

Spatial Contrastive Learning for Few-Shot Classification

Yassine Ouali, Céline Hudelot, Myriam Tami
Paris-Saclay University, CentraleSupélec, MICS

Few-shot learning: equip the learner with ability to rapidly learn new concept with few training samples.

?????

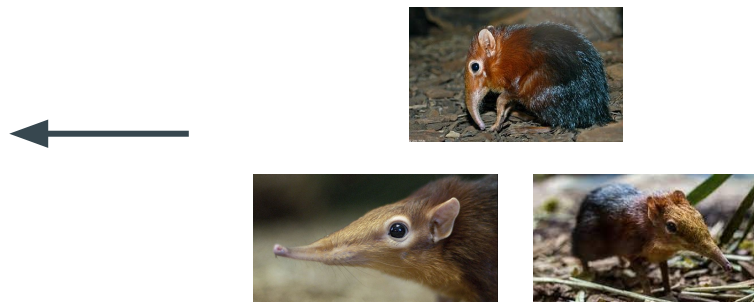


Few-shot learning: equip the learner with ability to rapidly learn new concept with few training samples.

Elephant Shrew



Few training samples



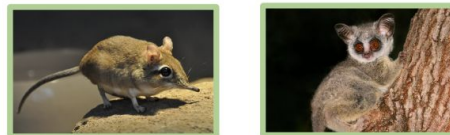
A few-shot classification task: N-way K-shot classification task

- N: number of classes
- K: number of examples per class (small)

5-way (classes) 1-shot (example per class) Task



Train set (support set)

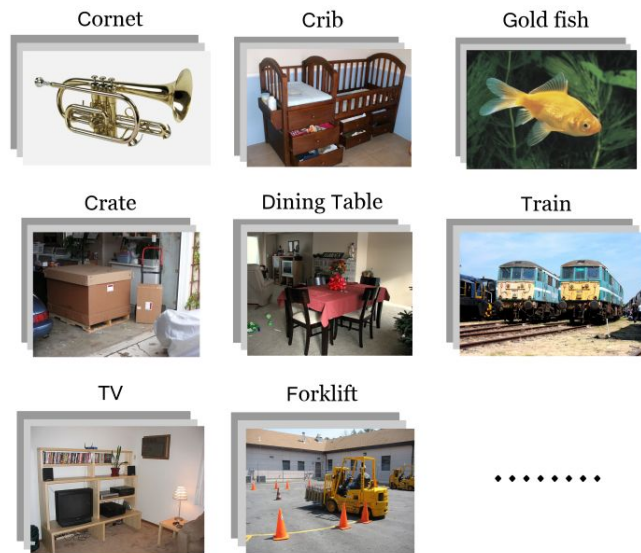


Test set (query set)

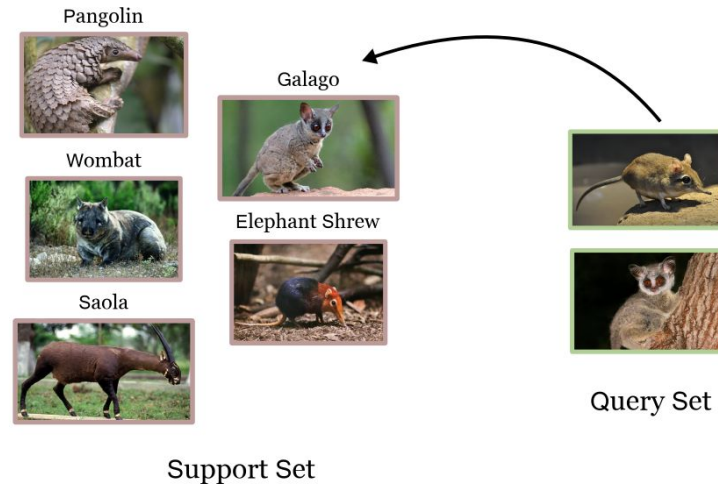
A few-shot classification task: N-way K-shot classification task

- N: number of classes
- K: number of examples per class (small)

Training Phase



Testing Phase



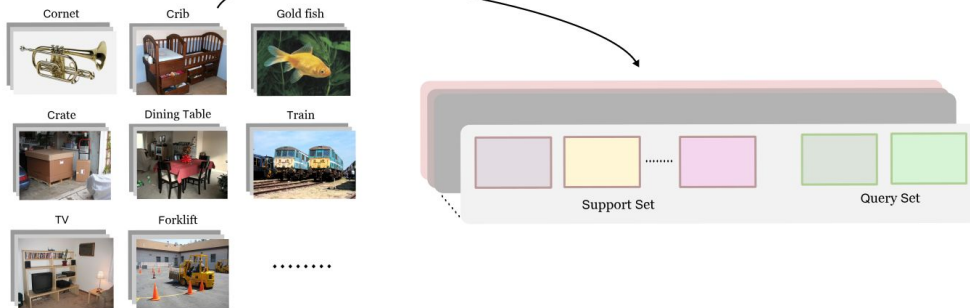
Meta-Learning: initialize then adapt

Training phase has matching conditions compared to the testing phase.

- The training set is transformed into many training tasks.
- The learner is then trained on a distribution of such tasks (episodes).

Training Phase

N-way K-shot



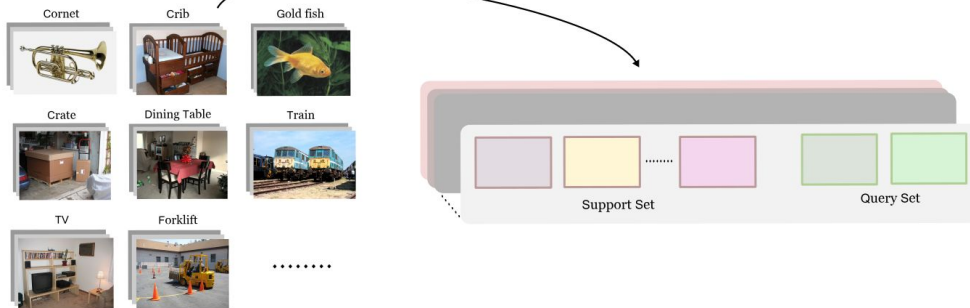
Meta-Learning: initialize then adapt

Training phase has matching conditions compared to the testing phase.

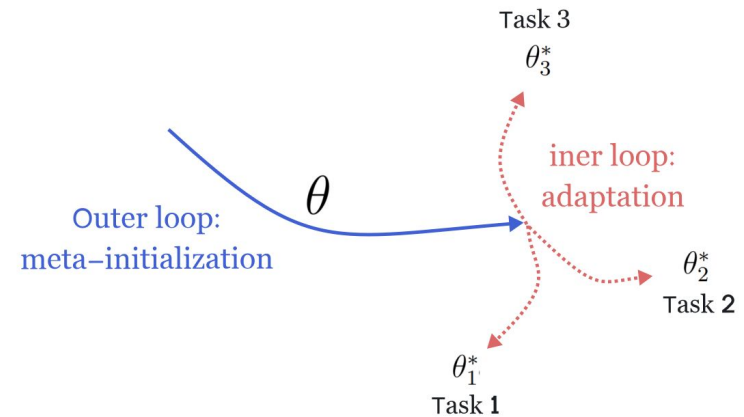
- The training set is transformed into many training tasks.
- The learner is then trained on a distribution of such tasks (episodes).

Training Phase

N-way K-shot



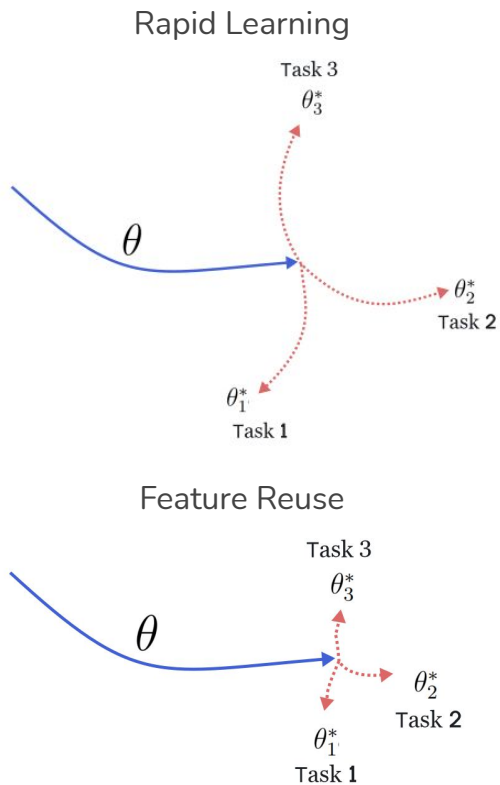
Model-Agnostic Meta-Learning



Meta-Learning: rapid learning or feature reuse

Is this inner/outer loop necessary?

