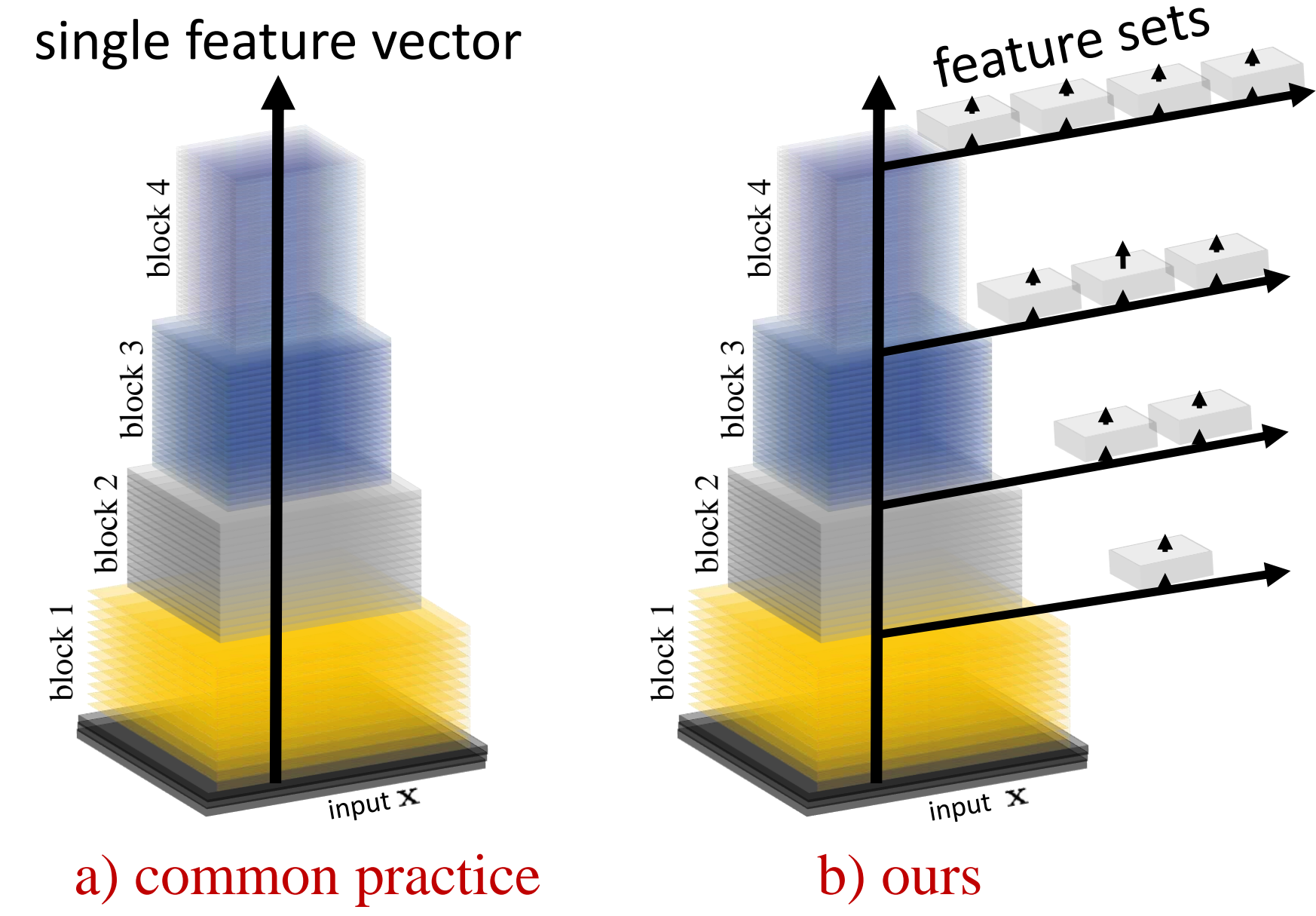
