

# Supervised Masked Knowledge Distillation for Few-Shot Transformers

Authors

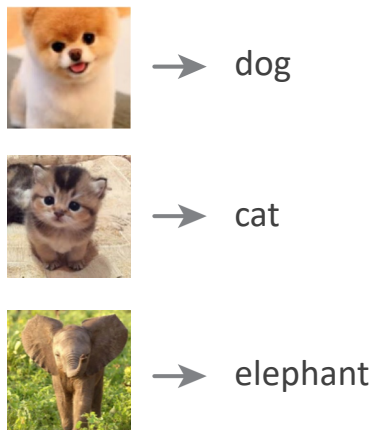
Han Lin\*, Guangxing Han\*, Jiawei Ma, Shiyuan Huang, Xudong Lin, Shih-Fu Chang

## Observations

- Unlike CNN-based models, Transformer architecture lacks **inductive bias**:
  - ✓ captures long-range token dependencies
  - ✗ data-hungry, easy to overfit to pre-trained datasets which are not large enough
- Overfitting to datasets with insufficiently training data may **hurt the generalization ability** of new classes (e.g., few-shot learning)
- Our problem settings:
  - Transformers on small dataset under few-shot learning

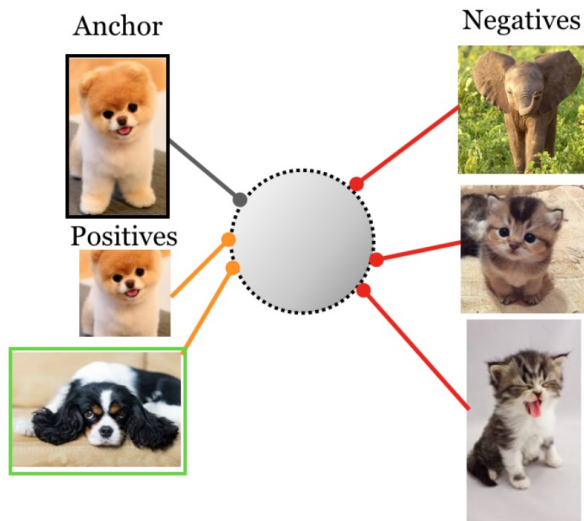
# Related Works

## Supervised Learning



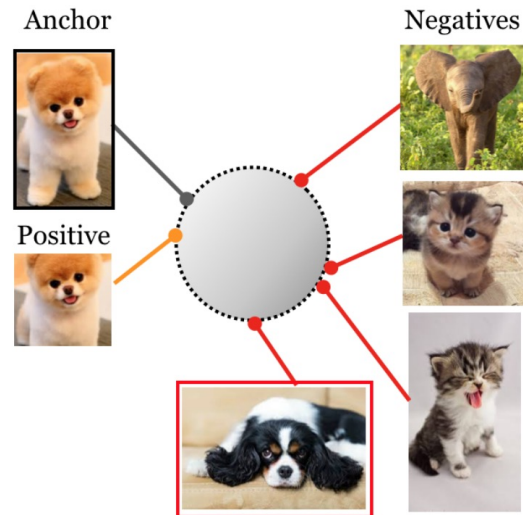
- ✗ Easy to overfit
- ✗ Bad generalization ability

