# MetaMix: Improved Meta-Learning with Interpolation-based Consistency Regularization

Yangbin Chen[1], Yun Ma[2], Tom Ko[3], Jianping Wang[1], Qing Li[2]

1. City University of Hong Kong  2. The Hong Kong Polytechnic University  3. Southern University of Science and Technology
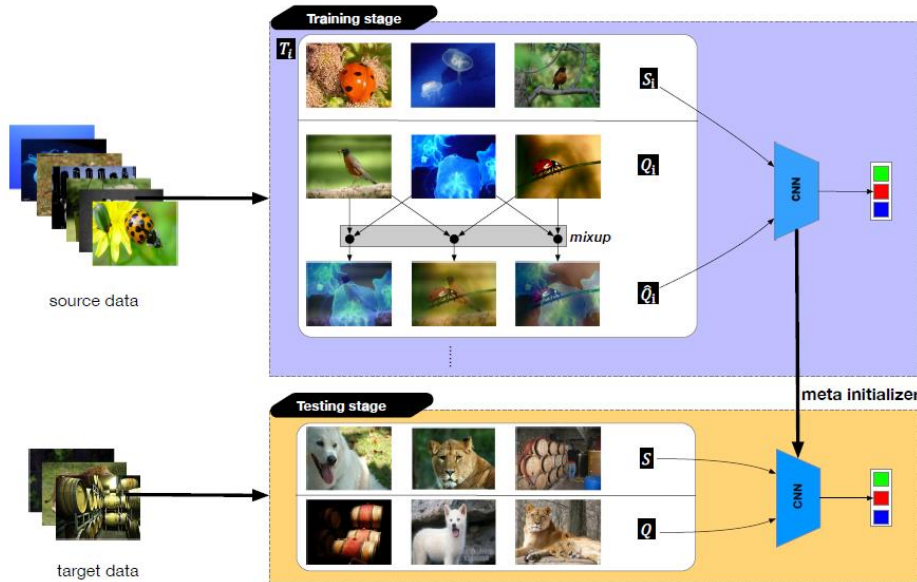
ICPR 2020

## BACKGROUND

- Few-Shot Learning (FSL) problem is a machine learning problem that learns with limited labelled data of the target tasks by incorporating external source data with a different distribution.
- Few-Shot Classification is a few-shot learning task defined as N-way, K-shot, where N is the number of classes in the target task and K is the number of labelled examples per class.
- Model-Agnostic Meta-Learning (MAML) and its variants aim to train a model, which can adapt quickly to any new tasks using only a few examples.

## MOTIVATION

- Conventional meta-learning algorithms face meta-overfitting problems, where the learned decision boundary stays too close to the limited labelled examples in few-shot classification tasks.
- The Empirical Risk Minimization (ERM) allows large neural networks to memorize (rather than generalize from) the training data.
- We aim to propose a regularization technique to solve the meta-overfitting problem.

## METHODOLOGY



**Algorithm 1 MetaMix with MAML**

**Require:** $p(\mathcal{T})$ : distribution over tasks
**Require:** $\mathcal{S}_i$ : support set; $\mathcal{Q}_i$ : query set
**Require:** $\alpha, \beta$ : learning rate
**Require:** $\check{\alpha}$ : Beta distribution parameter
**Require:** $mix_\lambda(a, b) = \lambda a + (1 - \lambda)b, \lambda \sim \mathbf{B}(\check{\alpha}, \check{\alpha})$

1: Randomly initialize model parameters $\theta$
2: **while** not done **do**
3:   Sample a batch of episodes $\mathcal{T}_i \sim p(\mathcal{T})$
4:   **for all** $\mathcal{T}_i$ **do**
5:     Sample a support set $\mathcal{S}_i = \{(x_j, y_j)\}_{j=1}^J$
6:     Evaluate $\nabla_\theta \mathcal{L}_{\mathcal{S}_i}(f_\theta)$ using $\mathcal{S}_i$ and $\mathcal{L}_{\mathcal{S}_i}(f_\theta)$
7:     Compute adapted parameters with gradient descent: $\theta'_i = \theta - \alpha \cdot \nabla_\theta \mathcal{L}_{\mathcal{S}_i}(f_\theta)$
8:     Sample a query set $\mathcal{Q}_i = \{(x_z, y_z)\}_{z=1}^Z$
9:     Randomly select pairs of examples $\{(x_m, y_m)\}_{m=1}^Z, \{(x_n, y_n)\}_{n=1}^Z$ from $\mathcal{Q}_i$
10:     $\hat{x}_z = mix_\lambda(x_m, x_n), \hat{y}_z = mix_\lambda(y_m, y_n)$
11:     Get new query set $\hat{\mathcal{Q}}_i = \{(\hat{x}_z, \hat{y}_z)\}_{z=1}^Z$
12:   **end for**
13:   Update $\theta \leftarrow \theta - \beta \cdot \nabla_\theta \sum_i \mathcal{L}_{\hat{\mathcal{Q}}_i}(f_{\theta'_i})$
14: **end while**

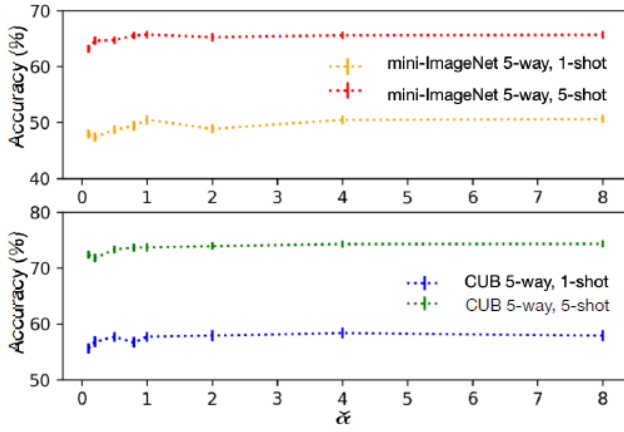## EXPERIMENT

- **Performance comparison of MetaMix and baseline approaches on 5-way classification tasks over three datasets**

| Models | mini-ImageNet | | CUB | | FC100 | |
| --- | --- | --- | --- | --- | --- | --- |
| | 1-shot | 5-shot | 1-shot | 5-shot | 1-shot | 5-shot |
| Matching Network | $50.47 \pm 0.80$ | $64.83 \pm 0.67$ | $57.70 \pm 0.87$ | $71.42 \pm 0.71$ | $36.97 \pm 0.67$ | $49.44 \pm 0.71$ |
| Prototypical Network | $49.33 \pm 0.82$ | $65.71 \pm 0.67$ | $51.34 \pm 0.86$ | $67.56 \pm 0.76$ | $36.83 \pm 0.69$ | $51.21 \pm 0.74$ |
| Relation Network | $50.48 \pm 0.80$ | $65.39 \pm 0.72$ | $59.47 \pm 0.96$ | $73.88 \pm 0.74$ | $36.40 \pm 0.69$ | $51.35 \pm 0.69$ |
| MAML | $48.18 \pm 0.78$ | $63.05 \pm 0.71$ | $54.32 \pm 0.91$ | $71.37 \pm 0.76$ | $35.96 \pm 0.71$ | $48.06 \pm 0.73$ |
| MetaMix+MAML | $\mathbf{50.51 \pm 0.86}$ | $\mathbf{65.73 \pm 0.72}$ | $\mathbf{57.70 \pm 0.92}$ | $\mathbf{73.66 \pm 0.74}$ | $\mathbf{37.09 \pm 0.74}$ | $\mathbf{49.31 \pm 0.72}$ |
| FOMAML | $45.22 \pm 0.77$ | $60.97 \pm 0.70$ | $53.12 \pm 0.93$ | $70.90 \pm 0.75$ | $34.97 \pm 0.70$ | $47.41 \pm 0.73$ |
| MetaMix+FOMAML | $\mathbf{47.78 \pm 0.77}$ | $\mathbf{63.55 \pm 0.70}$ | $\mathbf{54.81 \pm 0.97}$ | $\mathbf{72.90 \pm 0.74}$ | $\mathbf{36.48 \pm 0.67}$ | $\mathbf{49.48 \pm 0.71}$ |
| MetaSGD | $49.93 \pm 1.73$ | $64.01 \pm 0.90$ | $56.19 \pm 0.92$ | $69.14 \pm 0.75$ | $36.36 \pm 0.66$ | $49.96 \pm 0.72$ |
| MetaMix+MetaSGD | $\mathbf{50.60 \pm 1.80}$ | $\mathbf{64.47 \pm 0.88}$ | $\mathbf{57.64 \pm 0.88}$ | $\mathbf{70.50 \pm 0.70}$ | $\mathbf{37.44 \pm 0.71}$ | $\mathbf{51.41 \pm 0.69}$ |
| MTL | $61.37 \pm 0.82$ | $78.37 \pm 0.60$ | $71.90 \pm 0.86$ | $84.68 \pm 0.53$ | $42.17 \pm 0.79$ | $56.84 \pm 0.75$ |
| MetaMix+MTL | $\mathbf{62.74 \pm 0.82}$ | $\mathbf{79.11 \pm 0.58}$ | $\mathbf{73.04 \pm 0.86}$ | $\mathbf{86.10 \pm 0.50}$ | $\mathbf{43.58 \pm 0.73}$ | $\mathbf{58.27 \pm 0.73}$ |