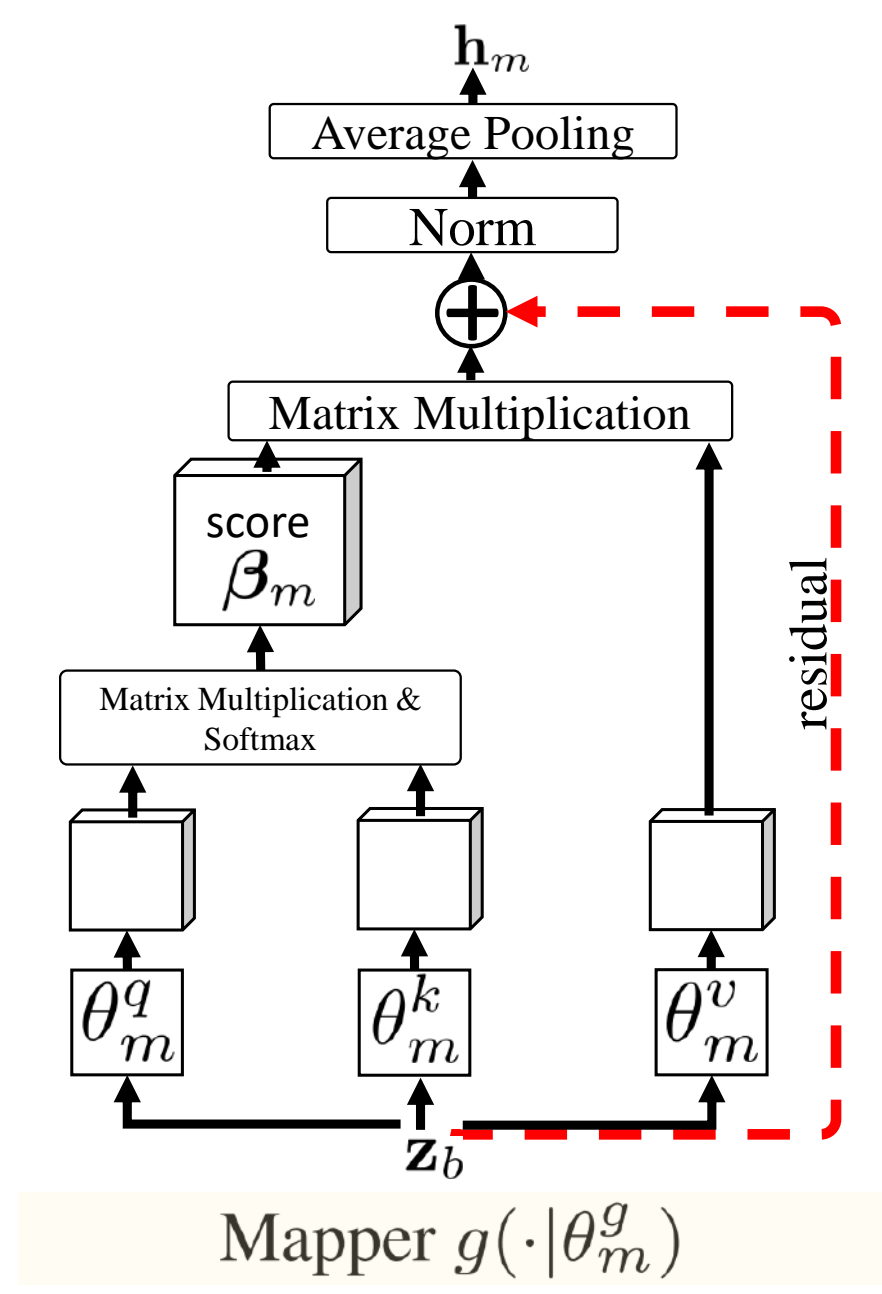