## Supervised-Contrastive Learning



- ✗ Construct negative samples needs large batch size

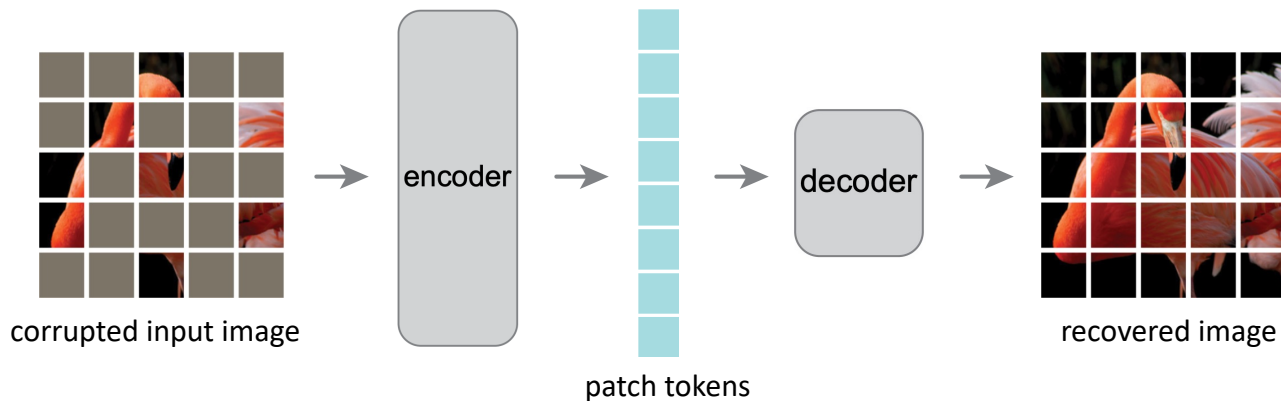
## Self-Supervised Learning



- ✓ Better generalization ability
- ✓ Can avoid negative samples
- ✗ Hard to capture high-level semantic features

## Related Works

- **Masked Image Modeling (MIM)**
  - Recovering masked pixels from a corrupted input image
  - Combined with self-distillation achieves better performance <sup>[3]</sup>



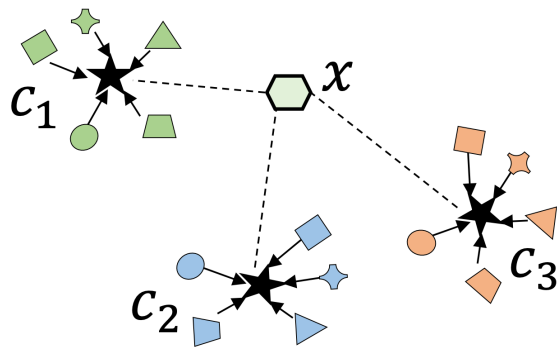
## Our Work

To tackle the problem of Transformers on small dataset:

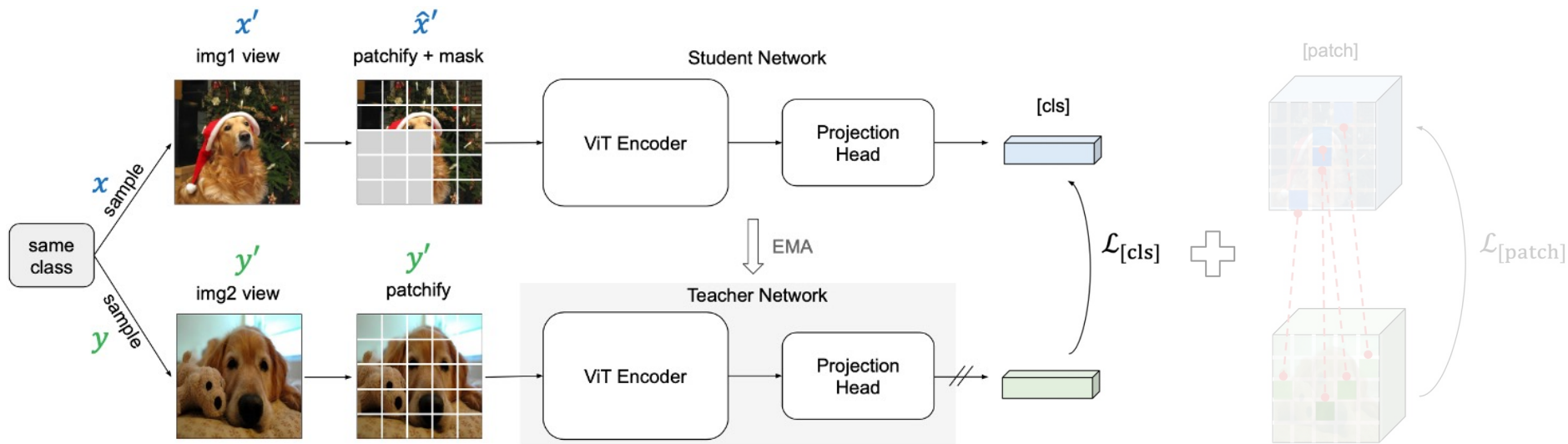


## Our Training & Few-Shot Evaluation Pipelines:

- Train on **base** classes:
  - Our method defined in the next page
- Few-shot evaluation on **novel** classes:
  - **Prototype** evaluation method
  - N-way K-shot (e.g. 3-way 5-shot)

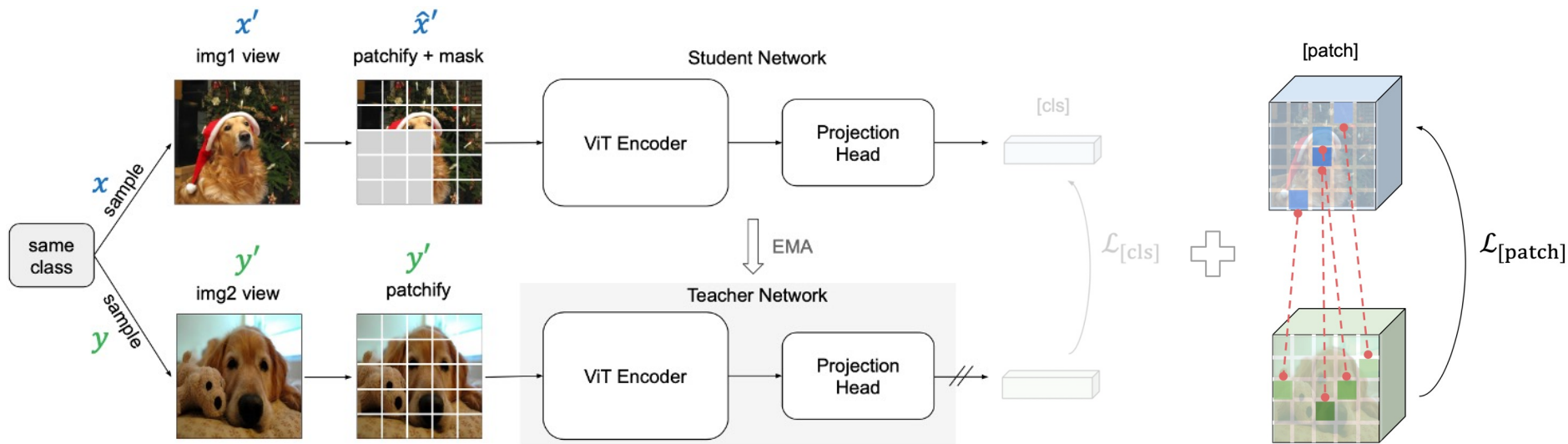


# Our SMKD Framework (Global Knowledge Distillation)



$$\mathcal{L}_{[cls]} = -P_t^{[cls]}(y') \log P_s^{[cls]}(\hat{x}')$$

# Our SMKD Framework (Local Knowledge Distillation)



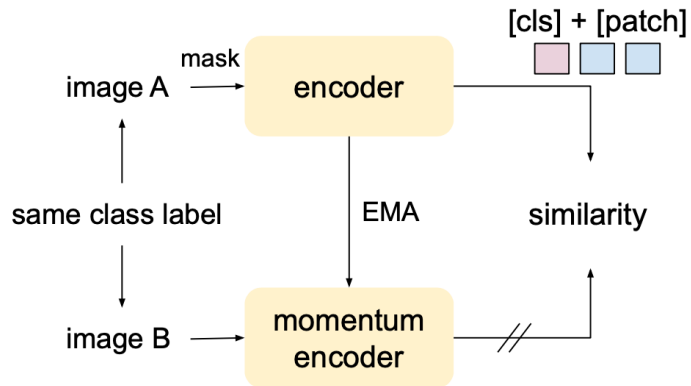
Finding dense correspondence of matched token pairs with highest similarities.

$$\mathcal{L}_{[patch]} = - \sum_{k=1}^N \omega_{k^+} \cdot P_t^{[patch]}(y'_k) \log P_s^{[patch]}(\hat{x}'_{k^+}) \quad k^+ = \arg \max_{l \in [N]} \frac{f_k^t \top f_l^s}{\|f_k^t\| \|f_l^s\|}$$



# Comparison with Self-Supervised / Supervised Contrastive Frameworks

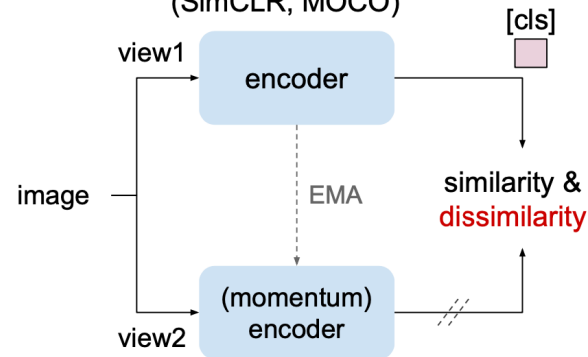
## Our supervised masked knowledge distillation



✓ Avoids the need for negative examples ("dissimilarity")

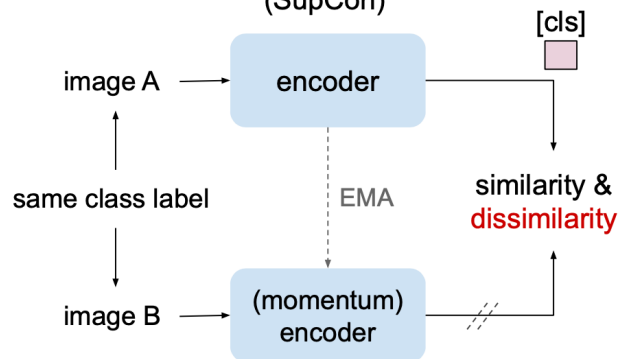
## self-supervised contrastive

(SimCLR, MOCO)



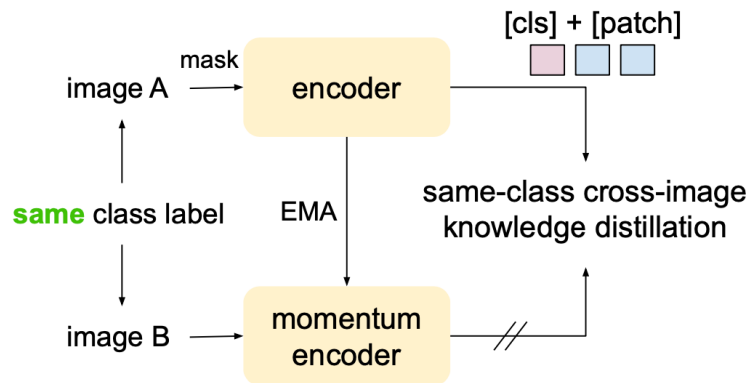
## supervised contrastive

(SupCon)

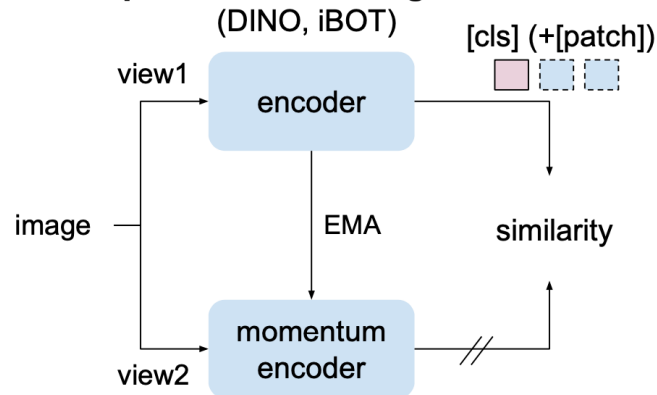


# Comparison with Self-Supervised Knowledge Distillation Framework

## Our supervised masked knowledge distillation



## self-supervised knowledge distillation



- ✓ Introduces **intra-class** knowledge distillation to self-supervised knowledge distillation framework

## Dataset Description

- We test on four widely-used few-shot classification benchmark datasets:

	Resolution	#Images	#Classes	#Images per class	(train, val, test) split
CIFAR-FS	$32 \times 32$	60000	100	600	(64, 16, 20)
FC100	$32 \times 32$	60000	100	600	(60, 20, 20)
<i>mini</i> -ImageNet	$224 \times 224$	60000	100	600	(64, 16, 20)
<i>tiered</i> -ImageNet	$224 \times 224$	779165	608	$\approx 1282$	(351, 97, 160)

- Comparison of dataset size:

CIFAR-FS  $\approx$  FC100  $\leq$  *mini*-ImageNet  $\ll$  *tiered*-ImageNet

# Visualizations

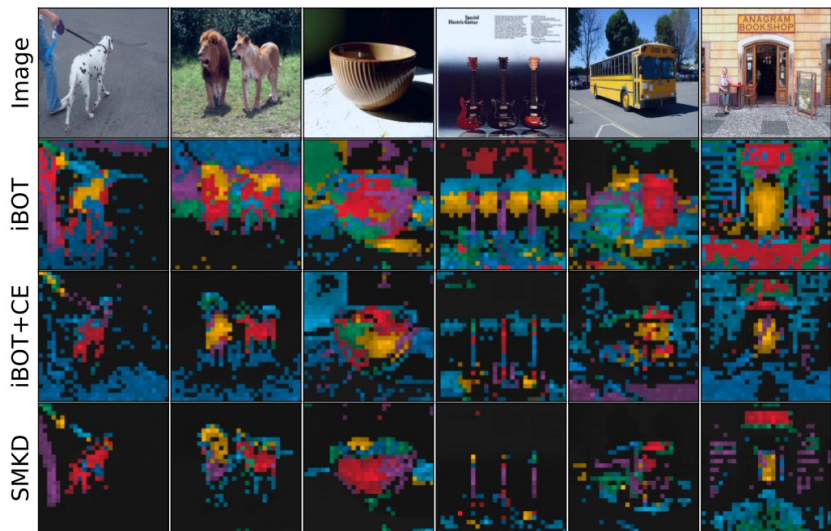


Figure 4. **Visualization of multi-head self-attention maps.** The self-attention of the  $[cls]$  tokens with different heads in the last attention layer of ViT are visualized in different colors. iBOT+CE represents the model first pre-trained with iBOT, then trained with CE loss. Our SMKD pays more attention to the foreground objects, especially the most discriminative parts.

