2006) optimize an embedding for metric learning by pushing similar pairs together and pulling dissimilar pairs apart. Of particular note is research in face recognition, where a same/different retrieval metric is used for many-way classification (Schroff et al., 2015). Our approach is more aligned with metric learning by meta-learning (Koch, 2015; Vinyals et al., 2016; Snell et al., 2017; Garcia & Bruna, 2018). These approaches learn a distance function by directly optimizing the task loss, such as cross-entropy for classification, through episodic optimization (Vinyals et al., 2016) for each setting of way and shot. Unlike metric learning on either neighbors (Goldberger et al., 2004; Schroff et al., 2015) or prototypes (Snell et al., 2017; Ren et al., 2018), our method adaptively interpolates between neighbor and uni-modal prototype representation by deciding the number of modes during clustering.

Cognitive Theories of Categorization Our approach is inspired by the study of categorization in cognitive science. Exemplar theory (Nosofsky, 1986) represents a category by storing its examples. Prototype theory (Reed, 1972) represents a category by summarizing its examples, by for

instance taking their mean. Vanpaemel et al. (2005) recognize that exemplars and prototypes are two extremes, and define intermediate models that represent a category by several clusters in their varying abstraction model. However, they do not define how to choose the clusters or their number, nor do they consider representation learning. Griffiths et al. (2007) unify exemplar and prototype categorization through the hierarchical Dirichlet process to model the transition from prototypes to exemplars as more data is collected. They obtain good fits for human data, but do not consider representation learning.

6. Conclusion

We made a case for the importance of considering the complexity of the data distribution in the regime of few-shot learning. By incorporating infinite mixture modeling with deep metric learning, we developed infinite mixture prototypes, a method capable of adapting its model capacity to the given data. Our multi-modal extension of prototypical networks additionally allows for fully unsupervised inference, and the natural incorporation of semi-supervised data during learning. As few-shot learning is applied to increasingly challenging tasks, models with adaptive complexity will become more important. Future work will look at extending IMP to the life-long setting, as well as integrating multiple input modalities.

Acknowledgements

We gratefully acknowledge support from DARPA grant 6938423 and KA is supported by NSERC. We thank Trevor Darrell and Ghassen Jerfel for advice and helpful discussions.