SETFEAT

- it is common practice to extract a single feature vector per input image
- we propose set-feature extractor (SetFeat) to represent images as *sets* of feature vectors



- we take inspiration from Feature Pyramid Networks to learn a richer feature space
- SetFeat embeds shallow self-attention mappers in existing architecture



- for attention map, we first compute

$$\beta_m = \text{Softmax} \left(q(z_{b_m} | \theta_m^q) k(z_{b_m} | \theta_m^k)^\top / \sqrt{d_k} \right),$$

where $\beta_m \in \mathbb{R}^{P \times P}$

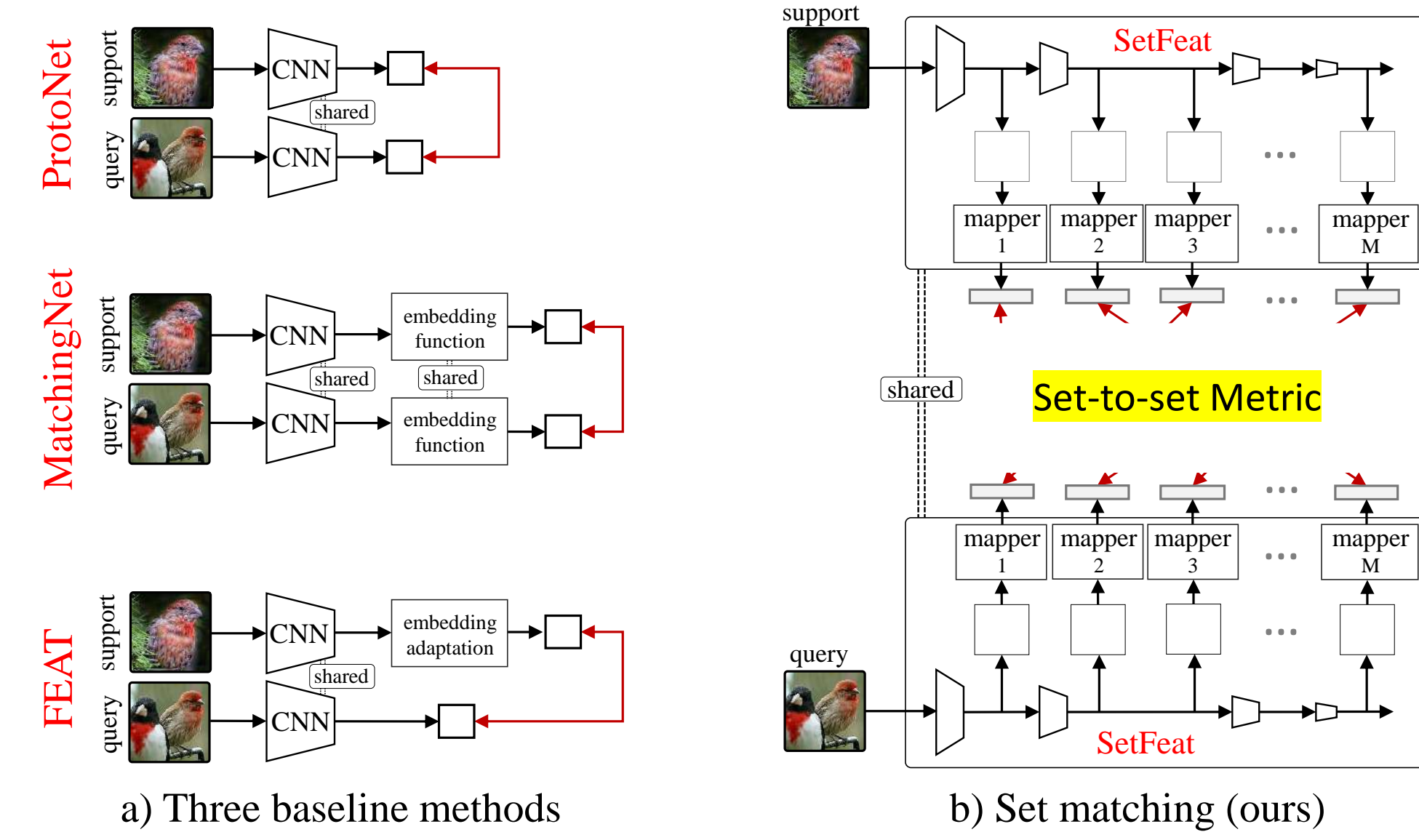
- then, we compute the attention over β_m

$$\mathbf{a}_m = \beta_m v(z_{b_m} | \theta_m^v),$$

where $\mathbf{a}_m \in \mathbb{R}^{P \times D^a}$ consists of P patches

SET-TO-SET METRICS

- SetFeat first extracts sets of features
- then, we need a set-to-set metric to compare the feature set of the query with the feature sets corresponding to each instance of the support set of each class