How much of the effectiveness of such methods is contingent on the inner/outer loop structure?

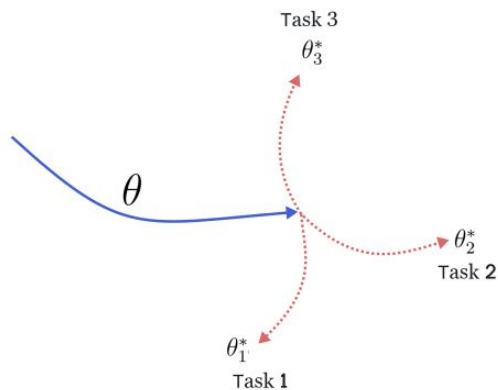


Meta-Learning: rapid learning or feature reuse

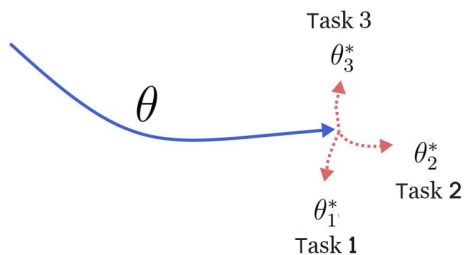
Is this inner/outer loop necessary?

How much of the effectiveness of such methods is contingent on the inner/outer loop structure?

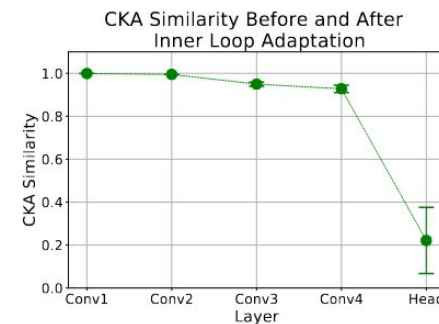
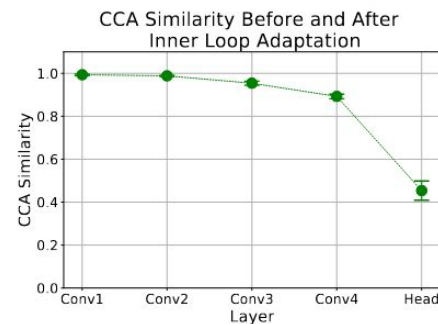
Rapid Learning



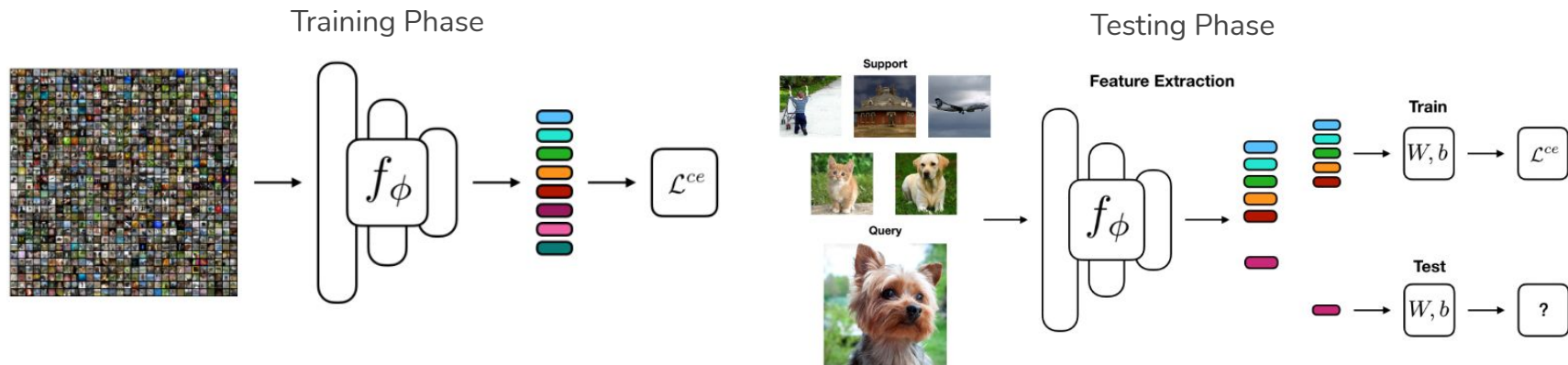
Feature Reuse

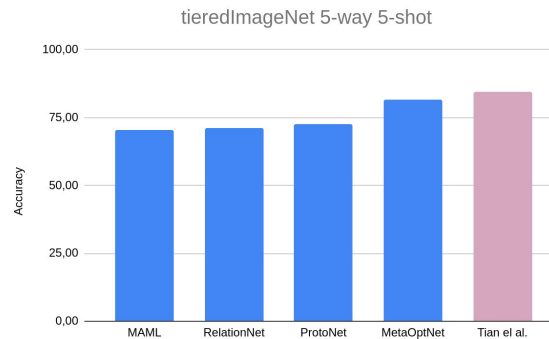
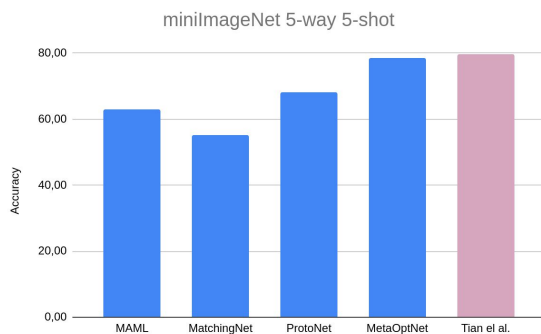
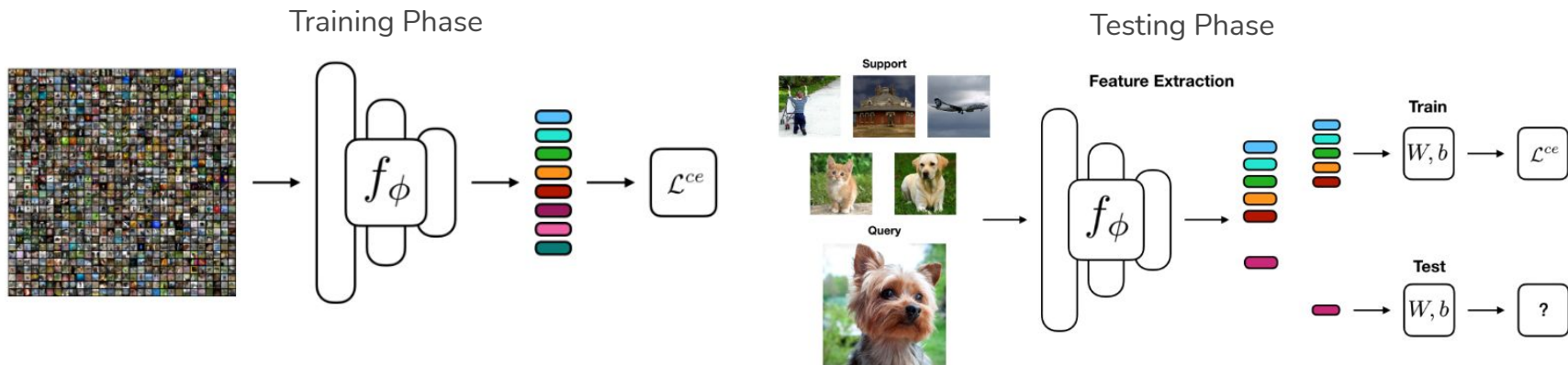


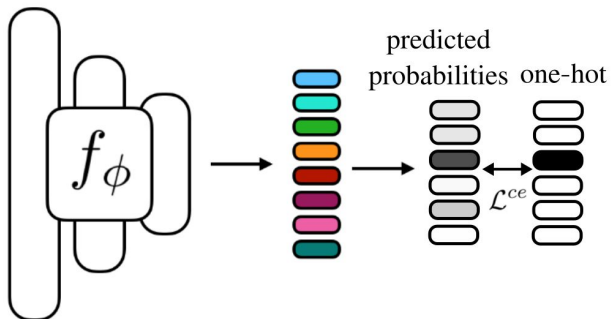
Similarity of the representation before and after adaptation



-> Feature reuse: Having general purpose representations is important for few-shot recognition

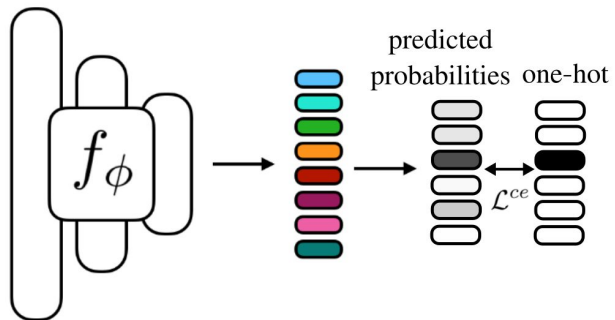






k-Nearest Neighbors Analysis:

- Train on the whole training set with CE loss.
- Find the nearest neighbors of test instances.



k-Nearest Neighbors Analysis:

- Train on the whole training set with CE loss.
- Find the nearest neighbors of test instances.

