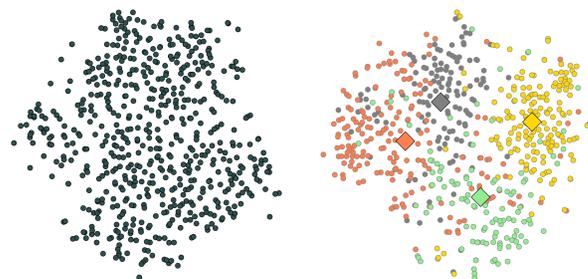


## MOTIVATION

- We introduce Mixture-based Feature Space Learning (MixtFSL) for obtaining a rich and robust feature representation
- our MixtFSL aims to learn a multimodal representation for the base classes

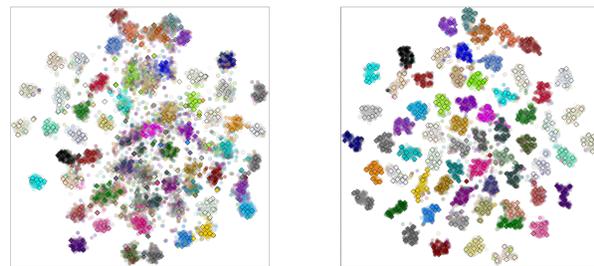


(a) without MixtFSL (b) our MixtFSL

- The idea is to learn both the *representation* and the *mixture model* jointly in an online manner

## MIXTFSL

- We present a robust two-stage scheme for training such a model.
- The training is done end-to-end in a fully differentiable fashion, without the need for an offline clustering method.

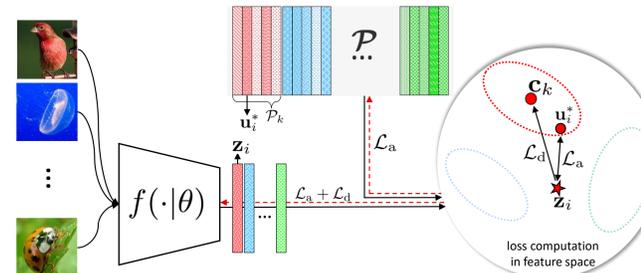


(a) initial training (b) progressive following

- We demonstrate, through an extensive experiments on four standard datasets and using four backbones, that our MixtFSL outperforms the state of the art in most of the cases tested.

## INITIAL TRAINING

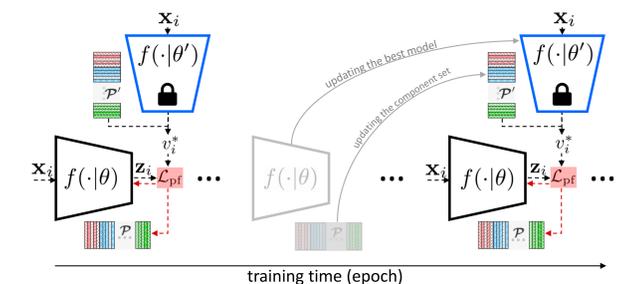
- The initial training of  $f(\cdot|\theta)$  and the learnable mixture model  $\mathcal{P}$  from the base class set  $\mathcal{X}^b$  is illustrated in figure bellow.



- Model parameters are updated using two losses: the “assignment” loss  $\mathcal{L}_a$ , which updates both the feature extractor and the mixture model such that feature vectors are assigned to their nearest mixture component; and the “diversity” loss  $\mathcal{L}_d$ , which updates the feature extractor to diversify the selection of components for a given class.

## PROGRESSIVE FOLLOWING

- The progressive following stage that aim to break the complex dynamic of simultaneously determining nearest components while training the representation  $f(\cdot|\theta)$  and mixture  $\mathcal{P}$ . The approach is shown in bellow.
- Using the “prime” notation ( $\theta'$  and  $\mathcal{P}'$  to specify the best feature extractor parameters and mixture component so far, resp.), the approach starts by taking a copy of  $f(\cdot|\theta')$  and  $\mathcal{P}'$ , and by using them to determine the nearest component of each training instance:



## RESULTS: MINIIMAGENET

- Evaluations of our MixtFSL on *mini-ImageNet* using different Conv4 and ResNet12.

Table 1. Evaluation on *miniImageNet* in 5-way. Bold/blue is best/second, and  $\pm$  is the 95% confidence intervals in 600 episodes.

