# Rethinking Generalization in Few-Shot Classification

*Hiller et al. NeurIPS 2022*

MLAI In-depth Seminar

# Plan for today

- Few-Shot Learning

- Existing Works

- Problem Define

- Method

- Results

- Conclusion & Discussion

# Few-Shot Learning

- Aims at classifying unlabeled samples (**query** set) into unseen classes given very few labeled samples (**support** set)

- Two main challenges
    - Unseen classes : *non-overlap between training and test classes*
    - Low-data problem : *very few labeled samples for the test unseen classes*
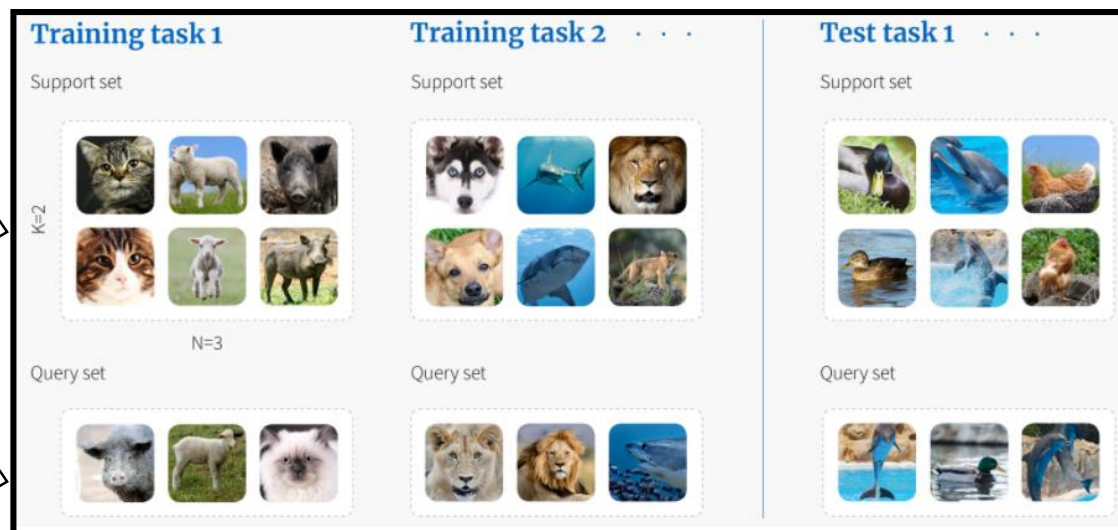
# Few-Shot Learning

- Aims at classifying unlabeled samples (**query** set) into unseen classes given very few labeled samples (**support** set)

- Two main challenges
  - Unseen classes : *non-overlap between training and test classes*
  - Low-data problem : *very few labeled samples for the test unseen classes*

- Problem Define

< Episodic training and testing >

Generally, Training set contains a large number of classes and correspoding labeled samples.

But, for matching training/test procedures, the dataset is constructed as a series of episodes.



$C_{train} \cap C_{test} = \emptyset$

$C_{train} = \{cat, goat, boar, dog, shark, lion\}$    $C_{test} = \{dug, dolphin, chicken\}$

Source:

# Few-Shot Learning

- Aims at classifying unlabeled samples (**query** set) into unseen classes given very few labeled samples (**support** set)
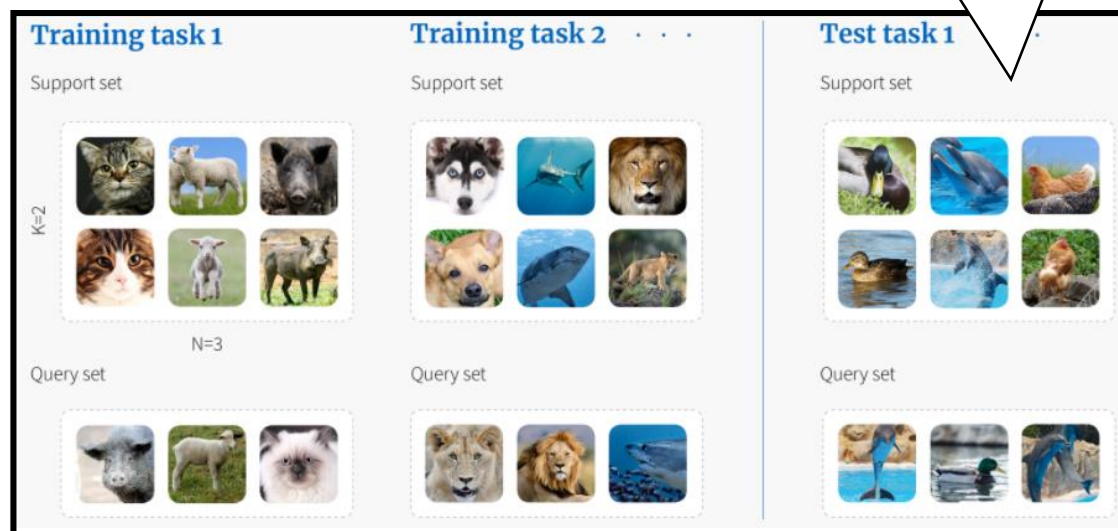
- Two main challenges
  - Unseen classes : *non-overlap between training and test classes*
  - Low-data problem : *very few labeled samples for the test unseen*

Note that, test set also contain $k$-shot labelled data subset (support set).
The performance is measure on the only unlabelled query set.