- **Effect of Beta distribution**



- **Effect of mixup on different sets**

| Set(s) | mini-ImageNet | | CUB | |
|---|---|---|---|---|
| | 1-shot | 5-shot | 1-shot | 5-shot |
| Q | **50.51 ± 0.86** | **65.73 ± 0.72** | **57.70 ± 0.92** | **73.66 ± 0.74** |
| S | 44.03 ± 0.79 | 53.74 ± 0.81 | 49.12 ± 0.96 | 63.27 ± 0.89 |
| Q+S | 48.36 ± 0.81 | 64.06 ± 0.72 | 54.32 ± 0.93 | 70.30 ± 0.75 |
| w/o MetaMix | 48.18 ± 0.78 | 63.05 ± 0.71 | 54.32 ± 0.91 | 71.37 ± 0.76 |

- **Effect of size of training data**

| Models | mini-ImageNet | | CUB | | FC100 | |
|---|---|---|---|---|---|---|
| | 1-shot | 5-shot | 1-shot | 5-shot | 1-shot | 5-shot |
| MAML(100%) | 48.18 ± 0.78 | 63.05 ± 0.71 | 54.32 ± 0.91 | 71.37 ± 0.76 | 35.96 ± 0.71 | 48.06 ± 0.73 |
| MetaMix+MAML(100%) | **50.51 ± 0.86** | **65.73 ± 0.72** | **57.70 ± 0.92** | **73.66 ± 0.74** | **37.09 ± 0.74** | **49.31 ± 0.72** |
| MAML(50%) | 46.34 ± 0.82 | 60.47 ± 0.73 | 50.78 ± 0.86 | 65.60 ± 0.81 | 35.38 ± 0.71 | 47.93 ± 0.78 |
| MetaMix+MAML(50%) | **48.04 ± 0.79** | **63.52 ± 0.67** | **53.22 ± 0.91** | **70.13 ± 0.70** | **36.35 ± 0.74** | **48.11 ± 0.69** |



(a) 1-shot task on mini-imagenet
(b) 5-shot task on mini-imagenet
(c) 1-shot task on CUB
(d) 5-shot task on CUB

# CONTRIBUTIONS

- We propose MetaMix as a regularization technique, which can be integrated with many meta-learning algorithms, including MAML and its variants, and improve their performance.
- MetaMix with MAML-based algorithms perform more robust with the reduction of training data, compared with original MAML-based algorithms.
- MetaMix with MTL achieves state-of-the-art performance.