References

- Aldous, D. J. Exchangeability and related topics. In *École d'Été de Probabilités de Saint-Flour XIII1983*, pp. 1–198. Springer, 1985.
- Chopra, S., Hadsell, R., and LeCun, Y. Learning a similarity metric discriminatively, with application to face verification. In *CVPR*, volume 1, pp. 539–546. IEEE, 2005.
- Cover, T. and Hart, P. Nearest neighbor pattern classification. *IEEE transactions on information theory*, 13(1):21–27, 1967.
- Fei-Fei, L., Fergus, R., and Perona, P. One-shot learning of object categories. *PAMI*, 2006.
- Finn, C., Abbeel, P., and Levine, S. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017.
- Garcia, V. and Bruna, J. Few-shot learning with graph neural networks. In *ICLR*, 2018.
- Goldberger, J., Roweis, S., Hinton, G., and Salakhutdinov, R. Neighbourhood components analysis. In *NIPS*, pp. 513–520, 2004.
- Griffiths, T. L., Canini, K. R., Sanborn, A. N., and Navarro, D. J. Unifying rational models of categorization via the hierarchical dirichlet process. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, 2007.
- Hadsell, R., Chopra, S., and LeCun, Y. Dimensionality reduction by learning an invariant mapping. In *CVPR*, volume 2, pp. 1735–1742. IEEE, 2006.
- Hjort, N. L., Holmes, C., Müller, P., and Walker, S. G. *Bayesian nonparametrics*, volume 28. Cambridge University Press, 2010.
- Jones, M. C., Marron, J. S., and Sheather, S. J. A brief survey of bandwidth selection for density estimation. *Journal of the American statistical association*, 91(433):401–407, 1996.
- Kaiser, L., Nachum, O., Roy, A., and Bengio, S. Learning to remember rare events. In *ICLR*, 2017.
- Kimura, T., Tokuda, T., Nakada, Y., Nokajima, T., Matsumoto, T., and Doucet, A. Expectation-maximization algorithms for inference in dirichlet processes mixture. *Pattern Analysis and Applications*, 16(1):55–67, 2013.
- Koch, G. Siamese neural networks for one-shot image recognition. In *NIPS Deep Learning Workshop*, 2015.
- Kulis, B. Metric learning: A survey. *Foundations and Trends in Machine Learning*, 5(4):287–364, 2013.
- Kulis, B. and Jordan, M. I. Revisiting k-means: New algorithms via bayesian nonparametrics. In *ICML*, 2012.
- Lake, B. M., Salakhutdinov, R., and Tenenbaum, J. B. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- Mensink, T., Verbeek, J., Perronnin, F., and Csurka, G. Distance-based image classification: Generalizing to new classes at near-zero cost. *IEEE transactions on pattern analysis and machine intelligence*, 35(11):2624–2637, 2013.
- Miller, E. G., Matsakis, N. E., and Viola, P. A. Learning from one example through shared densities on transforms. In *CVPR*, volume 1, pp. 464–471. IEEE, 2000.
- Neal, R. M. Markov chain sampling methods for dirichlet process mixture models. *Journal of computational and graphical statistics*, 9(2):249–265, 2000.
- Nichol, A. and Schulman, J. Reptile: a scalable metalearning algorithm. *arXiv preprint arXiv:1803.02999*, 2018.
- Nosofsky, R. M. Attention, similarity, and the identification–categorization relationship. *Journal of experimental psychology: General*, 115(1):39, 1986.
- Parzen, E. On estimation of a probability density function and mode. *The annals of mathematical statistics*, 33(3):1065–1076, 1962.
- Rasmussen, C. E. The infinite gaussian mixture model. In *Advances in neural information processing systems*, pp. 554–560, 2000.
- Ravi, S. and Larochelle, H. Optimization as a model for few-shot learning. In *ICLR*, 2017.
- Raykov, Y. P., Boukouvalas, A., Little, M. A., et al. Simple approximate map inference for dirichlet processes mixtures. *Electronic Journal of Statistics*, 10(2):3548–3578, 2016.
- Reed, S. K. Pattern recognition and categorization. *Cognitive psychology*, 3(3):382–407, 1972.

Ren, M., Triantafillou, E., Ravi, S., Snell, J., Swersky, K., Tenenbaum, J. B., Larochelle, H., and Zemel, R. S. Meta-learning for semi-supervised few-shot classification. In *ICLR*, 2018.

Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. Imagenet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015.

Schroff, F., Kalenichenko, D., and Philbin, J. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, pp. 815–823, 2015.

Snell, J., Swersky, K., and Zemel, R. Prototypical networks for few-shot learning. In *NIPS*, pp. 4080–4090, 2017.

Tieleman, T. and Hinton, G. Lecture 6.5-rmsprop, coursera: Neural networks for machine learning. *University of Toronto, Technical Report*, 2012.

Vanpaemel, W., Storms, G., and Ons, B. A varying abstraction model for categorization. In *Proceedings of the Annual Conference of the Cognitive Science Society*, volume 27, pp. 2277–2282, 2005.

Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al. Matching networks for one shot learning. In *NIPS*, pp. 3630–3638, 2016.

West, M., Mller, P., and Escobar, M. Hierarchical priors and mixture models, with application in regression and density estimation. *Aspects of Uncertainty: A Tribute to DV Lindley*, pp. 363–386, 1994.

A. Appendix

A.1. Implementation Details

For all few-shot experiments, we use the same base architecture as prototypical networks for the embedding network. It is composed of four convolutional blocks consisting of a 64-filter 3×3 convolution, a batch normalization layer, a ReLU nonlinearity, and a 2×2 max-pooling layer per block. This results in a 64-dimensional embedding vector for omniglot, and a 1600 dimensional embedding vector for mini-imagenet. Our models were trained via SGD with RMSProp (Tieleman & Hinton, 2012) with an α parameter of 0.9.