- Problem Define

< Episodic training and testing >



$C_{train} \cap C_{test} = \emptyset$

$C_{train} = \{cat, goat, boar, dog, shark, lion\}$          $C_{test} = \{dug, dolphin, chicken\}$
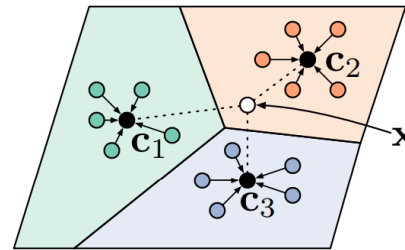
Source:

# Existing Works

- Metric-based approaches

  - Prototypical Network, Matching Network, Relation Network, …

  - Learn a good feature space where categories can distinguish with each other based on a distance metric, and perform distance-based prediction (nearest neighbor classifier).

$$p_\phi(y = k \mid \mathbf{x}) = \frac{\exp(-d(f_\phi(\mathbf{x}), \mathbf{c}_k))}{\sum_{k'} \exp(-d(f_\phi(\mathbf{x}), \mathbf{c}_{k'}))}$$

$$\mathbf{c}_k = \frac{1}{|S_k|} \sum_{(\mathbf{x}_i, y_i) \in S_k} f_\phi(\mathbf{x}_i)$$

Non-parameteric classifier

- Optimization-based approaches

  - MAML, Reptile, …

  - Learn a good initilaization so that the leaner could rapidly adapt to novel tasks within a few optimization steps.

# Existing Works

- Metric-based approaches

  - Prototypical Network, Matching Network, Relation Network, …

  - Learn a good feature space where categories can distinguish with each other based on a distance metric, and perform distance-based prediction (nearest neighbor classifier).

- Optimization-based approaches

  - MAML, Reptile, …

  - Learn a good initilaization so that the leaner could rapidly adapt to novel tasks within a few optimization steps.

**Algorithm 1** Model-Agnostic Meta-Learning

**Require:** $p(\mathcal{T})$: distribution over tasks
**Require:** $\alpha, \beta$: step size hyperparameters
1: randomly initialize $\theta$
2: **while** not done **do**
3:     Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
4:     **for all** $\mathcal{T}_i$ **do**
5:         Evaluate $\nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$ with respect to $K$ examples
6:         Compute adapted parameters with gradient descent: $\theta_i' = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$
7:     **end for**
8:     Update $\theta \leftarrow \theta - \beta \nabla_\theta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta_i'})$
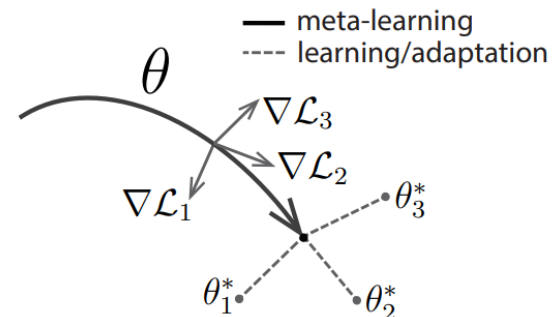9: **end while**

*Figure 1.* Diagram of our model-agnostic meta-learning algorithm (MAML), which optimizes for a representation $\theta$ that can quickly adapt to new tasks.

Source:

- Single image-level annotations can not depict the complex real-world scenes. They **only describe a small subset of an image's content**.

  - Neural Network representations lose any information that is not necessary for performing the training task.
    - ✓ However, this information may be necessary for transfer to new tasks or domain.
    - ✓ **Supervision collapse**

  - This might be problematic,
    - ✓ especially when the training and test time classes are differ.
    - ✓ When a test image containing multiple objects

Source:

- Supervision collapse



(a) Input training image  (b) Class Activation Map of existing method  (c) Input test image  (d) Class Activation Map of existing method  (e) Class Activation Map of our CAN

- The extracted feature may attend to the objects from **seen** classes (person, chair) which have large number of labeled samples in the training set, while ignore the target object from **unseen** class (curtain).