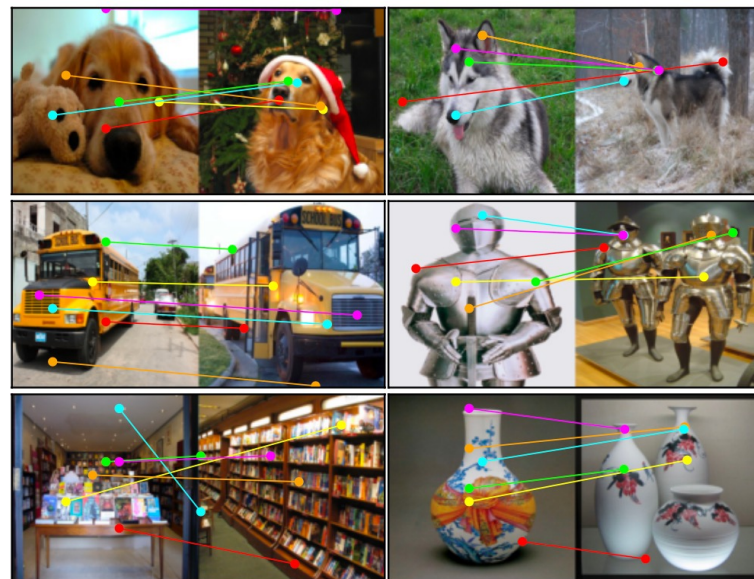


Figure 5. **Visualization of dense correspondence.** We use the patches with the highest self-attention of the  $[cls]$  token on each attention head (6 in total) of the last layer of ViT-S as queries. Best-matched patches with the highest similarities are connected with lines.

## Comparison with SOTAs

- Our method with simple Prototype and Linear Classifier evaluation methods could beat all models with CNN backbone.
- Our method grows more effective on datasets with **smaller resolutions** and **fewer training images**.
- Our method, combined with tricks from HCT (spectral tokens pooling & small patch size), achieves a new SOTA on **mini-ImageNet, CIFAR-FS, and FC100**, and is comparable with current SOTA (HCTransformers) on **tiered-ImageNet**.

Table 1. Results on mini-ImageNet and tiered-ImageNet. Top three methods are colored in red, blue and green respectively.

Method	Backbone	#Params	miniImageNet,5-way		tieredImageNet,5-way	
			1-shot	5-shot	1-shot	5-shot
DeepEMD [73]	<i>ResNet-12</i>	12.4M	65.91 ± 0.82	82.41 ± 0.56	71.16 ± 0.87	86.03 ± 0.58
IE [50]	<i>ResNet-12</i>	12.4M	67.28 ± 0.80	84.78 ± 0.52	72.21 ± 0.90	87.08 ± 0.58
PAL [41]	<i>ResNet-12</i>	12.4M	69.37 ± 0.64	84.40 ± 0.44	72.25 ± 0.72	86.95 ± 0.47
DC [71]	<i>ResNet-12</i>	12.4M	68.57 ± 0.55	82.88 ± 0.42	78.19 ± 0.25	89.90 ± 0.41
COSOC [40]	<i>ResNet-12</i>	12.4M	69.28 ± 0.49	85.16 ± 0.42	73.57 ± 0.43	87.57 ± 0.10
FEAT [72]	<i>WRN-28-10</i>	36.5M	65.10 ± 0.20	81.11 ± 0.14	70.41 ± 0.23	84.38 ± 0.16
Meta-QDA [75]	<i>WRN-28-10</i>	36.5M	67.38 ± 0.55	84.27 ± 0.75	74.29 ± 0.66	89.41 ± 0.77
OM [46]	<i>WRN-28-10</i>	36.5M	66.78 ± 0.30	85.29 ± 0.41	71.54 ± 0.29	87.79 ± 0.46
SUN [16]	<i>ViT</i>	12.5M	67.80 ± 0.45	83.25 ± 0.30	72.99 ± 0.50	86.74 ± 0.33
FewTURE [29]	<i>Swin-Tiny</i>	29.0M	72.40 ± 0.78	86.38 ± 0.49	76.32 ± 0.87	89.96 ± 0.55
HCTransformers [28]	$3 \times \text{ViT-S}$	63.0M	<b>74.74 ± 0.17</b>	<b>89.19 ± 0.13</b>	<b>79.67 ± 0.20</b>	<b>91.72 ± 0.11</b>
Ours (Prototype)	<i>ViT-S</i>	21M	<b>74.28 ± 0.18</b>	88.82 ± 0.09	<b>78.83 ± 0.20</b>	91.02 ± 0.12
Ours (Classifier)	<i>ViT-S</i>	21M	74.10 ± 0.17	<b>88.89 ± 0.09</b>	78.81 ± 0.21	<b>91.21 ± 0.11</b>
Ours + HCT [28]	$3 \times \text{ViT-S}$	63M	<b>75.32 ± 0.18</b>	<b>89.57 ± 0.09</b>	<b>79.74 ± 0.20</b>	<b>91.68 ± 0.11</b>

Table 2. Results on CIFAR-FS and FC100. Top three methods are colored in red, blue and green respectively.