- Match-sum** aggregates the distance between matching mappers

$$d_{ms}(\mathbf{x}_q, \mathcal{S}^n) = \sum_{i=1}^M d(\mathbf{h}_i(\mathbf{x}_q), \bar{\mathbf{h}}_i(\mathcal{S}^n)).$$

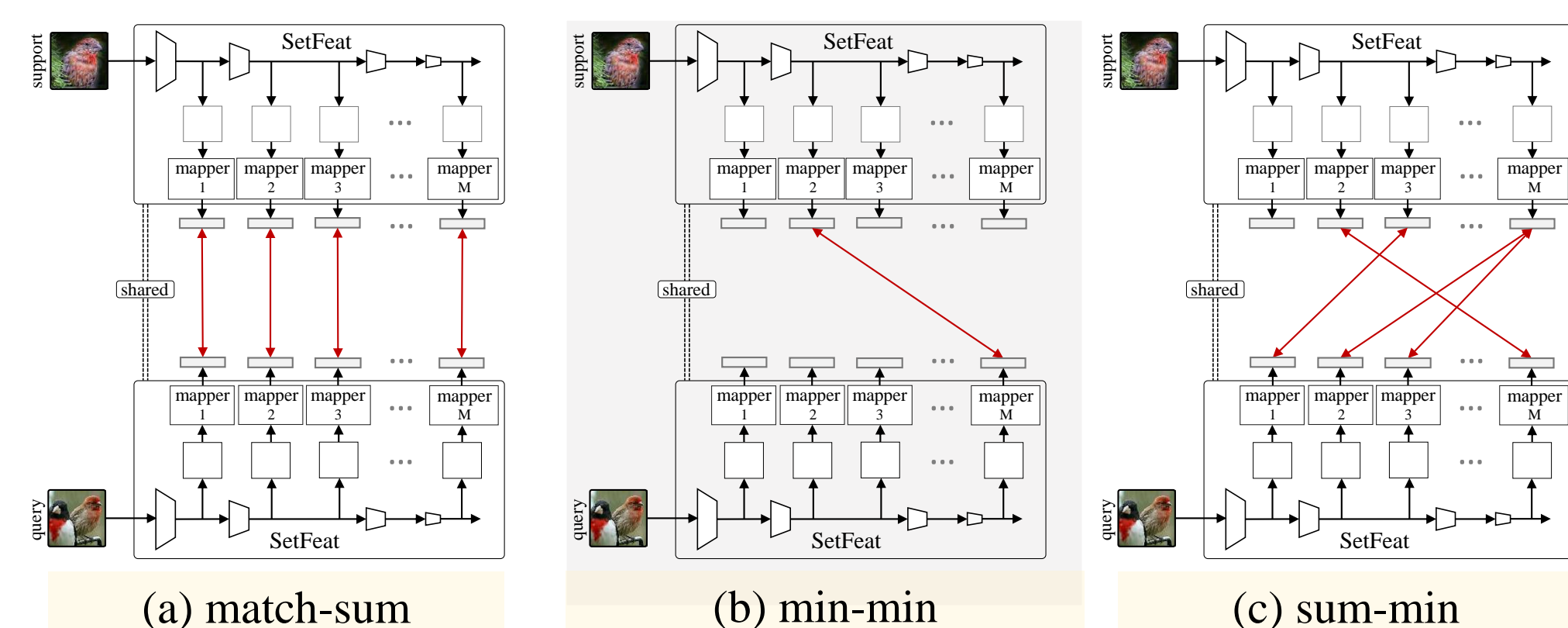
- Min-min** uses the minimum distance across all possible pairs of elements

$$d_{mm}(\mathbf{x}_q, \mathcal{S}^n) = \min_{i=1}^M \min_{j=1}^M d(\mathbf{h}_i(\mathbf{x}_q), \bar{\mathbf{h}}_j(\mathcal{S}^n)).$$