Test Image



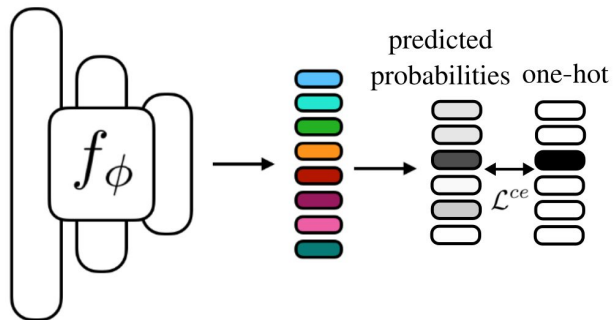
cuirass



golden retriever



hourglass



k-Nearest Neighbors Analysis:

- Train on the whole training set with CE loss.
- Find the nearest neighbors of test instances.

Test Image



cuirass



vase



malamute



hunting dog



hunting dog



malamute



school bus



king cra



malamute



golden retriever



vase



ferret



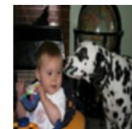
vase



trifle



crate



dalmatian



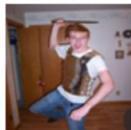
vase



trifle



hourglass



cuirass



theater curtain



vase



electric guitar



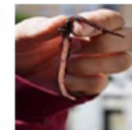
trifle



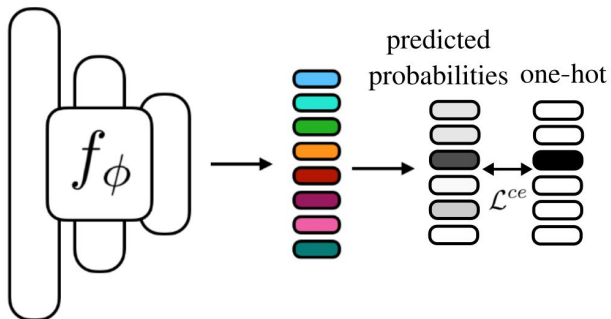
electric guitar



dalmatian



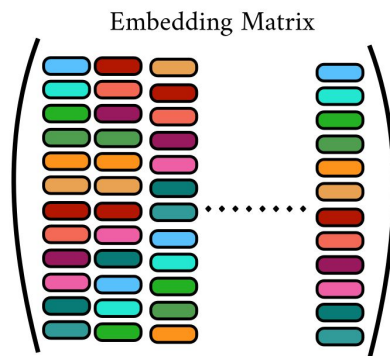
nematode



k-Nearest Neighbors Analysis:

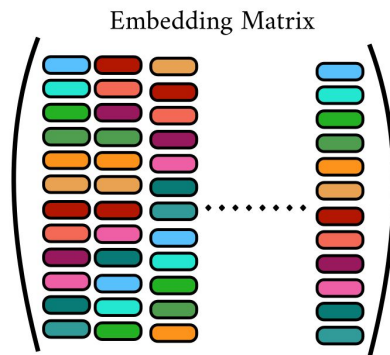
- Train on the whole training set with CE loss.
- Find the nearest neighbors of test instances.

| Test Image | Nearest Neighbors | | | | | | | Predicted class | True class | |
|--|--|--|--|--|---|--|--|---|--|---|
|  cuirass |  vase |  malamute |  hunting dog |  hunting dog |  malamute |  school bus |  king cra |  malamute |  malamute |  cuirass |
|  golden retriever |  vase |  ferret |  vase |  trifle |  crate |  dalmatian |  vase |  trifle |  vase |  golden retriever |
|  hourglass |  cuirass |  theater curtain |  vase |  electric guitar |  trifle |  electric guitar |  dalmatian |  nematode |  electric guitar |  hourglass |



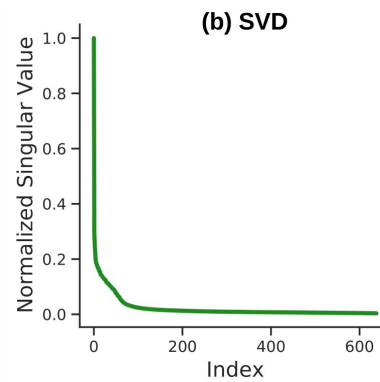
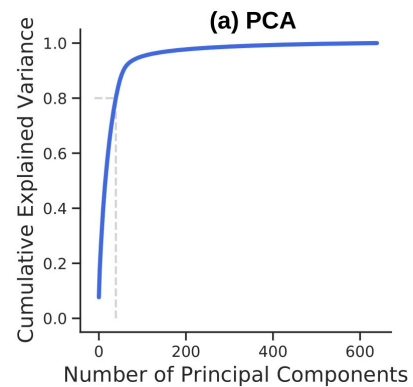
Spectral Analysis:

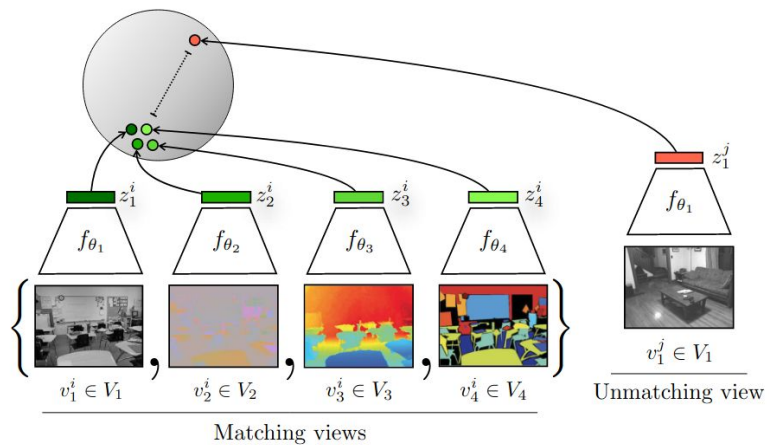
- The explained cumulative variance.
- Normalized singular values.

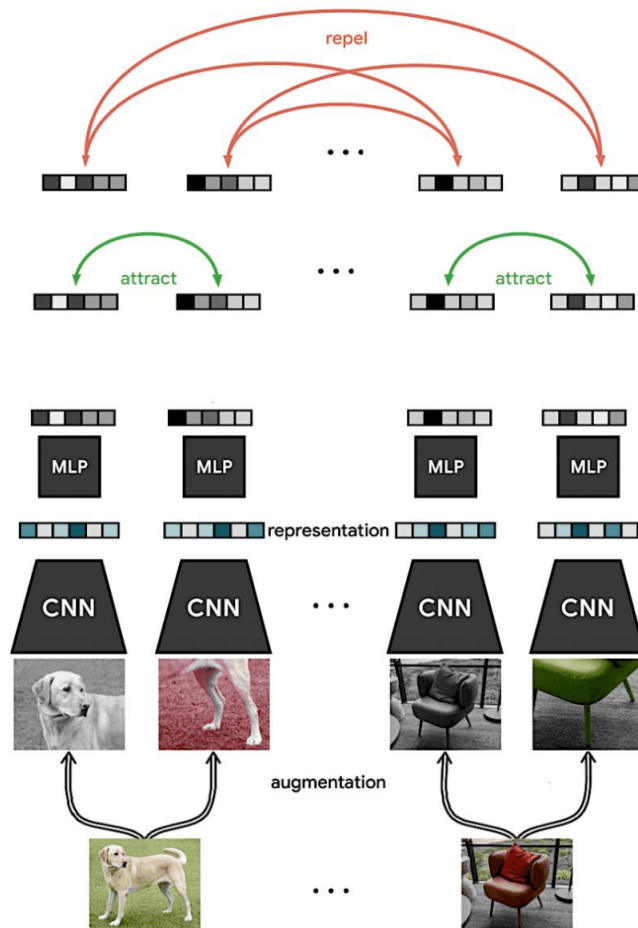
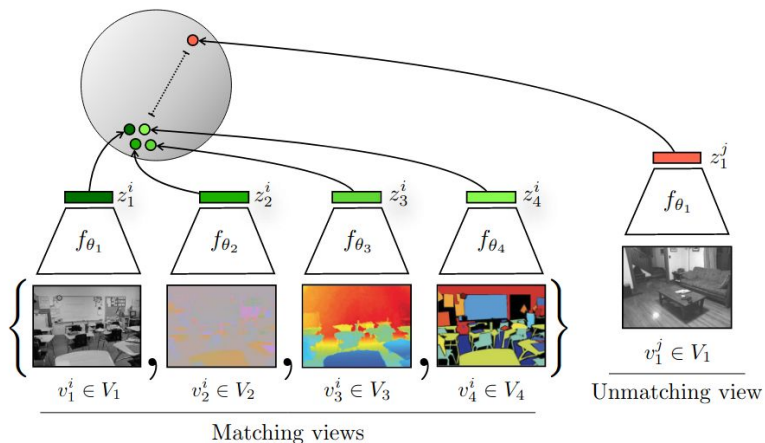


Spectral Analysis:

- The explained cumulative variance.
- Normalized singular values.