Source:   Cross Attention Network for Few-shot Classification, NeurIPS 2019

Alleviate supervision collapse for few-shot generalization

- Supervision collapse
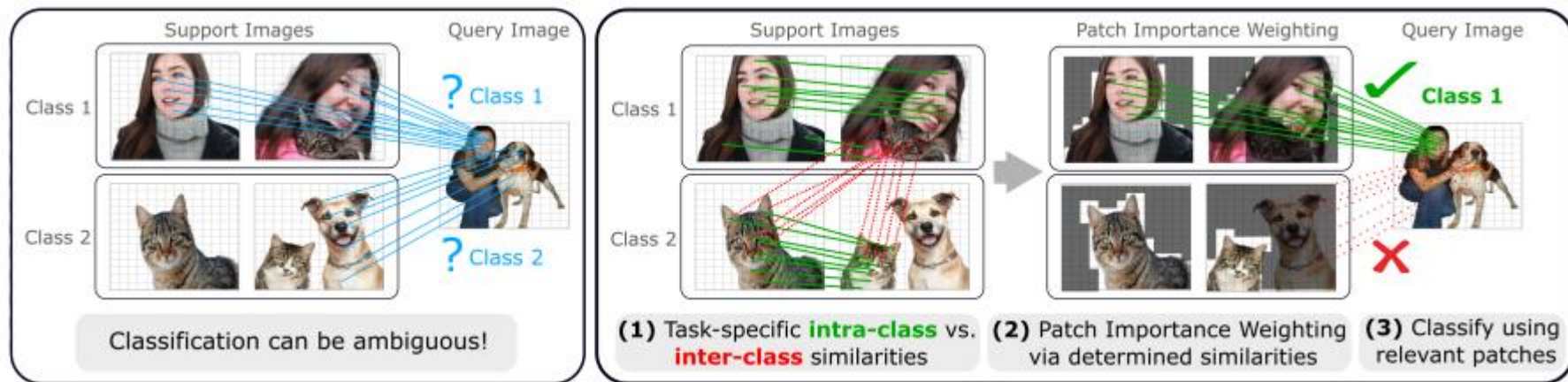


- The errors within that wrong class often have widely different appearance.

- Author's interpretation: the network picks up on image patterns during training that allow images of each class to be tightly grouped in the feature space, minimizing other ways that the image might be similar to other classes in preparation for a confident classification.

Source:     CrossTransformers: spatially-aware few-shot transfer, NeurIPS 2020

- Supervision collapse



- Labels assigned to real-world images with multiple entities only correctly describe a subset of the depicted content, leading to ambiguous classification results.

# Method

- Summary of method

  - Leverage self-supervised pretraining (MIM) for FSL

  - Patch-level similarity (local correspondence) based classification (with ViT)

  - Token importance re-weighting for better classification
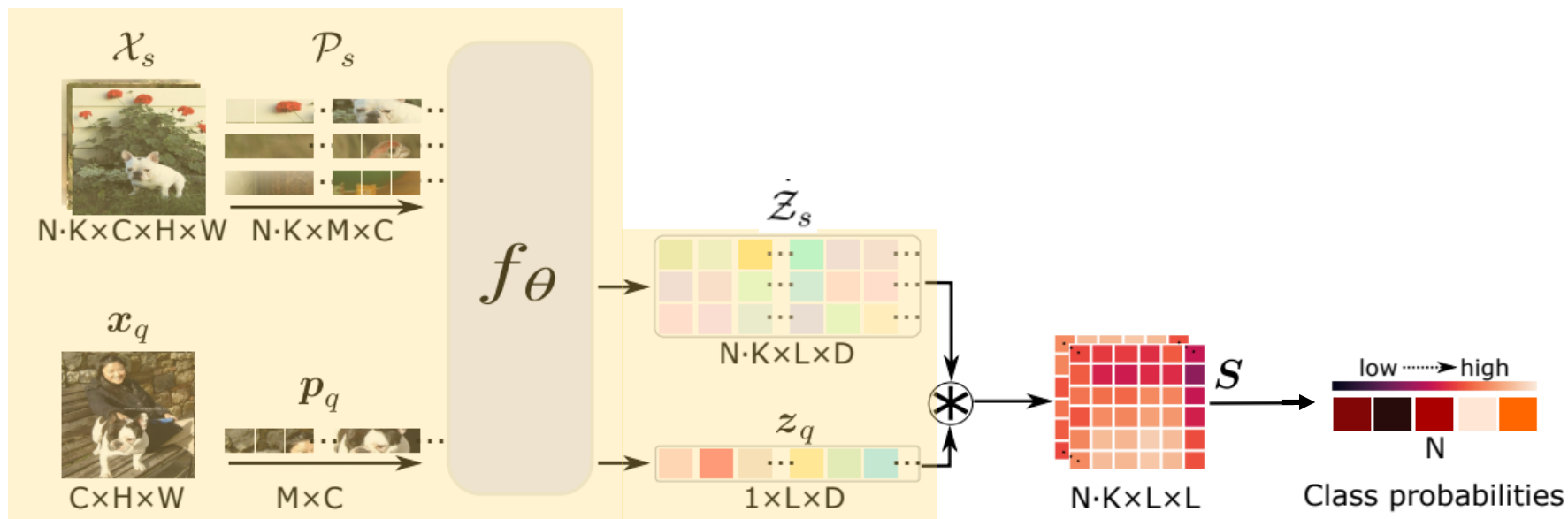
Source:

# Method

- Problem formulation

  - N-way  K-shot  classification with episodic training and testing  $C_{train} \cap C_{test} = \emptyset$

  - An episode (train) is composed of a

    - ✓ Support set $\mathcal{X}_s = \{(x_s^{nk}, y_s^{nk}) | n = 1, \ldots, N; k = 1, \ldots, K; y_s^{nk} \in C_{train}\}$ ,
      where $x_s^{nk}$ denotes the $k$-th sample of class $n$ with label $y_s^{nk}$.

    - ✓ Query set $\mathcal{X}_q = \{(x_q^n, y_q^n) | n = 1, \ldots, N\}$ ,
      where $x_q^n$ denotes a query sample of class $n$ with label $y_q^n$.
      (assume one query sample per class for easy understanding)

# Method

- Overview


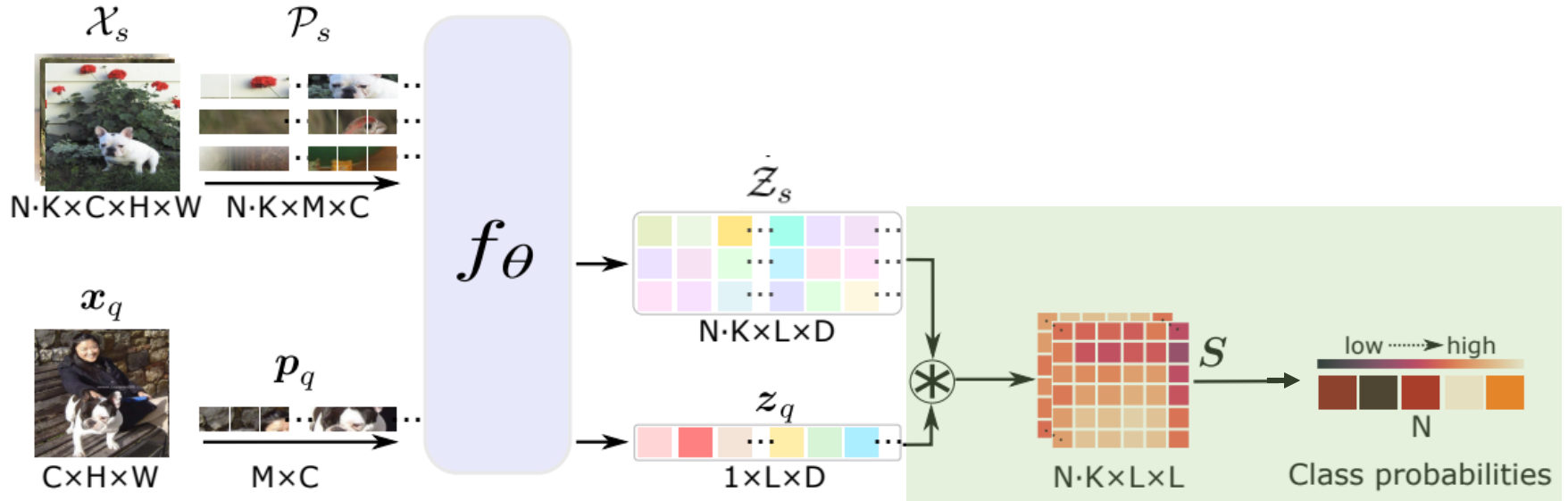
Extract support set and query embeddings
from MIM pre-trained ViT backbone $f_\theta(\cdot)$