- Sum-min** aggregates with a sum the minimum distances between the mappers

$$d_{sm}(\mathbf{x}_q, \mathcal{S}^n) = \sum_{i=1}^M \min_{j=1}^M d(\mathbf{h}_i(\mathbf{x}_q), \bar{\mathbf{h}}_j(\mathcal{S}^n)).$$

- illustration of our set-to-set metric



EVALUATIONS

- MiniImageNet** evaluations of SetFeat4 results in **+2.03%** improvement in 1-shot

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 [50]	Conv4-64	49.42 \pm 0.78	68.20 \pm 0.66	
MAML [18]		48.07 \pm 1.75	63.15 \pm 0.91	
RelationNet [53]		50.44 \pm 0.82	65.32 \pm 0.70	
Baseline++ [8]		48.24 \pm 0.75	66.43 \pm 0.63	
IMP [3]		49.60 \pm 0.80	68.10 \pm 0.80	
MemoryNet [7]		53.37 \pm 0.48	66.97 \pm 0.35	
Neg-Margin [35]		52.84 \pm 0.76	70.41 \pm 0.66	
MixtFSL [2]		52.82 \pm 0.63	70.67 \pm 0.57	
FEAT [68]		55.15 \pm 0.20	71.61 \pm 0.16	
MELR [16]		55.35 \pm 0.43	72.27 \pm 0.35	
BOIL [39]		49.61 \pm 0.16	66.45 \pm 0.37	
Match-sum		SF4-64	55.74 \pm 0.65	72.18 \pm 0.70
Min-min			56.22\pm0.89	72.70\pm0.65
Sum-min			57.18\pm0.89	73.67\pm0.71

- TieredImageNet** evaluation of SetFeat12 results in **+1.42%** improvement in 1-shot

Table 2. TieredImageNet evaluation. Bold/red is best/second best, and \pm indicates the 95% conf. intervals over 600 episodes of 5-way.

Method	Backbone	1-shot	5-shot
OptNet [31]	ResNet12	65.99 \pm 0.72	81.56 \pm 0.53
MTL [52]		65.62 \pm 1.80	80.61 \pm 0.90
DNS [48]		66.22 \pm 0.75	82.79 \pm 0.48
Simple [55]		69.74 \pm 0.72	84.41 \pm 0.55
TapNet [70]		63.08 \pm 0.15	80.26 \pm 0.12
ProtoNet [†] [50]		68.23 \pm 0.23	84.03 \pm 0.16
FEAT [68]		70.80 \pm 0.23	84.79 \pm 0.16
MixtFSL [2]		70.97 \pm 1.03	86.16 \pm 0.67
Distill [55]		71.52 \pm 0.69	86.03 \pm 0.49
DeepEMD [74]		71.16 \pm 0.87	86.03 \pm 0.58
DMF [66]		71.89 \pm 0.52	85.96 \pm 0.35
MELR [16]		72.14 \pm 0.51	87.01 \pm 0.35
Distill [45]		72.21\pm0.90	87.08\pm0.58
Match-sum		SF12	71.22 \pm 0.86
Min-min		71.75 \pm 0.90	86.40 \pm 0.56
Sum-min		73.63\pm0.88	87.59\pm0.57

[†]taken from [31]; Mappers dimension: SF12 $\in \mathbb{R}^{512}$

- CUB** evaluation of SetFeat4 results in **+1.83%** improvement in 1-shot

Table 3. Fine-grained evaluation using CUB in 5-way. \pm is the 95% confidence intervals on 600 episodes ([†]taken from [54]).