Method	Backbone	1-shot	5-shot
ProtoNet [64]	Conv4	49.42 $\pm$ 0.78	68.20 $\pm$ 0.66
MAML [19]	Conv4	48.07 $\pm$ 1.75	63.15 $\pm$ 0.91
RelationNet [67]	Conv4	50.44 $\pm$ 0.82	65.32 $\pm$ 0.70
Baseline++ [8]	Conv4	48.24 $\pm$ 0.75	66.43 $\pm$ 0.63
IMP [2]	Conv4	49.60 $\pm$ 0.80	68.10 $\pm$ 0.80
MemoryNetwork [5]	Conv4	<b>53.37</b> $\pm$ 0.48	66.97 $\pm$ 0.35
Arcmax [1]	Conv4	51.90 $\pm$ 0.79	69.07 $\pm$ 0.59
Neg-Margin [41]	Conv4	<b>52.84</b> $\pm$ 0.76	<b>70.41</b> $\pm$ 0.66
MixtFSL (ours)	Conv4	52.82 $\pm$ 0.63	<b>70.67</b> $\pm$ 0.57
<hr/>			
DNS [62]	RN-12	62.64 $\pm$ 0.66	78.83 $\pm$ 0.45
Var.FSL [87]	RN-12	61.23 $\pm$ 0.26	77.69 $\pm$ 0.17
MTL [66]	RN-12	61.20 $\pm$ 1.80	75.50 $\pm$ 0.80
SNAIL [46]	RN-12	55.71 $\pm$ 0.99	68.88 $\pm$ 0.92
AdaResNet [48]	RN-12	56.88 $\pm$ 0.62	71.94 $\pm$ 0.57
TADAM [49]	RN-12	58.50 $\pm$ 0.30	76.70 $\pm$ 0.30
MetaOptNet [37]	RN-12	62.64 $\pm$ 0.61	78.63 $\pm$ 0.46
Simple [69]	RN-12	62.02 $\pm$ 0.63	79.64 $\pm$ 0.44
TapNet [83]	RN-12	61.65 $\pm$ 0.15	76.36 $\pm$ 0.10
Neg-Margin [41]	RN-12	<b>63.85</b> $\pm$ 0.76	<b>81.57</b> $\pm$ 0.56
MixtFSL (ours)	RN-12	<b>63.98</b> $\pm$ 0.79	<b>82.04</b> $\pm$ 0.49

## TIEREDIMAGENET AND FC100

- Evaluations of our MixtFSL on *tieredImageNet* using different ResNet12 and ResNet18.

Table 2. Evaluation on *tieredImageNet* and FC100 in 5-way classification. Bold/blue is best/second best, and  $\pm$  indicates the 95% confidence intervals over 600 episodes.

	Method	Backbone	1-shot	5-shot
tieredImageNet	DNS [62]	RN-12	66.22 $\pm$ 0.75	82.79 $\pm$ 0.48
	MetaOptNet [37]	RN-12	65.99 $\pm$ 0.72	81.56 $\pm$ 0.53
	Simple [69]	RN-12	<b>69.74</b> $\pm$ 0.72	<b>84.41</b> $\pm$ 0.55
	TapNet [83]	RN-12	63.08 $\pm$ 0.15	80.26 $\pm$ 0.12
	Arcmax* [1]	RN-12	68.02 $\pm$ 0.61	83.99 $\pm$ 0.62
	MixtFSL (ours)	RN-12	<b>70.97</b> $\pm$ 1.03	<b>86.16</b> $\pm$ 0.67
	Arcmax [1]	RN-18	<b>65.08</b> $\pm$ 0.19	<b>83.67</b> $\pm$ 0.51
	ProtoNet [64]	RN-18	61.23 $\pm$ 0.77	80.00 $\pm$ 0.55
	MixtFSL (ours)	RN-18	<b>68.61</b> $\pm$ 0.91	<b>84.08</b> $\pm$ 0.55
	TADAM [49]	RN-12	40.1 $\pm$ 0.40	<b>56.1</b> $\pm$ 0.40
FC100	MetaOptNet [37]	RN-12	41.1 $\pm$ 0.60	55.5 $\pm$ 0.60
	ProtoNet† [64]	RN-12	37.5 $\pm$ 0.60	52.5 $\pm$ 0.60
	MTL [66]	RN-12	<b>43.6</b> $\pm$ 1.80	55.4 $\pm$ 0.90
	MixtFSL (ours)	RN-12	<b>44.89</b> $\pm$ 0.63	<b>60.70</b> $\pm$ 0.67
	Arcmax [1]	RN-18	40.84 $\pm$ 0.71	57.02 $\pm$ 0.63
	MixtFSL (ours)	RN-18	<b>41.50</b> $\pm$ 0.67	<b>58.39</b> $\pm$ 0.62