Method	Backbone	#Params	CIFAR-FS,5-way		FC100,5-way	
			1-shot	5-shot	1-shot	5-shot
BML [79]	<i>ResNet-12</i>	12.4M	73.45 ± 0.47	88.04 ± 0.33	45.00 ± 0.41	63.03 ± 0.41
IE [50]	<i>ResNet-12</i>	12.4M	77.87 ± 0.85	89.74 ± 0.57	47.76 ± 0.77	65.30 ± 0.76
PAL [41]	<i>ResNet-12</i>	12.4M	77.10 ± 0.70	88.00 ± 0.50	47.20 ± 0.60	64.00 ± 0.60
TPMN [66]	<i>ResNet-12</i>	12.4M	75.50 ± 0.90	87.20 ± 0.60	46.93 ± 0.71	63.26 ± 0.74
MN+MC [74]	<i>ResNet-12</i>	12.4M	74.63 ± 0.91	86.45 ± 0.59	46.40 ± 0.81	61.33 ± 0.71
ConstellationNet [70]	<i>ResNet-12</i>	12.4M	75.40 ± 0.20	86.80 ± 0.20	43.80 ± 0.20	59.70 ± 0.20
PSST [11]	<i>WRN-28-10</i>	36.5M	77.02 ± 0.38	88.45 ± 0.35	-	-
Meta-QDA [75]	<i>WRN-28-10</i>	36.5M	75.95 ± 0.59	88.72 ± 0.79	-	-
SUN [16]	<i>ViT</i>	12.5M	78.37 ± 0.46	88.84 ± 0.32	-	-
FewTURE [29]	<i>Swin-Tiny</i>	29.0M	77.76 ± 0.81	88.90 ± 0.59	47.68 ± 0.78	63.81 ± 0.75
HCTransformers [28]	$3 \times \text{ViT-S}$	63.0M	<b>78.89 ± 0.18</b>	<b>90.50 ± 0.09</b>	<b>48.27 ± 0.15</b>	<b>66.42 ± 0.16</b>
Ours (Prototype)	<i>ViT-S</i>	21M	<b>80.08 ± 0.18</b>	<b>90.63 ± 0.13</b>	<b>50.38 ± 0.16</b>	<b>68.37 ± 0.16</b>
Ours (Classifier)	<i>ViT-S</i>	21M	<b>79.82 ± 0.18</b>	<b>90.91 ± 0.13</b>	<b>50.28 ± 0.16</b>	<b>68.50 ± 0.16</b>

Thank you!

Any suggestions and comments are welcome!

## Reference

- [1] Khosla, Prannay, et al. "Supervised contrastive learning." Advances in Neural Information Processing Systems 33 (2020): 18661-18673.
- [2] He, Kaiming, et al. "Masked autoencoders are scalable vision learners." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.
- [3] Zhou, Jinghao, et al. "Image BERT Pre-training with Online Tokenizer." International Conference on Learning Representations. 2021.
- [4] Ma, Jiawei, et al. "Partner-assisted learning for few-shot image classification." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2021.
- [5] He, Kaiming, et al. "Masked autoencoders are scalable vision learners." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.
- [6] He, Kaiming, et al. "Momentum contrast for unsupervised visual representation learning." Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2020.
- [7] Chen, Ting, et al. "A simple framework for contrastive learning of visual representations." International conference on machine learning. PMLR, 2020.
- [8] Chen, Xinlei, Saining Xie, and Kaiming He. "An empirical study of training self-supervised vision transformers." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2021.