Method	Backbone	1-shot	5-shot
MatchingNet [59]	Conv4-64	61.16 \pm 0.89	72.86 \pm 0.70
ProtoNet [50]		64.42 \pm 0.48	81.82 \pm 0.35
MAML [17]		55.92 \pm 0.95	72.09 \pm 0.76
RelationNet [53]		62.45 \pm 0.98	76.11 \pm 0.69
FEAT [68]		68.87 \pm 0.22	82.90 \pm 0.15
MELR [16]		70.26\pm0.50	85.01\pm0.32
Match-sum	SF4-64	67.35 \pm 0.93	83.82 \pm 0.61
Min-min		70.15 \pm 0.93	84.94 \pm 0.64
Sum-min		72.09\pm0.92	87.05\pm0.58

ABLATIONS

- Mapper configurations.** different ways of embedding ten mappers throughout the backbone

Table 4. Ablation of different mapper-level combinations using miniImageNet. The results are validation accuracy with min-sum.

Mappers	SetFeat4-64	SetFeat4-512
	1-shot	5-shot
ProtoNet*	53.51	71.57
0-0-0-1	53.55	71.51
1-2-3-4 (concat)	53.56	71.82
1-1-1-1	51.11	69.41
0-0-0-10	52.90	69.49
2-2-3-3	54.73	71.98
1-2-3-4	54.71	71.35

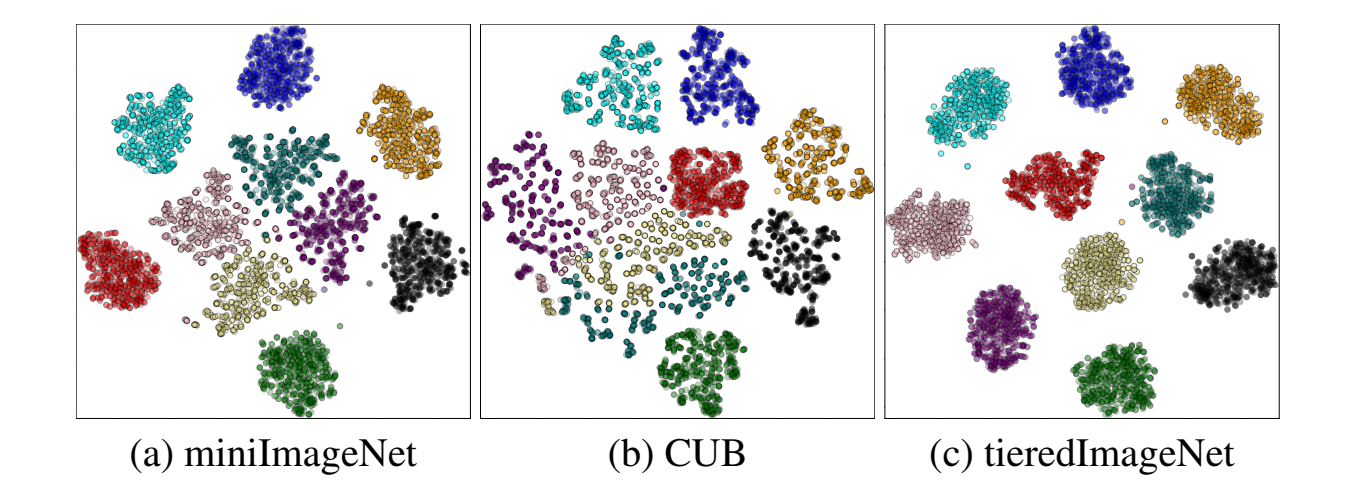
* with Conv4-512

- Top-m analysis.** the results improve as we move towards sum-min, which uses all of the mappers

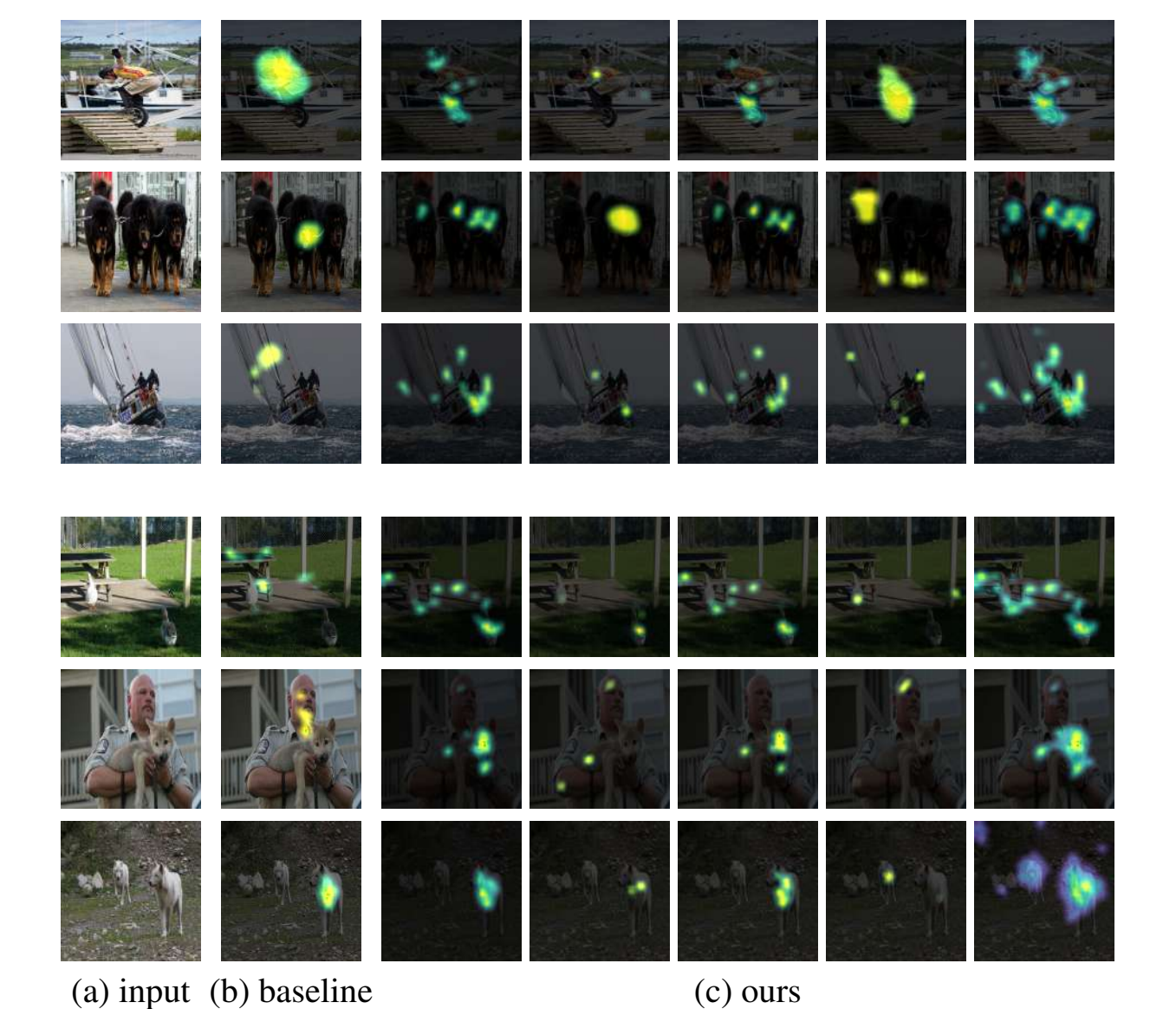
Table 6. Ablation of top-m mapper in the min-sum metric using SetFeat4 and SetFeat12* on CUB. The results are validation set.

Method	SetFeat4	SetFeat12*
	1-shot	5-shot
top-1 (min-min)	70.15	84.94
top-2	70.84	85.30
top-4	70.34	85.95
top-8	71.47	86.88
top-10 (sum-min)	72.09	87.05

- t-SNE visualization.** the distributions of mapper embeddings are generally disjoint



- Visualizing mappers saliency.** our approach devotes attention to many more parts of the images



Acknowledgements. This project was supported by NSERC-Canada and CIFAR Chair.