Few-shot classification with Transformers Using Reweighted Embedding similarity
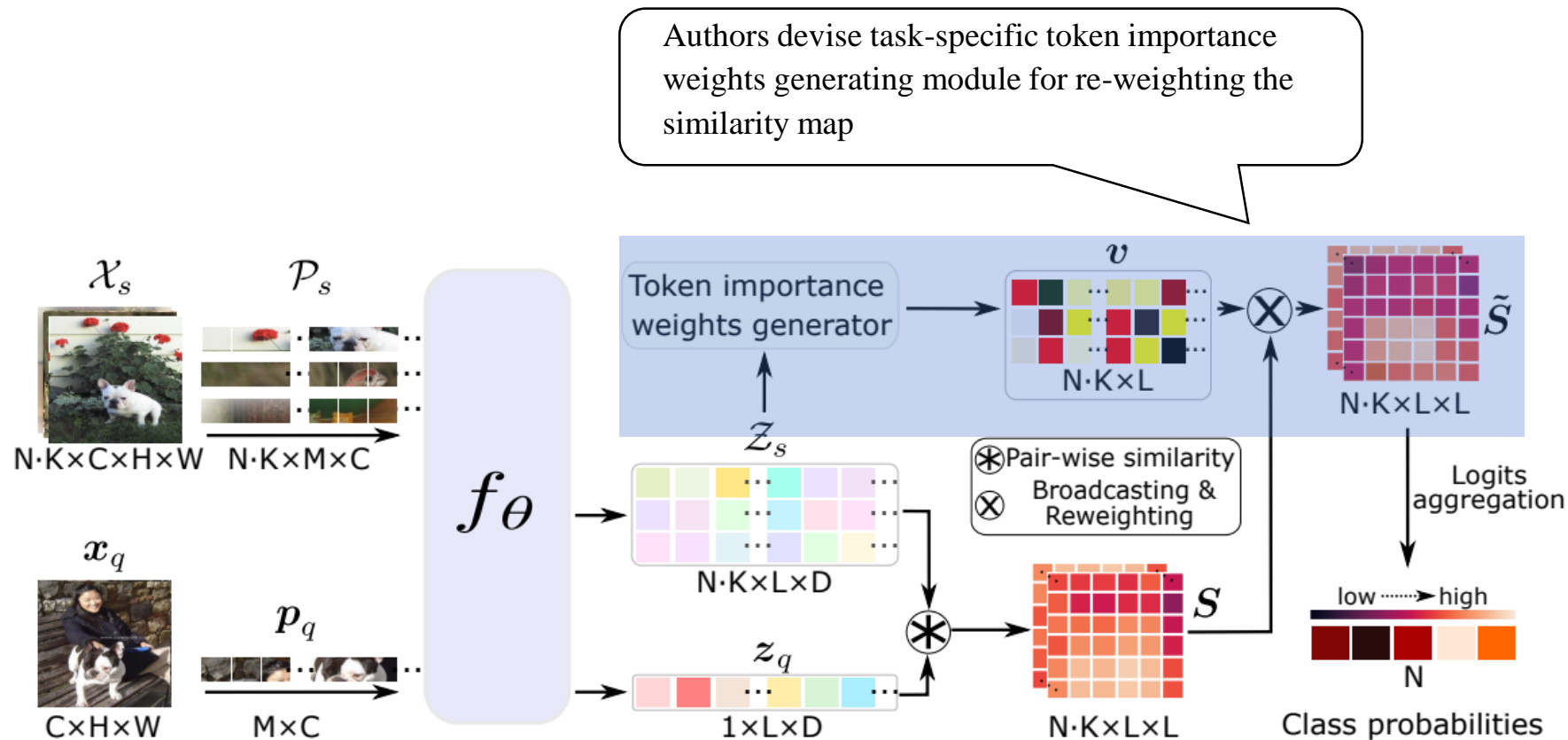
- Overview



Actually, we can directly derive the class probabilities from patch-wise similarity map

- Overview

Authors devise task-specific token importance weights generating module for re-weighting the similarity map

# Method

- In-detail



In ViT, L=M,
In Swim , L<M

1. Split each input image $x \in \mathbb{R}^{H \times W \times C}$ into
   a sequence of $M = (H \cdot W)/P^2$ patches $\boldsymbol{p} = \{p^i\}_{i=1}^M$, with each patch $p^i \in \mathbb{R}^{P^2 \times C}$

2. Fed those tokens into ViT encoder and obtain:
   $Z_s = f_\theta(P_s)$ with $Z_s = \{\boldsymbol{z}_s^{nk} | n = 1, ..., N, k = 1, ..., K\}$, $\boldsymbol{z}_s^{nk} = \{z_s^{nkl} | l = 1, ..., L; z_s^{nkl} \in \mathbb{R}^D\}$
   $z_q = f_\theta(\boldsymbol{p}_q)$ with $z_q = \{z_q^l | l = 1, ..., L; z_q^l \in \mathbb{R}^D\}$

# Method

- In-detail



Local image regions representing similar entities will exhibit higher scores.

3. Based on patch embedding, pair-wise patch simiarlity matrix $S$ is obtained by

$$s_{nk}^{l_s,l_q} = sim(z_s^{nkl_s}, z_q^{l_q}), \text{ where } l_s = 1, ..., L \text{ and } l_q = 1, ..., L$$

- In-detail



4. Task-specific token importance weights $\boldsymbol{v} \in \mathbb{R}^{N \cdot K \cdot L \times 1}$ is inferred via online optimization (inner loop) => will be described in next slide.