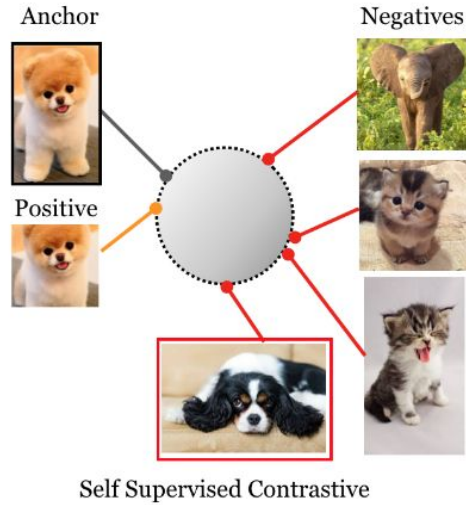


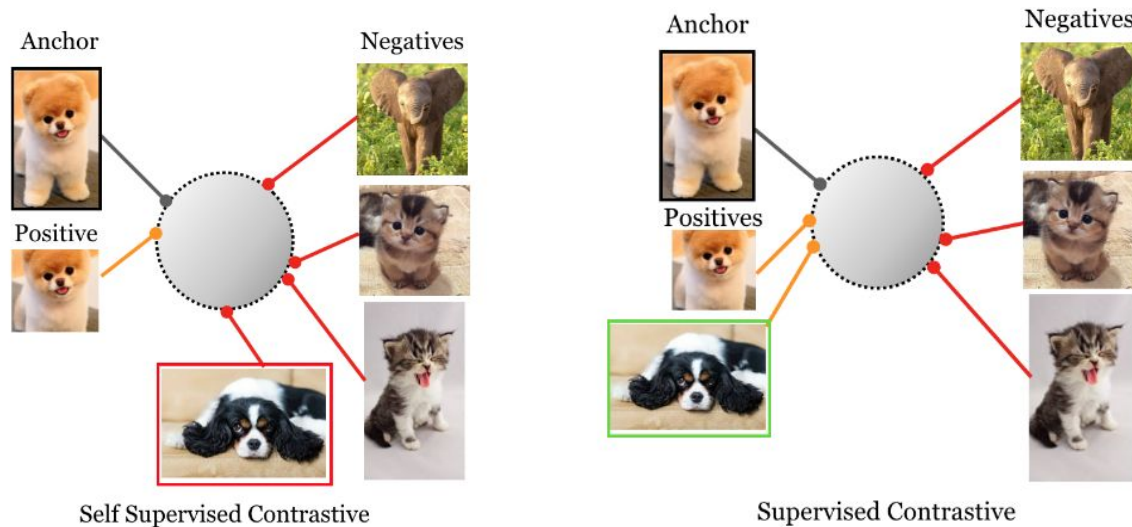


Self-Supervised contrastive loss:

- For each anchor i , attract its positive j against the rest of the examples (negatives).
- The positive examples is a transformed version of the anchor i .

$$-\log \frac{\exp(\text{sim}_g(\mathbf{f}_i, \mathbf{f}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{i \neq k} \cdot \exp(\text{sim}_g(\mathbf{f}_i, \mathbf{f}_k)/\tau)}$$

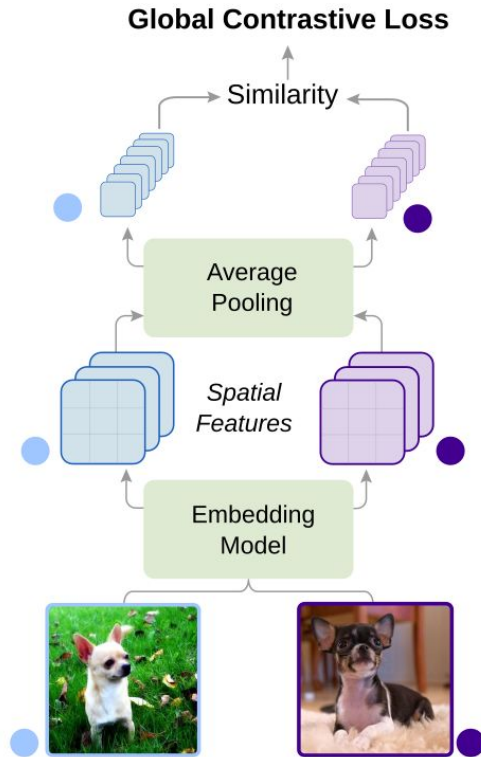


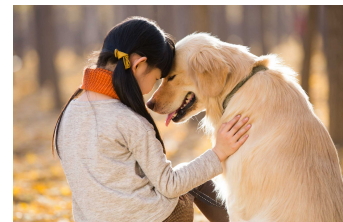
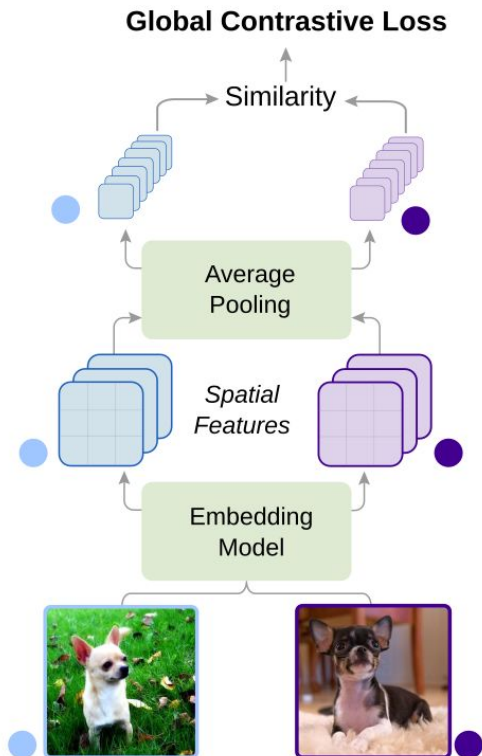


Supervised contrastive loss:

The positives are examples of the same class as the anchor.

The negatives are examples of different class than the anchor.