For Omniglot, the initial learning rate was set to 1e-3, and cut by a factor of two every 2,000 iterations, starting at 4,000 iterations. Optimization is stopped at 160,000 iterations. We use gradient accumulation and accumulate gradients over eight episodes before making an update when performing 5-way training. Both σ_l and σ_u are initialized to 5.0. σ_l is learned jointly during training while we found learning σ_u

on Omniglot to be unstable and so it is therefore fixed. α was set to 0.1.

For mini-ImageNet, the initial learning rate was set to 1e-3, then halved every 20,000 iterations, starting at 20,000 iterations. Optimization is stopped at 100,000 iterations. Both σ_u and σ_l are initialized to 15.0 and both are learned jointly. We found that on average, σ_l stabilized around 12, and σ_u stabilized around 25. α was set to 10^{-5} . Clusters were still regularly created even with such a small α .

A.2. Controlling for the Number of Gradients Taken During Optimization

Consider the gradient of the loss: it has the dimensions of shot \times way because every example has a derivative with respect to every class. In this manner, by default, the episode size determines the number of gradients in an update. Quantitatively, 20-way episodes accumulate 16 times as many gradients as 5-way episodes. By sampling 16 5-way episodes and accumulating the gradients to make an update, we achieve significantly better results, matching the results obtained with 20-way episodes within statistical significance in most settings. Note that agreement across conditions may not be perfectly exact because subtle adjustments to hyperparameters might be necessary. See Table 8 for the quantitative results of these control experiments.

A.3. Alternative Infinite Mixture Model Algorithms

Here we discuss two alternatives to IMP for performing inference in infinite mixture models. We will first discuss an approximation to a Gibbs sampler for estimating the MAP of a Chinese restaurant process (CRP) (Aldous, 1985). We will then discuss an expectation maximization procedure which maintains soft assignments throughout inference.

The generative model of the CRP consists of sampling assignments z_1, \dots, z_J which could take on cluster values $c = 1, \dots, C$ from the CRP prior with hyperparameter α , which controls the concentration of clusters, and number of cluster members N_c . Cluster parameters μ_c are sampled from a base distribution $H(\theta_0; \mu_0, \sigma_0)$, and instances x_j are then sampled from the associated Gaussian distribution $\mathcal{N}(\mu_{z_j}, \sigma_{z_j})$. θ consists of the means μ and sigmas σ .

The CRP generative model is defined as

$$p(z_{J+1} = c | z_{1:J}, \alpha) = \frac{N_c}{N + \alpha} \text{ for } c \in \{1 \dots C\} \text{ and}$$

$$p(z_{J+1} = C + 1 | z_{1:J}, \alpha) = \frac{\alpha}{N + \alpha}$$

for assignments z of examples x to clusters c , cluster counts N_c , and parameter α to control assignments to new clusters. N is the total number of examples observed so far.

One popular sampling procedure for parameter estimation is

Table 8. Results on Omniglot for different gradient accumulations. Bolded results are not significantly different from each other, showing that equalizing the number of gradients can equalize the accuracy.

SHOT	BATCH-WAY	EPISODE-WAY	5-WAY		20-WAY	
			1-SHOT	5-SHOT	1-SHOT	5-SHOT
1	20	20	98.5	99.6	95.0	98.8
1	20	5	98.3	99.5	94.8	98.6
1	5	5	97.7	99.4	92.1	98.0
5	20	20	97.8	99.6	93.2	98.6
5	20	5	97.9	99.6	92.9	98.5
5	5	5	96.8	99.4	89.8	97.7

Gibbs sampling (Neal, 2000). In Gibbs sampling, we draw from a conditional distribution on the cluster assignments until convergence. The conditional draws are:

$$p(z_{J+1} = c | z_{1:J}, \alpha) \propto \begin{cases} N_{c,-j} \int P(x_j | \theta) dH_{-j,c}(\theta) & \text{for } c \leq C \\ \alpha \int P(x_j | \theta) dH_0(\theta) & \text{for } c = C + 1 \end{cases} \quad (7)$$

For the case of a spherical Gaussian likelihood, let us define $\mathcal{N}_c = \mathcal{N}(x_i; \mu_c, \sigma)$ as the likelihood of assigning x_i to cluster c and $\mathcal{N}_0 = \mathcal{N}(x_i; \mu_0, \sigma + \sigma_0)$ as the likelihood of assigning x_i to a new cluster drawn from the base distribution (Gaussian with mean μ_0 and σ_0). We can then write:

$$\begin{aligned} p(z_i = c | \mu) &= \frac{N_{k,-n} \mathcal{N}_c}{\alpha \mathcal{N}_0 + \sum_{j=1}^C N_{j,-n} \mathcal{N}_j} \\ p(z_i = C + 1 | \mu) &= \frac{\alpha \mathcal{N}_0}{\alpha \mathcal{N}_0 + \sum_{j=1}^C N_{j,-n} \mathcal{N}_j} \\ p(\sigma_c | z) &= \frac{\sigma \sigma_0}{\sigma + \sigma_0 N_c} \\ p(\mu_c | z) &= \mathcal{N} \left(\mu_c; \frac{\sigma \mu_0 + \sigma_0 \sum_{i, z_i=c} x_i}{\sigma + \sigma_0 N_c}, \sigma_c \right) \end{aligned}$$

Unfortunately, because inference must be performed during every episode of our learning procedure, and there are many episodes, Gibbs sampling until convergence is impractical. We therefore use the approach from (Raykov et al., 2016) to approximate the procedure with a single pass over all data in the episode. This approximates the MAP by considering only the most probable cluster assignment for each data point, and updating cluster parameters based on these assignments. A full discussion is given in Raykov et al. (2016), and we include their method here for reference (Algorithm 2). While their method is fully-unsupervised, we employ a cross-entropy loss on the query points given the updated means and counts for the *labeled* clusters, for end-to-end optimization of classification, and initialize clusters with the class-wise means as in IMP.

Results for 5-way 1-shot Omniglot and mini-ImageNet are in Table 9. Unlabeled points are often incorrectly assigned

to the labeled clusters, which both reduces the variance of that cluster, and increases its likelihood via the prior. The hard assignments lead to unstable clustering, making learning substantially more challenging.

We additionally implemented a simple expectation maximization approach (Algorithm 3). Here we maintain soft assignments z throughout, and use the updates to the cluster means μ_c as in (Kimura et al., 2013). Our three main differences are to: 1. include labeled points for initialization; 2. instead of having a fixed truncation parameter T for the maximum number of available clusters, we instantiate new clusters when the probability of a new cluster exceeds a certain threshold ϵ ; 3. we do not estimate variances, as this led to very unstable results. Instead of estimating variances based on assignments, we use the same variance learning technique as IMP, which provides significant improvement. The best value of α was one for which no new clusters were created in both Omniglot and mini-ImageNet.

Table 9. Ablation experiments comparing different inference schemes for infinite mixture prototypes. Accuracies are for semi-supervised 5-way 1-shot episodes, with 5 unlabeled examples per class, and 5 distractors.

METHOD	OMNIGLOT	MINI-IMAGENET
MAP-DP (μ, σ)	70.0 \pm 0.1	UNSTABLE
EM	95.9 \pm 0.2	41.0 \pm 0.6
HARD DP-MEANS	98.0 \pm 0.2	45.2 \pm 1.0
IMP	99.0 \pm 0.1	49.6 \pm 0.6

We additionally tested the hypothesis that the CRP prior was leading to worse performance by ablating it. With the prior ablated, the EM approach improves to 48.6% accuracy on mini-ImageNet, and 98.0% accuracy on Omniglot. While this is still below IMP’s performance, this gives some explanation for why the EM inference procedure fails.

The experiments in this section examine the semi-supervised 5-way 1-shot setting, with 5 unlabeled examples of each character and 5 distractor classes (see Section 4.2 of the paper for more experimental detail). In this setting, there

is no effect of multi-modality in the labeled examples, and so any improvements by IMP are attributed to the way it clusters *unlabeled* data relative to these inference methods.

Algorithm 2 MAP-DP approach for inference. n_s is the number of labeled classes (way). $q(i, c)$ is $\log p(i, c)$, the joint probability of cluster C and assignment i . $\mathcal{N}(x; \mu, \sigma)$ is the Gaussian density. α is the concentration hyperparameter of the CRP.

```

initialize  $\{\mu_1, \dots, \mu_{n_s}\}$   $\triangleright$  Initialize a cluster for each labeled class by taking class-wise means
initialize  $\{\sigma_1, \dots, \sigma_{n_s}\}$   $\triangleright$  Initialize cluster variances based on equation 4.
initialize  $\{z_1, \dots, z_I\}$   $\triangleright$  Initialize cluster assignments for labeled data points. All unlabeled cluster assignments start at 0.
 $C = n_s$   $\triangleright$  Initialize number of clusters  $C$ 
 $\triangleright$  Begin clustering pass
for each example  $i$  do
  for each cluster  $c \in \{1, \dots, C\}$  do
     $N_c \leftarrow \sum_i z_{i,c}$ 
     $\sigma_c \leftarrow \frac{\sigma\sigma_0}{\sigma + \sigma_0 N_c}$ 
     $\mu_c \leftarrow \frac{\sigma\mu_0 + \sigma_0 \sum_i z_{i,c} h_\phi(x_i)}{\sigma + \sigma_0 N_c}$ 
    estimate  $q_{i,c} \propto \log(N_{c,-i}) + \log(\mathcal{N}(x_i; \mu_c, \sigma_c))$ 
  end for
  estimate  $q_{i,C+1} \propto \log(\alpha) + \log(\mathcal{N}_0(x_i; \mu_0, \sigma_0))$ 
   $z_i \leftarrow \operatorname{argmin}(q_{i,1}, \dots, q_{i,C+1})$ 
  if  $z_i = C + 1$  then
     $C \leftarrow C + 1$ 
  end if
end for

```

Algorithm 3 EM approach for inference. n_s is the number of labeled classes (way). $q(i, c)$ is $\log p(i, c)$, the joint probability of cluster C and assignment i . $\mathcal{N}(x; \mu, \sigma)$ is the Gaussian density. α is the concentration hyperparameter of the CRP. ϵ threshold for generating new cluster.

```

initialize  $\{\mu_1, \dots, \mu_{n_s}\}$   $\triangleright$  Initialize a cluster for each labeled class by taking class-wise means
initialize  $\{\sigma_1, \dots, \sigma_{n_s}\}$   $\triangleright$  Initialize cluster variances based on equation 4.
initialize  $\{z_1, \dots, z_I\}$   $\triangleright$  Initialize cluster assignments for labeled data points. All unlabeled cluster assignments start at 0.
 $C = n_s$   $\triangleright$  Initialize number of clusters  $C$ 
 $\triangleright$  Begin clustering pass
for each example  $i$  do
  for each cluster  $c \in \{1, \dots, C\}$  do
    estimate  $q_{i,c} \propto \log(N_{c,-i}) + \log(\mathcal{N}(x_i; \mu_c, \sigma_c))$ 
  end for
  estimate  $q_{i,C+1} \propto \log(\alpha) + \log(\mathcal{N}_0(x_i; \mu_0, \sigma_0))$ 
   $z_{i,c} \leftarrow \operatorname{softmax}(q_{i,1}, \dots, q_{i,C+1})$ 
  if  $z_{i,C+1} > \epsilon$  then
     $C \leftarrow C + 1$ 
     $\mu_C \sim \mathcal{N}(x_i, \mu_0, \sigma_0)$   $\triangleright$  Sample from the base distribution conditioned on the single observation  $x_i$ 
  end if
end for

```