5. Reweighting similarity map as $\tilde{S} = S + [v \cdot 1^{1 \times L}]$ and aggregate the information as below

$$\hat{\boldsymbol{y}}_q = \text{softmax}\left(\left\{\hat{y}_q^n\right\}_{n=1}^N\right) = \text{softmax}\left(\left\{\log \sum_{k=1}^K \sum_{l_q=1}^L \sum_{l_s=1}^L \exp\left(\tilde{s}_{nk}^{l_s,l_q}/\tau_S\right)\right\}_{n=1}^N\right)$$

# Method

- In-detail: inner loop token importance weight generator.



- Support set is coppied as $Z_s$ (original) and $Z_{sq}$ (pseudo-query)

- Obtrain the similarity matrix $S_s \in \mathbb{R}^{N \cdot K \cdot L \times N \cdot K \cdot L}$

- Reweighted similarity matrix is computed as $\widetilde{S_s} = S_s + [v^0 + 1^{1 \times N \cdot K \cdot L}]$

- Based on reweighted similarity matrix, perform classification using support set labels

$$\arg\min_{\boldsymbol{v}} \sum_{n=1}^{N} \sum_{k=1}^{K} \mathcal{L}_{\text{CE}} \left( \boldsymbol{y}_s^{nk}, \ \hat{\boldsymbol{y}}_s^{nk}(\boldsymbol{v}) \right)$$

Same with previous equation, but N*K times rep.

# Method

- Train & Inferemce procedure



1. Self-supervised pretraining for $f_\theta$

   ✓ With Masked Image Modeling (ibot) objective

2. Meta fine-tuning for $f_\theta$ and $v$

   ✓ With patch-wise similarity based classification objective

- Note, $v$ is additionally updated in inference time with labelled support set.

# Results

- Dataset
  - miniImageNet
  - tieredImageNet
  - CIFAR-FS
  - FC-100

- Architecture
  - ViT-small
  - Swin-tiny

**<Training strategy comparison>**



Messeage 1. supervised PT < SSL PT

Messeage 2. meta fine-tuning helps FSL largely. Especially for SSL PT backbone.

Source:

- Learned embedding analysis



**Patch embeddings of a support class**

trained w/o $v$      trained w/ $v$

**<same class different instances>**

Messeage : the embeddings retain the instance information and separation is improved when using token importance reweighting.

I think their argument is inconsistent…



**Patch embeddings of the entire support set**

w/ $v$, step 0      w/ $v$, step 15

**<different classes different instances>**

Messeage : token reweighting make the embedding space more disentangled between intra-class samples as well as inter-class samples.

lack of explanation for visualization technique

Source:

- Analysis on token importance re-weighting mechanism



Messeage : FewTRUE select characteristic regions of the depicted objects, and to exclude unimportant or out-of-task information.



(a)

(b)

Messeage : increasing inner loop steps up to 20 aligns with increased performance, both during validation and testing.

Source:

**Few-shot classification with standard benchmarks**

- SOTA on benchmarks

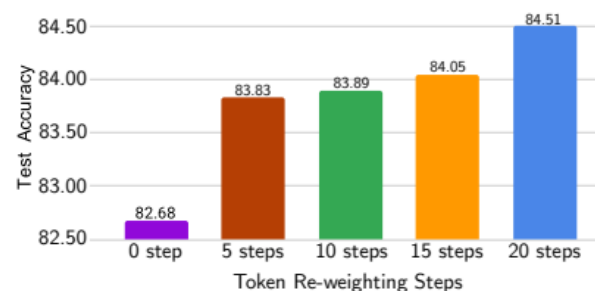| Model | Backbone | ≈ # Params | miniImageNet 1-shot | 5-shot | tieredImageNet 1-shot | 5-shot |
|---|---|---|---|---|---|---|
| ProtoNet [41] | ResNet-12 | 12.4 M | $62.29_{\pm0.33}$ | $79.46_{\pm0.48}$ | $68.25_{\pm0.23}$ | $84.01_{\pm0.56}$ |
| FEAT [53] | ResNet-12 | 12.4 M | $66.78_{\pm0.20}$ | $82.05_{\pm0.14}$ | $70.80_{\pm0.23}$ | $84.79_{\pm0.16}$ |
| DeepEMD [54] | ResNet-12 | 12.4 M | $65.91_{\pm0.82}$ | $82.41_{\pm0.56}$ | $71.16_{\pm0.87}$ | $86.03_{\pm0.58}$ |
| IEPT [56] | ResNet-12 | 12.4 M | $67.05_{\pm0.44}$ | $82.90_{\pm0.30}$ | $72.24_{\pm0.50}$ | $86.73_{\pm0.34}$ |
| MELR [12] | ResNet-12 | 12.4 M | $67.40_{\pm0.43}$ | $83.40_{\pm0.28}$ | $72.14_{\pm0.51}$ | $87.01_{\pm0.35}$ |
| FRN [49] | ResNet-12 | 12.4 M | $66.45_{\pm0.19}$ | $82.83_{\pm0.13}$ | $72.06_{\pm0.22}$ | $86.89_{\pm0.14}$ |
| CG [58] | ResNet-12 | 12.4 M | $67.02_{\pm0.20}$ | $82.32_{\pm0.14}$ | $71.66_{\pm0.23}$ | $85.50_{\pm0.15}$ |
| DMF [52] | ResNet-12 | 12.4 M | $67.76_{\pm0.46}$ | $82.71_{\pm0.31}$ | $71.89_{\pm0.52}$ | $85.96_{\pm0.35}$ |
| InfoPatch [25] | ResNet-12 | 12.4 M | $67.67_{\pm0.45}$ | $82.44_{\pm0.31}$ | - | - |
| BML [60] | ResNet-12 | 12.4 M | $67.04_{\pm0.63}$ | $83.63_{\pm0.29}$ | $68.99_{\pm0.50}$ | $85.49_{\pm0.34}$ |
| CNL [58] | ResNet-12 | 12.4 M | $67.96_{\pm0.98}$ | $83.36_{\pm0.51}$ | $73.42_{\pm0.95}$ | $87.72_{\pm0.75}$ |
| Meta-NVG [55] | ResNet-12 | 12.4 M | $67.14_{\pm0.80}$ | $83.82_{\pm0.51}$ | $74.58_{\pm0.88}$ | $86.73_{\pm0.61}$ |
| PAL [29] | ResNet-12 | 12.4 M | $69.37_{\pm0.64}$ | $84.40_{\pm0.44}$ | $72.25_{\pm0.72}$ | $86.95_{\pm0.47}$ |
| COSOC [28] | ResNet-12 | 12.4 M | $69.28_{\pm0.49}$ | $85.16_{\pm0.42}$ | $73.57_{\pm0.43}$ | $87.57_{\pm0.10}$ |
| Meta DeepBDC [51] | ResNet-12 | 12.4 M | $67.34_{\pm0.43}$ | $84.46_{\pm0.28}$ | $72.34_{\pm0.49}$ | $87.31_{\pm0.32}$ |
| LEO [39] | WRN-28-10 | 36.5 M | $61.76_{\pm0.08}$ | $77.59_{\pm0.12}$ | $66.33_{\pm0.05}$ | $81.44_{\pm0.09}$ |
| CC+rot [15] | WRN-28-10 | 36.5 M | $62.93_{\pm0.45}$ | $79.87_{\pm0.33}$ | $70.53_{\pm0.51}$ | $84.98_{\pm0.36}$ |
| FEAT [53] | WRN-28-10 | 36.5 M | $65.10_{\pm0.20}$ | $81.11_{\pm0.14}$ | $70.41_{\pm0.23}$ | $84.38_{\pm0.16}$ |
| PSST [8] | WRN-28-10 | 36.5 M | $64.16_{\pm0.44}$ | $80.64_{\pm0.32}$ | - | - |
| MetaQDA [57] | WRN-28-10 | 36.5 M | $67.83_{\pm0.64}$ | $84.28_{\pm0.69}$ | $74.33_{\pm0.65}$ | $89.56_{\pm0.79}$ |
| OM [36] | WRN-28-10 | 36.5 M | $66.78_{\pm0.30}$ | $85.29_{\pm0.41}$ | $71.54_{\pm0.29}$ | $87.79_{\pm0.46}$ |
| FewTURE (ours) | ViT-Small | 22 M | $68.02_{\pm0.88}$ | $84.51_{\pm0.53}$ | $72.96_{\pm0.92}$ | $86.43_{\pm0.67}$ |
| FewTURE (ours) | Swin-Tiny | 29 M | $\mathbf{72.40_{\pm0.78}}$ | $\mathbf{86.38_{\pm0.49}}$ | $\mathbf{76.32_{\pm0.87}}$ | $\mathbf{89.96_{\pm0.55}}$ |