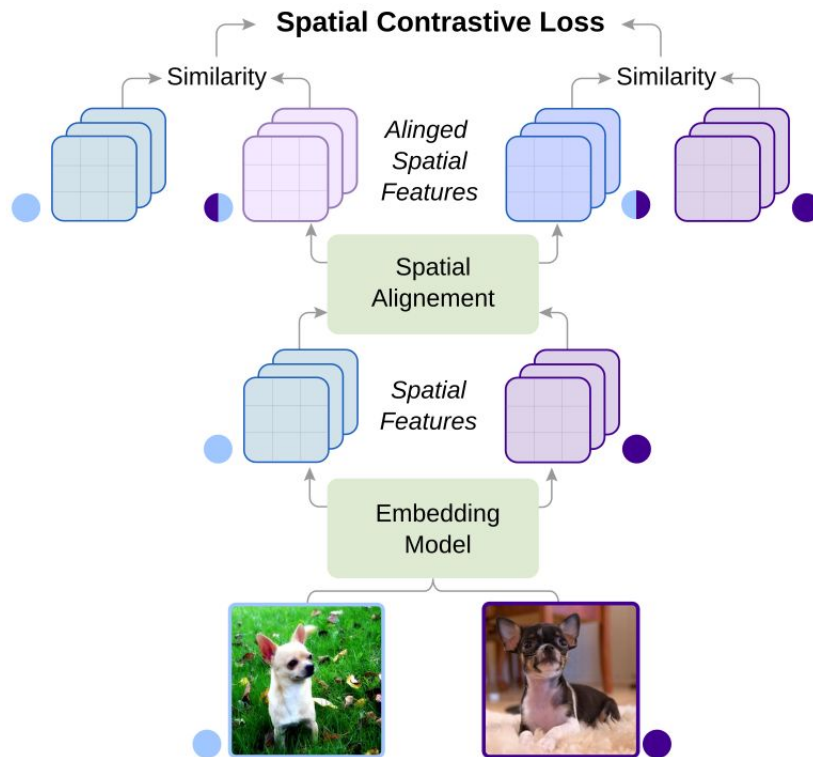
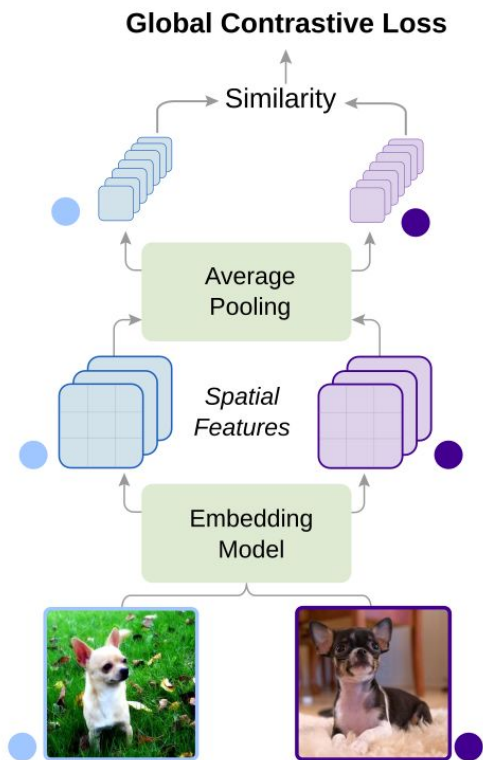


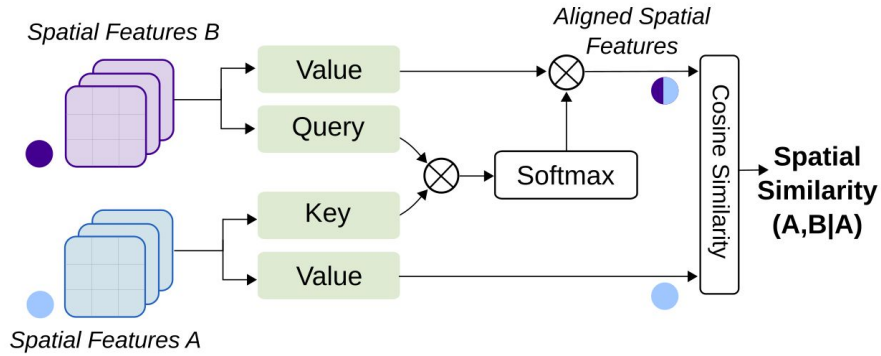
We might align irrelevant objects:

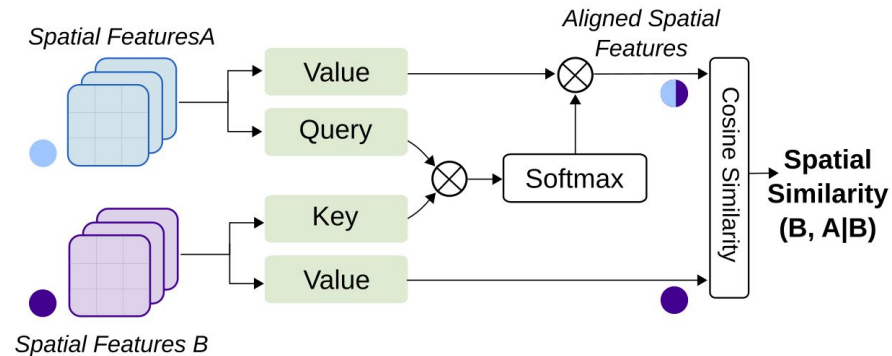
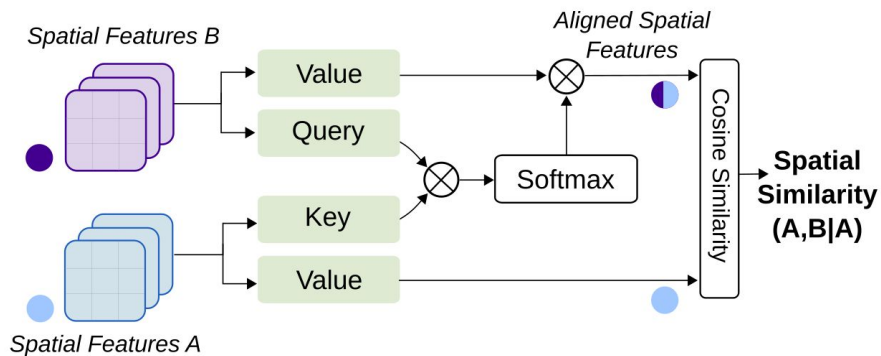
- Dog & person
- Dog & grass

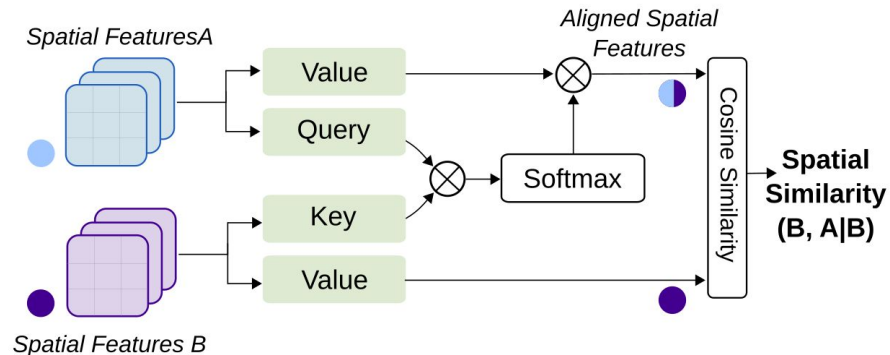
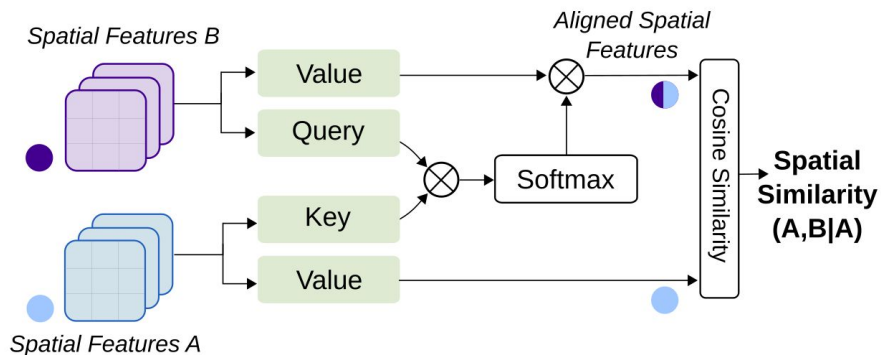


-> Use the spatial features



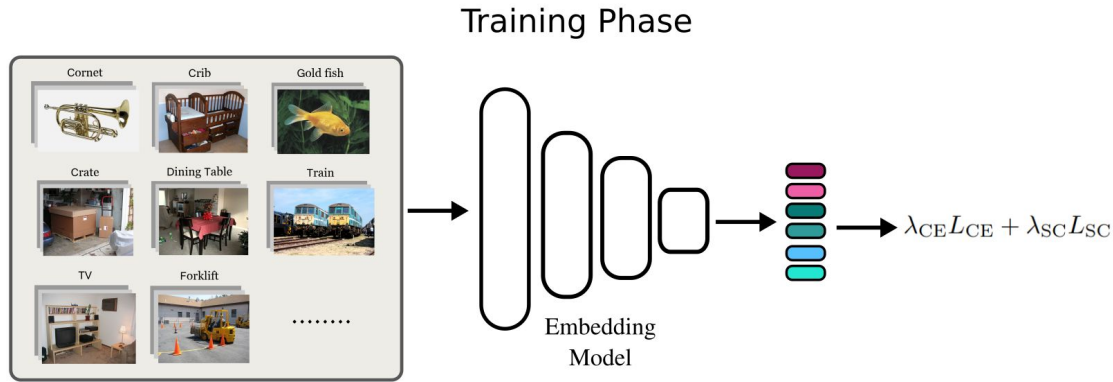


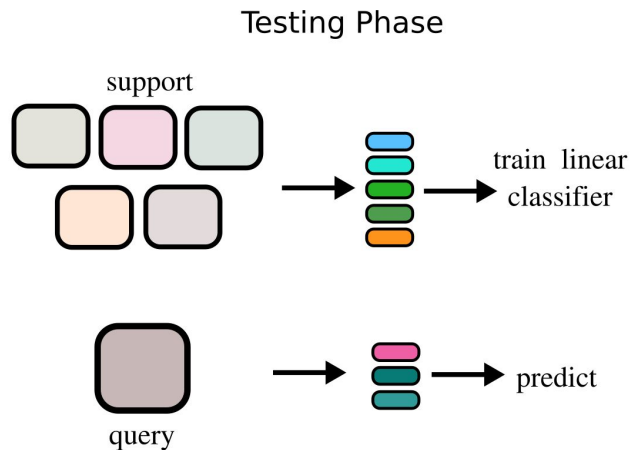
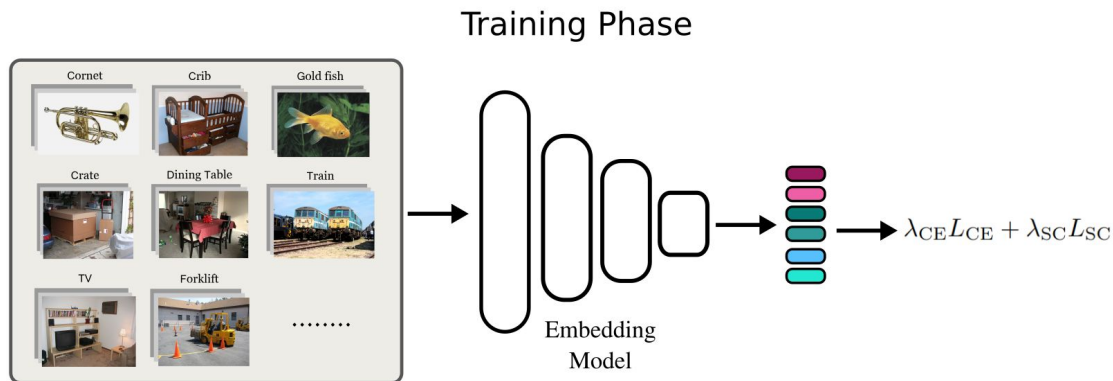




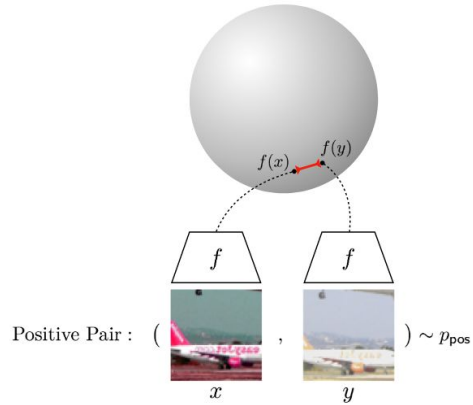
$$\text{Spatial Similarity (A, B)} = \text{Spatial Similarity (A, B|A)} + \text{Spatial Similarity (B, A|B)}$$

$$l_{ij} = -\log \frac{\exp(\text{sim}_s(\mathbf{z}_i^s, \mathbf{z}_j^s)/\tau')}{\sum_{k=1}^{2N} \mathbb{1}_{i \neq k} \cdot \exp(\text{sim}_s(\mathbf{z}_i^s, \mathbf{z}_k^s)/\tau')},$$

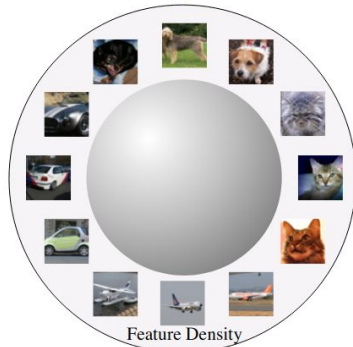




Contrastive Distillation: an additional training step

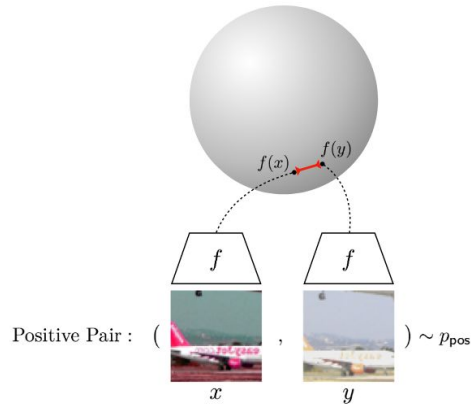


Alignment: Similar samples have similar features.

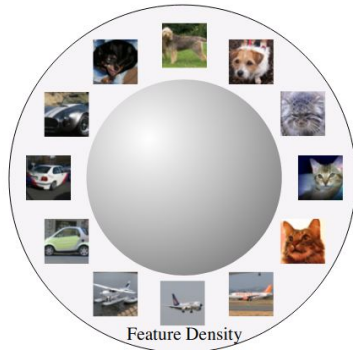


Uniformity: Preserve maximal information.

Contrastive Distillation: an additional training step



Alignment: Similar samples have similar features.



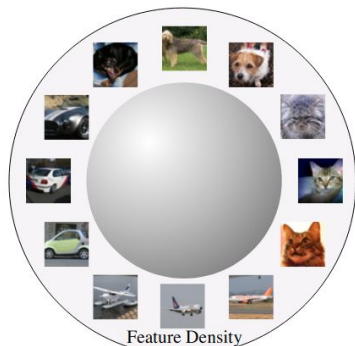
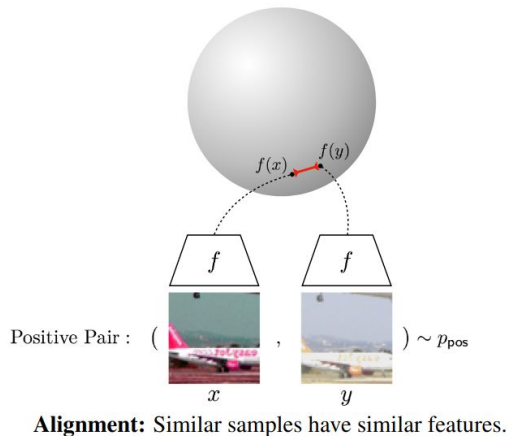
Uniformity: Preserve maximal information.

A possible over clustering of the features of the same class.

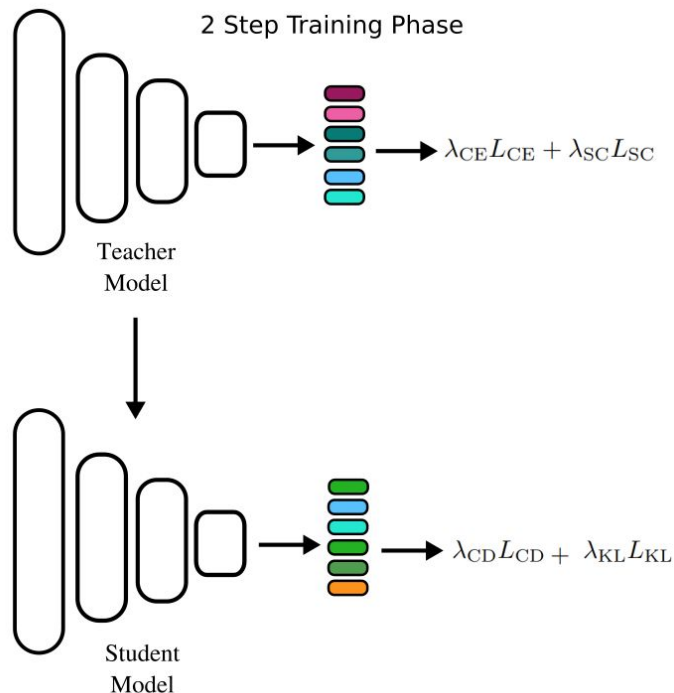
Contrastive Distillation: an additional training step

Possible over clustering of the features of the same class.

-> Contrastive alignment distillation without uniformity
$$L_{CD} = \frac{1}{N} \sum_{i=1}^N \|\mathbf{z}_i^{gt} - \mathbf{z}_i^{gs}\|_2^2$$



Uniformity: Preserve maximal information.



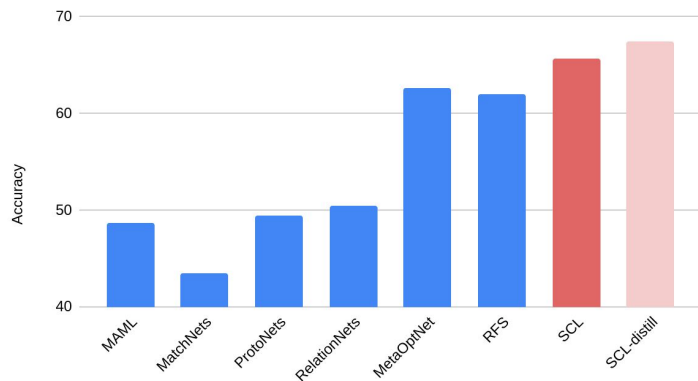
| Loss Function | Aug. | <i>mini</i> -ImageNet, 5-way | | CIFAR-CS, 5-way | |
|--------------------|------|------------------------------|-------------------|-------------------|-------------------|
| | | 1-shot | 5-shot | 1-shot | 5-shot |
| CE | | 61.8 ± 0.7 | 79.7 ± 0.6 | 71.3 ± 0.9 | 86.1 ± 0.6 |
| CE | ✓ | 61.8 ± 0.8 | 78.6 ± 0.5 | 71.9 ± 0.9 | 86.3 ± 0.5 |
| CE + SS-GC | ✓ | 62.7 ± 0.7 | 81.0 ± 0.6 | 70.9 ± 0.9 | 84.5 ± 0.6 |
| CE + SS-SC | ✓ | 64.0 ± 0.8 | 81.5 ± 0.5 | 72.1 ± 0.8 | 86.2 ± 0.6 |
| CE + SS-GC + SS-SC | ✓ | 62.8 ± 0.8 | 81.1 ± 0.6 | 69.0 ± 0.9 | 85.0 ± 0.6 |
| CE + GC | ✓ | 65.0 ± 0.8 | 81.6 ± 0.5 | 74.0 ± 0.8 | 87.3 ± 0.6 |
| CE + SC | ✓ | 65.7 ± 0.8 | 82.5 ± 0.5 | 75.0 ± 0.9 | 87.4 ± 0.6 |
| CE + GC + SC | ✓ | 65.0 ± 0.8 | 81.3 ± 0.5 | 76.0 ± 0.7 | 87.5 ± 0.5 |

-> Optimizing the spatial features yields in better results.

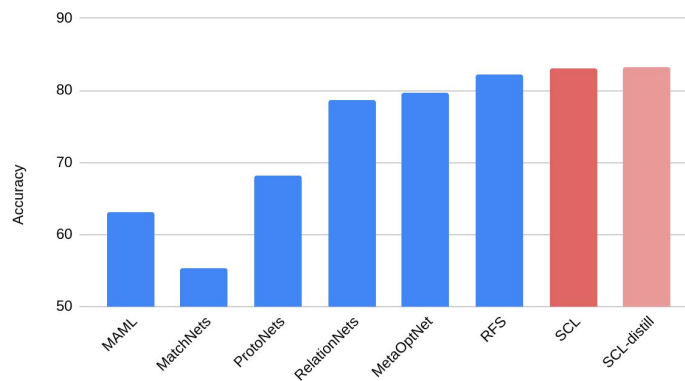
| Loss Function | Aug. | <i>mini</i> -ImageNet, 5-way | | CIFAR-CS, 5-way | |
|--------------------|------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|
| | | 1-shot | 5-shot | 1-shot | 5-shot |
| CE | | 61.8 ± 0.7 | 79.7 ± 0.6 | 71.3 ± 0.9 | 86.1 ± 0.6 |
| CE | ✓ | 61.8 ± 0.8 | 78.6 ± 0.5 | 71.9 ± 0.9 | 86.3 ± 0.5 |
| CE + SS-GC | ✓ | 62.7 ± 0.7 | 81.0 ± 0.6 | 70.9 ± 0.9 | 84.5 ± 0.6 |
| CE + SS-SC | ✓ | 64.0 ± 0.8 | 81.5 ± 0.5 | 72.1 ± 0.8 | 86.2 ± 0.6 |
| CE + SS-GC + SS-SC | ✓ | 62.8 ± 0.8 | 81.1 ± 0.6 | 69.0 ± 0.9 | 85.0 ± 0.6 |
| CE + GC | ✓ | 65.0 ± 0.8 | 81.6 ± 0.5 | 74.0 ± 0.8 | 87.3 ± 0.6 |
| CE + SC | ✓ | 65.7 ± 0.8 | 82.5 ± 0.5 | 75.0 ± 0.9 | 87.4 ± 0.6 |
| CE + GC + SC | ✓ | 65.0 ± 0.8 | 81.3 ± 0.5 | 76.0 ± 0.7 | 87.5 ± 0.5 |

-> Optimizing the spatial features yields in better results.

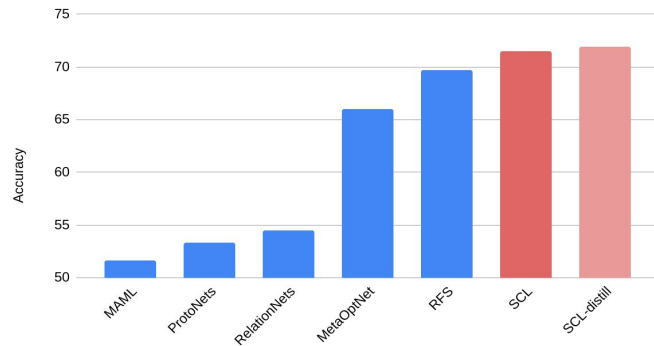
mini-ImageNet, 5-way, 1-shot



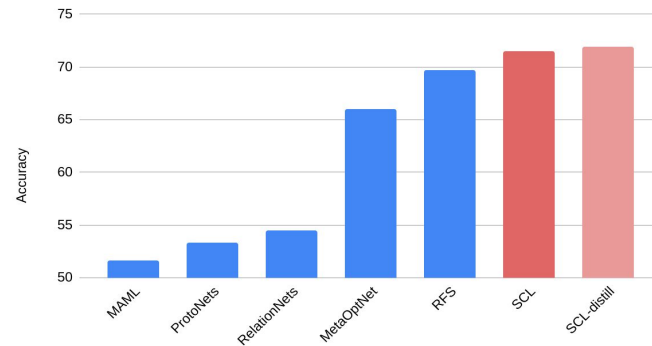
mini-ImageNet, 5-way, 5-shot



tiered-ImageNet, 5-way, 1-shot

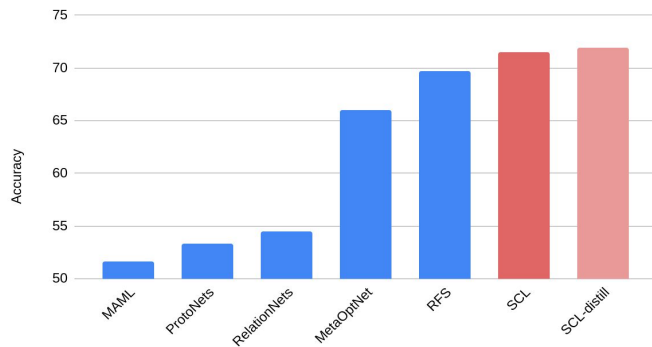


tiered-ImageNet, 5-way, 1-shot

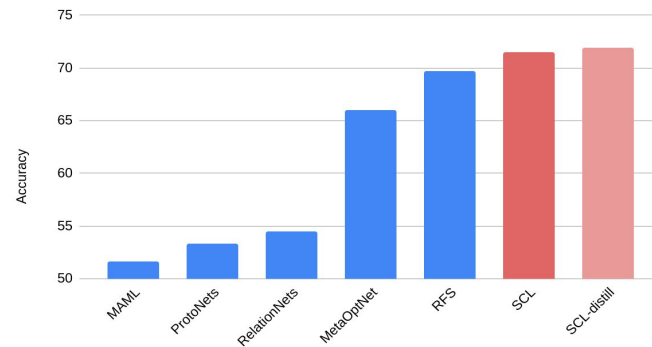


Results: tieredImageNet & cross-domain few-shot

tiered-ImageNet, 5-way, 1-shot

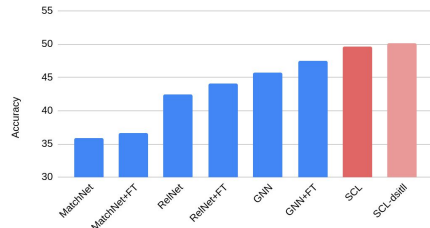


tiered-ImageNet, 5-way, 1-shot

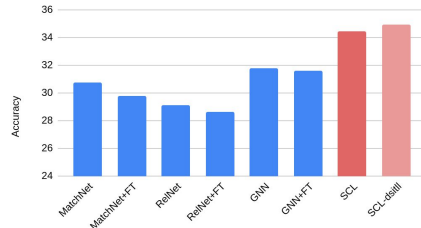


Cross-domain few-shot classification: Train on minilImageNet and test on a new domain (Cub, Cars, Places, or Plantae)

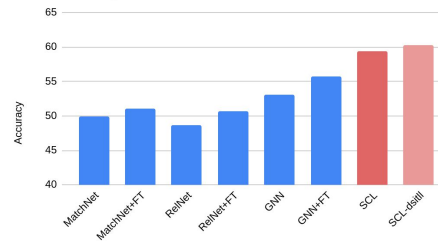
CUB, 5-way, 1-shot



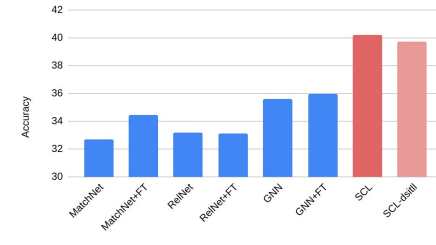
Cars, 5-way, 1-shot

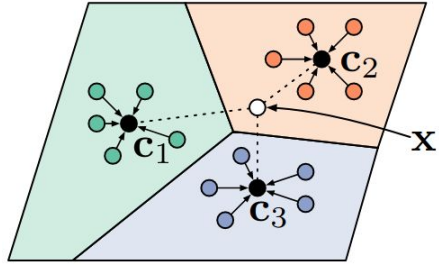


Places, 5-way, 1-shot

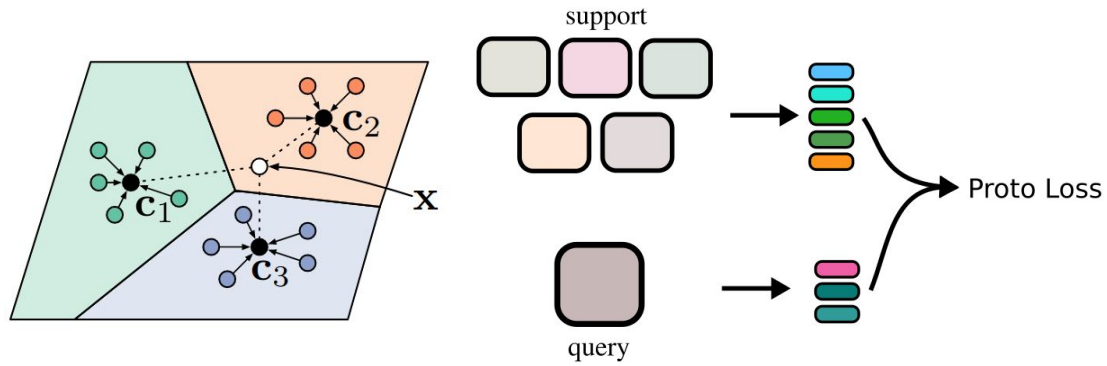


Plantae, 5-way, 1-shot

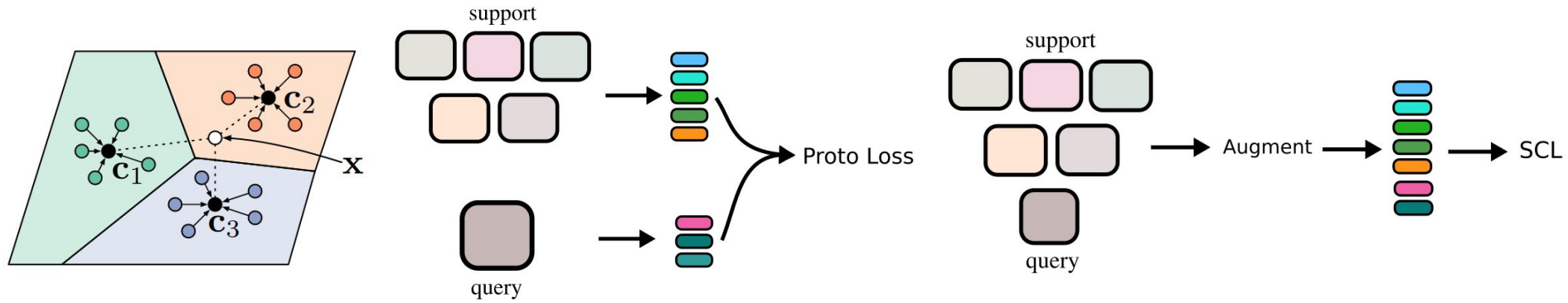




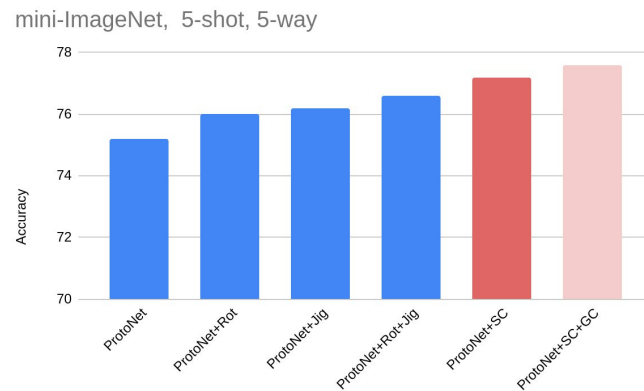
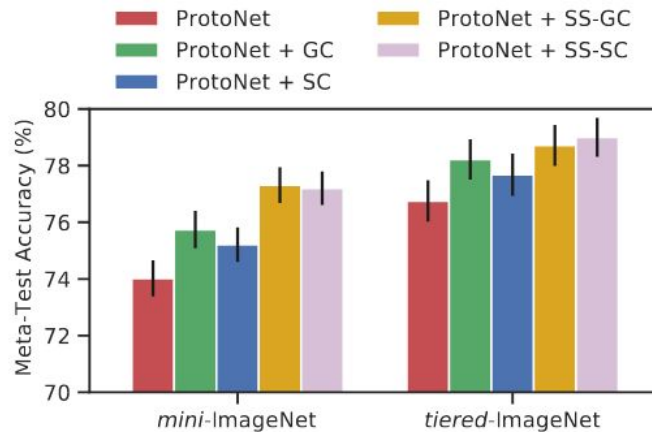
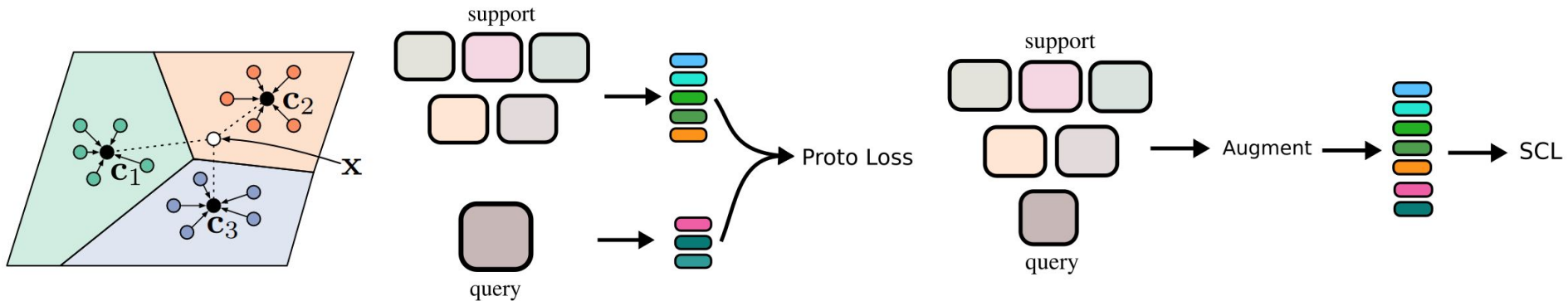
Results: prototypical networks



Results: prototypical networks



Results: prototypical networks



In this work, we:

- Explored the usage of contrastive learning for few-Shot classification to learn more general purpose representations.
- Proposed a novel Spatial contrantrasitive loss to further promotes class-independent discriminative patterns.
- Proposed a contrastive distillation step to relax the placement of the embeddings and enhance the learned representations.
- Demonstrated the effectiveness of the proposed method on different datasets, settings and frameworks.

For more details please see the [paper](#) / [code](#).

Thank You!

For more details please visit the project's webpage:



https://vassouali.github.io/SCL_page/

Paris-Saclay University, CentraleSupélec, MICS