\*our implementation †taken from [37]

## CUB AND CROSS-DOMAIN

- Evaluations of our MixtFSL on CUB in object recognition and cross-domain adaptation using ResNet18.

Table 3. Fine-grained and on cross-domain from *miniImageNet* to CUB evaluation in 5-way using ResNet-18. Bold/blue is best/second, and  $\pm$  is the 95% confidence intervals on 600 episodes.

Method	CUB		miniIN $\rightarrow$ CUB
	1-shot	5-shot	5-shot
GNN-LFT <sup>o</sup> [70]	51.51 $\pm$ 0.8	73.11 $\pm$ 0.7	–
Robust-20 [13]	58.67 $\pm$ 0.7	75.62 $\pm$ 0.5	–
RelationNet <sup>‡</sup> [67]	67.59 $\pm$ 1.0	82.75 $\pm$ 0.6	57.71 $\pm$ 0.7
MAML <sup>‡</sup> [18]	68.42 $\pm$ 1.0	83.47 $\pm$ 0.6	51.34 $\pm$ 0.7
ProtoNet <sup>‡</sup> [64]	71.88 $\pm$ 0.9	<b>86.64</b> $\pm$ 0.5	62.02 $\pm$ 0.7
Baseline++ [8]	67.02 $\pm$ 0.9	83.58 $\pm$ 0.5	64.38 $\pm$ 0.9
Arcmax [1]	71.37 $\pm$ 0.9	85.74 $\pm$ 0.5	64.93 $\pm$ 1.0
Neg-Margin [41]	<b>72.66</b> $\pm$ 0.9	<b>89.40</b> $\pm$ 0.4	<b>67.03</b> $\pm$ 0.8
MixtFSL (ours)	<b>73.94</b> $\pm$ 1.1	86.01 $\pm$ 0.5	<b>68.77</b> $\pm$ 0.9

<sup>‡</sup>taken from [68] <sup>o</sup>backbone is ResNet-10

## AN EXTENSION OF MIXTFSL

- Two changes are necessary to adapt our MixtFSL to exploit the “centroid alignment” of [1].
- First, we employ the learned mixture model  $\mathcal{P}$  to find the related base classes.
- Second, they used a classification layer  $\mathbf{W}$  in  $c(\mathbf{x}|\mathbf{W}) \equiv \mathbf{W}^\top f(\mathbf{x}|\theta)$  (followed by softmax).

Table 6. Comparison of our MixtFSL with alignment (MixtFSL-Align) in 5-way classification. Here, bold is the best performance.

	Method	Backbone	1-shot	5-shot
miniIN	Cent. Align.* [1]	RN-12	63.44 $\pm$ 0.67	80.96 $\pm$ 0.61
	MixtFSL-Align. (ours)	RN-12	<b>64.38</b> $\pm$ 0.73	<b>82.45</b> $\pm$ 0.62
	Cent. Align.* [1]	RN-18	59.85 $\pm$ 0.67	80.62 $\pm$ 0.72
	MixtFSL-Align. (ours)	RN-18	<b>60.44</b> $\pm$ 1.02	<b>81.76</b> $\pm$ 0.74
tieredImageNet	Cent. Align.* [1]	RN-12	71.08 $\pm$ 0.93	86.32 $\pm$ 0.66
	MixtFSL-Align. (ours)	RN-12	<b>71.83</b> $\pm$ 0.99	<b>88.20</b> $\pm$ 0.55
	Cent. Align.* [1]	RN-18	69.18 $\pm$ 0.86	<b>85.97</b> $\pm$ 0.51
	MixtFSL-Align. (ours)	RN-18	<b>69.82</b> $\pm$ 0.81	85.57 $\pm$ 0.60

\* our implementation

## ACKNOWLEDGEMENT