| Model | Backbone | ≈ # Params | CIFAR-FS 1-shot | 5-shot | FC100 1-shot | 5-shot |
|---|---|---|---|---|---|---|
| ProtoNet [41] | ResNet-12 | 12.4 M | - | - | $41.54_{\pm0.76}$ | $57.08_{\pm0.76}$ |
| MetaOpt [21] | ResNet-12 | 12.4 M | $72.00_{\pm0.70}$ | $84.20_{\pm0.50}$ | $41.10_{\pm0.60}$ | $55.50_{\pm0.60}$ |
| MABAS [20] | ResNet-12 | 12.4 M | $73.51_{\pm0.92}$ | $85.65_{\pm0.65}$ | $42.31_{\pm0.75}$ | $58.16_{\pm0.78}$ |
| RFS [45] | ResNet-12 | 12.4 M | $73.90_{\pm0.80}$ | $86.90_{\pm0.50}$ | $44.60_{\pm0.70}$ | $60.90_{\pm0.60}$ |
| BML [60] | ResNet-12 | 12.4 M | $73.45_{\pm0.47}$ | $88.04_{\pm0.33}$ | - | - |
| CG [14] | ResNet-12 | 12.4 M | $73.00_{\pm0.70}$ | $85.80_{\pm0.50}$ | - | - |
| Meta-NVG [55] | ResNet-12 | 12.4 M | $74.63_{\pm0.91}$ | $86.45_{\pm0.59}$ | $46.40_{\pm0.81}$ | $61.33_{\pm0.71}$ |
| RENet [19] | ResNet-12 | 12.4 M | $74.51_{\pm0.46}$ | $86.60_{\pm0.32}$ | - | - |
| TPMN [50] | ResNet-12 | 12.4 M | $75.50_{\pm0.90}$ | $87.20_{\pm0.60}$ | $46.93_{\pm0.71}$ | $63.26_{\pm0.74}$ |
| MixFSL [1] | ResNet-12 | 12.4 M | - | - | $44.89_{\pm0.63}$ | $60.70_{\pm0.60}$ |
| CC+rot [15] | WRN-28-10 | 36.5 M | $73.62_{\pm0.31}$ | $86.05_{\pm0.22}$ | - | - |
| PSST [8] | WRN-28-10 | 36.5 M | $77.02_{\pm0.38}$ | $88.45_{\pm0.35}$ | - | - |
| Meta-QDA [57] | WRN-28-10 | 36.5 M | $75.83_{\pm0.88}$ | $88.79_{\pm0.75}$ | - | - |
| FewTURE (ours) | ViT-Small | 22 M | $76.10_{\pm0.88}$ | $86.14_{\pm0.64}$ | $46.20_{\pm0.79}$ | $63.14_{\pm0.73}$ |
| FewTURE (ours) | Swin-Tiny | 29 M | $\mathbf{77.76_{\pm0.81}}$ | $\mathbf{88.90_{\pm0.59}}$ | $\mathbf{47.68_{\pm0.78}}$ | $\mathbf{63.81_{\pm0.75}}$ |

Powerful, but, different backbone architectures can be justified?

Source:

- Additional experiments

  - Token pruning

    Table 3: Pruning the number of tokens. Test accuracy for 5-way 5-shot on *mini*ImageNet [48].

    | # tokens | Test Acc. |
    |---|---|
    | 100% | $84.05 \pm 0.53$ |
    | 75% | $83.15 \pm 0.57$ |
    | 50% | $83.81 \pm 0.59$ |
    | 25% | $81.79 \pm 0.57$ |
    | 10% | $81.05 \pm 0.62$ |

    > Pruning via attention map value thresholding

  - Ablation on classifier

    Table 4: Changing the classifier. Test accuracy for 5-way 5-shot on *mini*ImageNet [48].

    | Classifier | Test Acc. |
    |---|---|
    | Prototyp. w/ Euclid. Dist. | $82.80 \pm 0.59$ |
    | Prototyp. w/ Cosine. Dist. | $79.90 \pm 0.65$ |
    | Linear (optimized online) | $82.37 \pm 0.57$ |
    | FewTURE ( 0 rew. steps) | $82.68 \pm 0.55$ |
    | FewTURE (15 rew. steps) | $\mathbf{84.05 \pm 0.53}$ |

    > Comparison with
    > - Other non-parameter classifier (proto)
    > - learnable parameteric classifier (linear)

Source:

# Conclusion & Discussion

- Summary
  - Using ViT with SSL pre-training (MIM) for FSL for the first-time.
    - ✓ (To mitigate supervision collapse)
  - Propose a new patch-wise simiarlity based non-parametric classifier
  - Propose online token importance re-weighting module

- Discussion
  - The "supervision collapse" and "multiple entities" problem is not convincing…
    - ✓ Authors gave just illustration.
  - The experiments are only performed on standard single-object-centric dataset..
    - ✓ Few-shot variants of ImageNet / CIFAR100.
  - The idea underlying methodology is not completely new.
    - ✓ Consideration on local correspondence is already addressed by CAN (Hou et al. 2019)
    - ✓ SSL objective is also used in previous work CTX (Doersch et al. 2020)

Source: