



UNIVERSITÉ
JEAN MONNET
SAINT-ÉTIENNE

IMPROVING FEW-SHOT LEARNING THROUGH MULTI-TASK REPRESENTATION LEARNING THEORY



Quentin BOUNIOT





INTRODUCTION



Apple



Blueberry



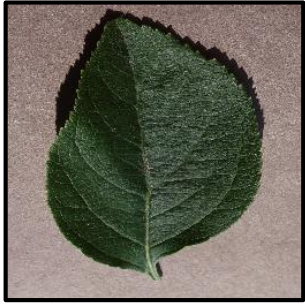
Apple



Blueberry



?



Apple



Blueberry



?





Apple



Blueberry



?



Da Vinci



Botero



Apple



Blueberry



?



Da Vinci



Botero



?



Apple



Blueberry



?



Da Vinci



Botero



?

Meta-learning = Learning to Learn

- Meta-learning 101
- Multi-task Representation Learning Theory
- From Theory to Practice
- Take Home Message



META-LEARNING 101





META-LEARNING 101

- What is Meta-learning ?
 - ▶ Meta-Training : solve a set of **source tasks**
 - +
 - ▶ Meta-Testing : use knowledge from meta-training to solve **previously unseen tasks** more efficiently

- What is Meta-learning ?
 - ▶ Meta-Training : solve a set of source tasks
 - +
 - ▶ Meta-Testing : use knowledge from meta-training to solve previously unseen tasks more efficiently
- ▶ Meta-Learning can be used for a lot of problems (classification, regression, RL, ...)

- What is Meta-learning ?
 - ▶ Meta-Training : solve a set of source tasks
 - +
 - ▶ Meta-Testing : use knowledge from meta-training to solve previously unseen tasks more efficiently
- ▶ Meta-Learning can be used for a lot of problems (classification, regression, RL, ...)
- How is it related to Few-shot Learning ?
 - ▶ The meta-learner *learns to learn* a new task with few shots.

INTRODUCING EPISODES

Training Support Set

Testing Query Set

Meta-Training



Meta-Testing



INTRODUCING EPISODES

Training Support Set

Testing Query Set

Meta-Training



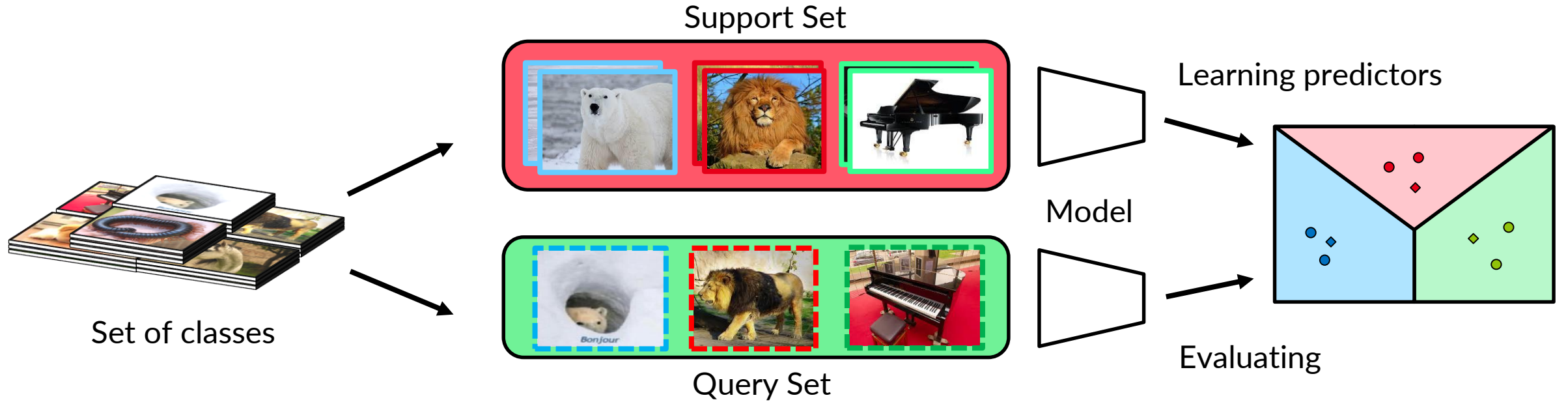
← Episode i

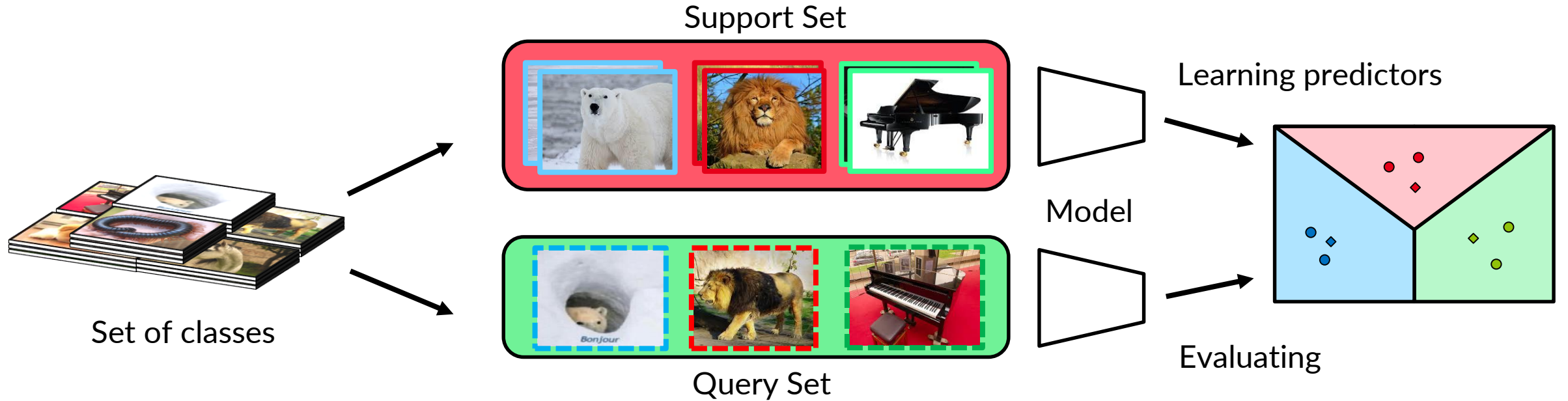
← Episode $i + 1$

Meta-Testing



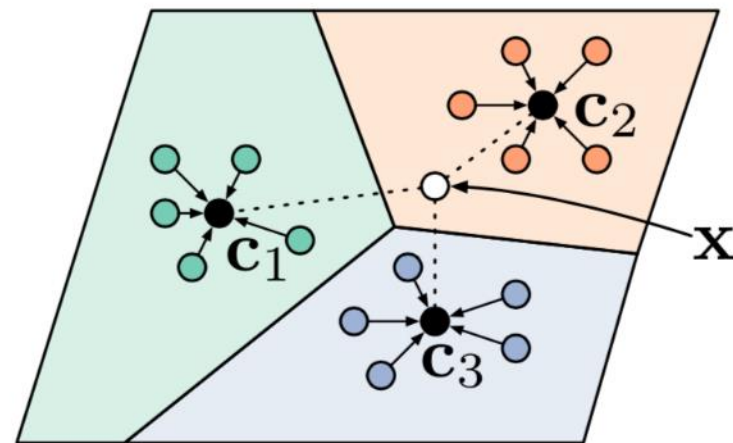
► *N*-way *k*-shot episode: task with *N* different classes and *k* images for each class.



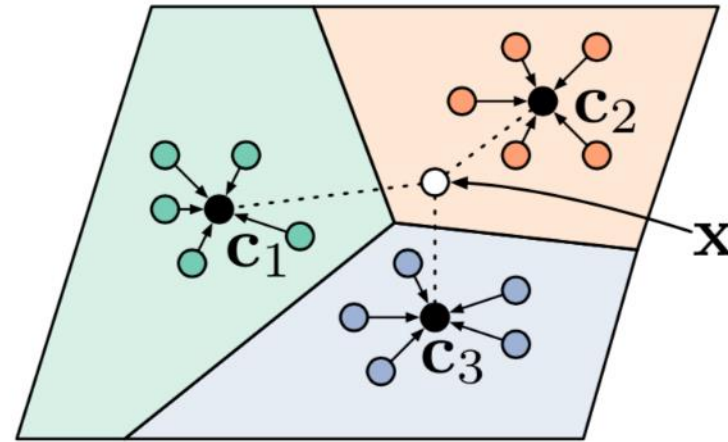


- Disjoint sets of classes between meta-training and meta-testing classes
- Construction of episodes from dataset

METHODS I: METRIC-BASED PROTOTYPICAL NETWORK (PROTONET)



Snell J. et al. (2017), *Prototypical Networks for Few-shot Learning*. In NeurIPS 2017.
 Allen K. et al. (2019), *Infinite Mixture Prototypes for few-shot learning*. In ICML 2019.

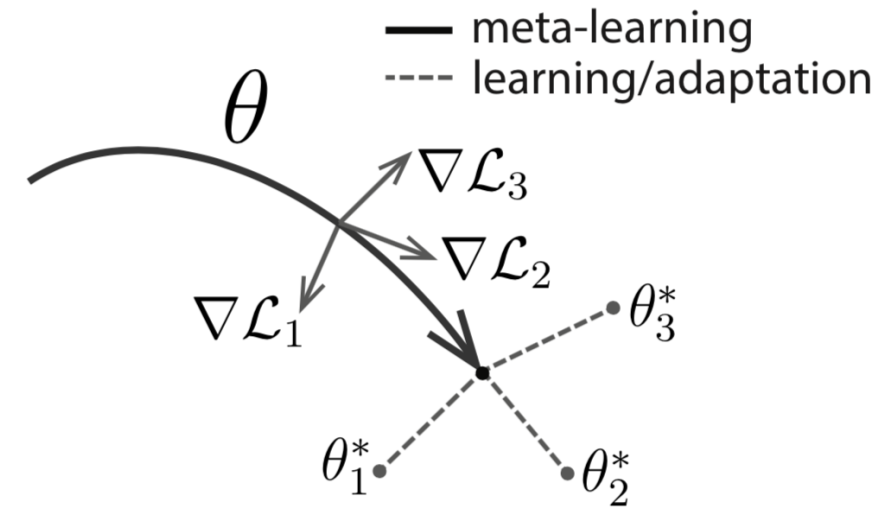
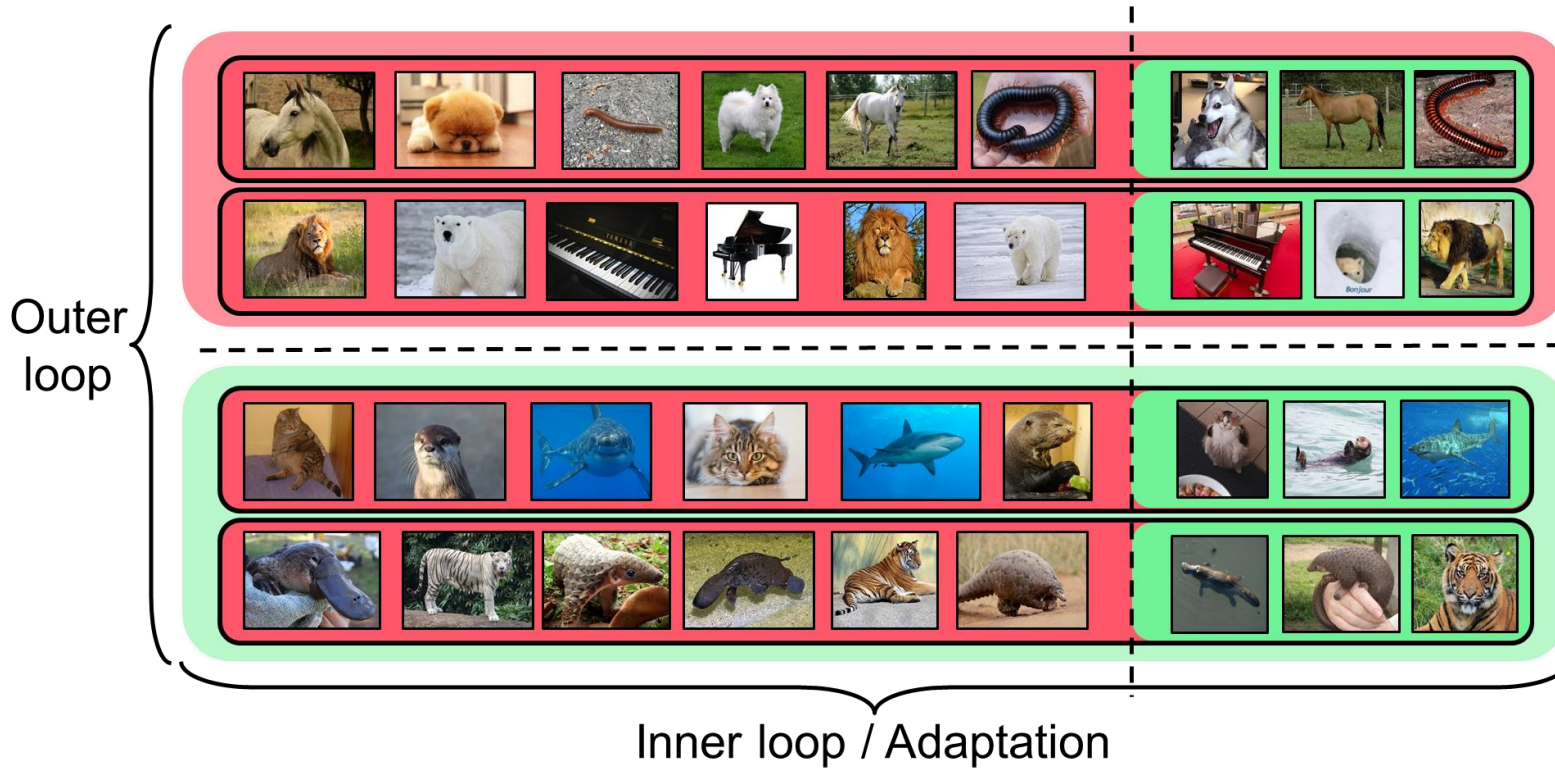


- **Embedding function** to encode query and support samples.
- Support samples fused into **prototypes** c_i for each class
- Probability distribution using **inverse of distances** to prototypes.
- **Contrastive loss** according to distance function.

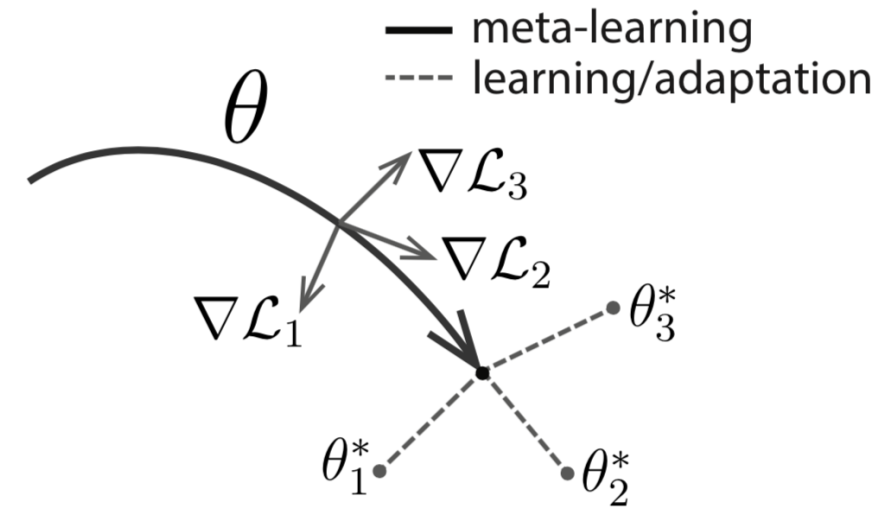
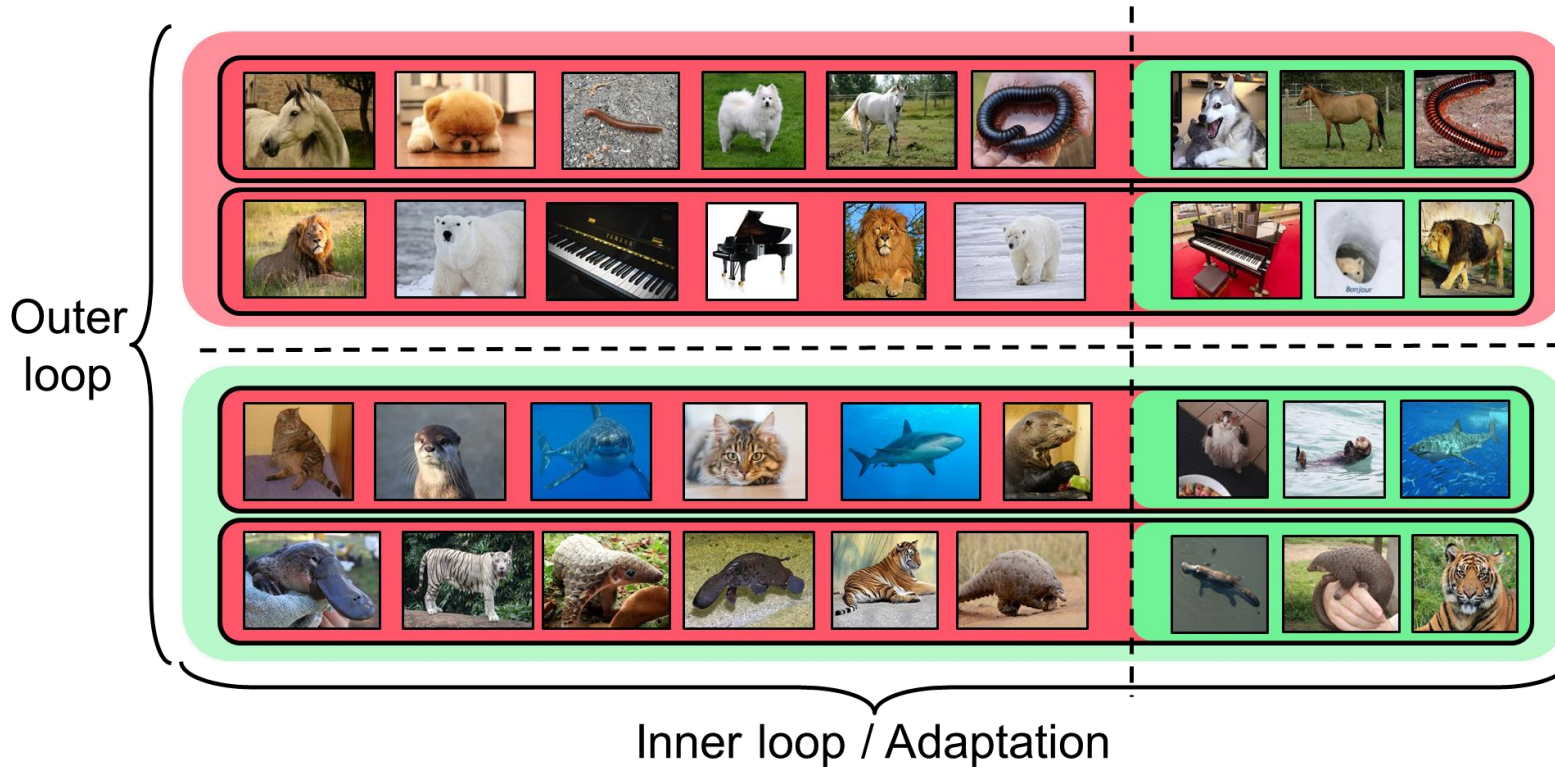
Snell J. et al. (2017), *Prototypical Networks for Few-shot Learning*. In NeurIPS 2017.

Allen K. et al. (2019), *Infinite Mixture Prototypes for few-shot learning*. In ICML 2019.

METHODS II: GRADIENT-BASED MODEL AGNOSTIC META-LEARNING (MAML)



METHODS II: GRADIENT-BASED MODEL AGNOSTIC META-LEARNING (MAML)



Inner Loop:

- Performs a few gradient updates over the k labelled examples (the support set) of **current episode/task**.

Outer Loop:

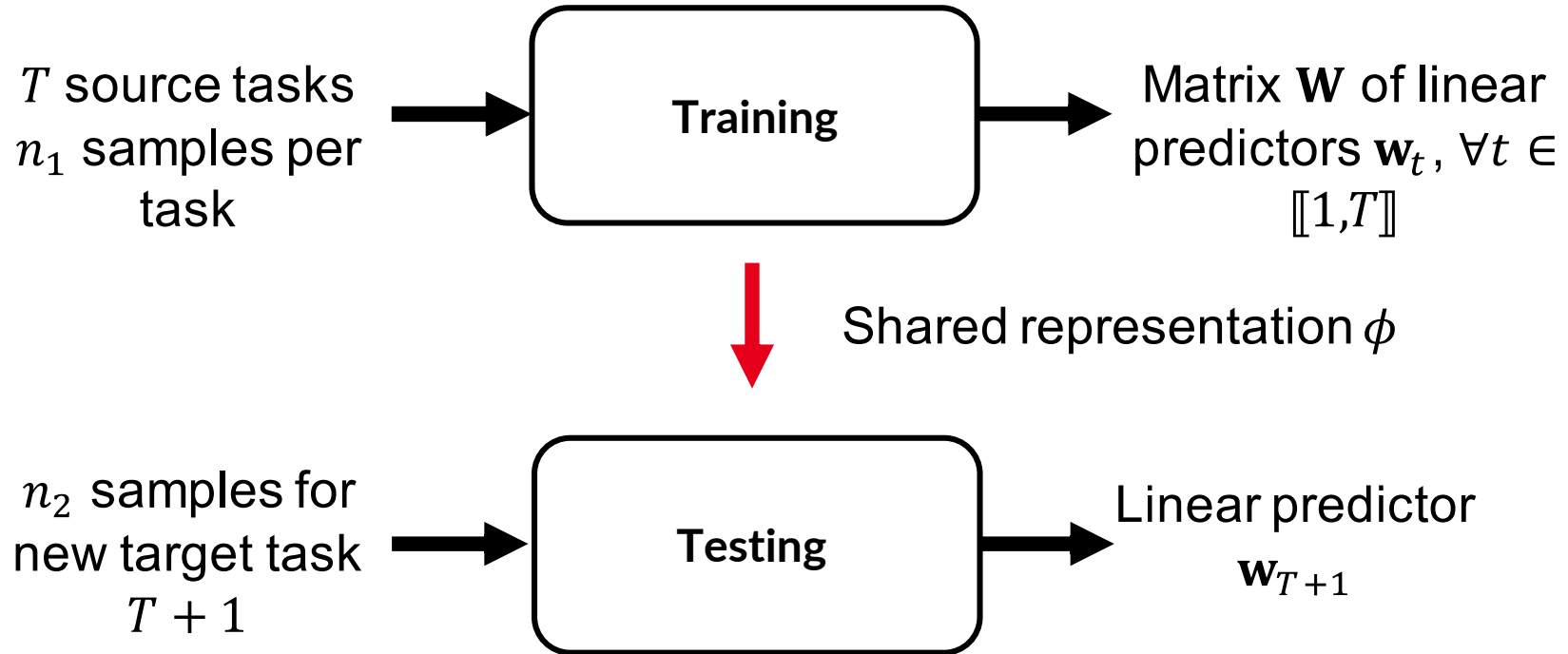
- Updates the **initialization** of the parameters (often called the *meta-initialization*).

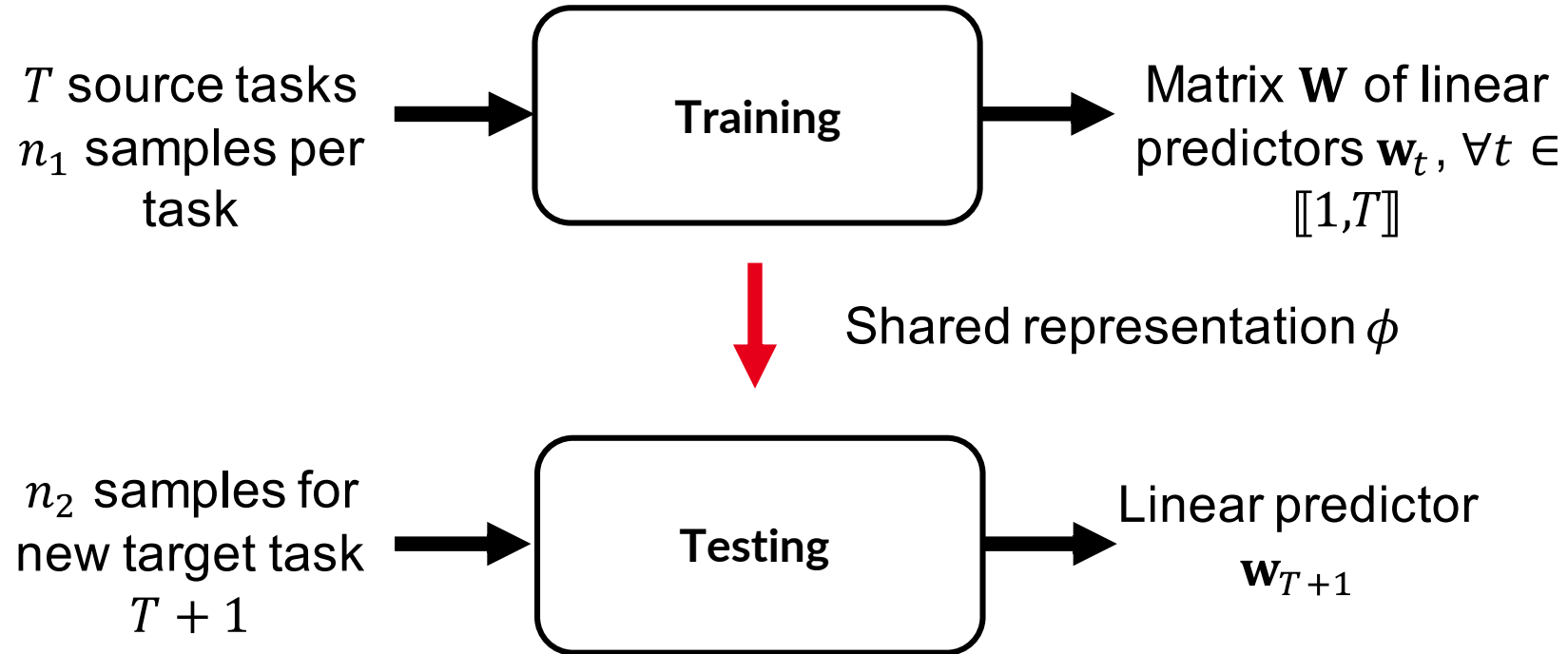


MULTI-TASK REPRESENTATION LEARNING THEORY







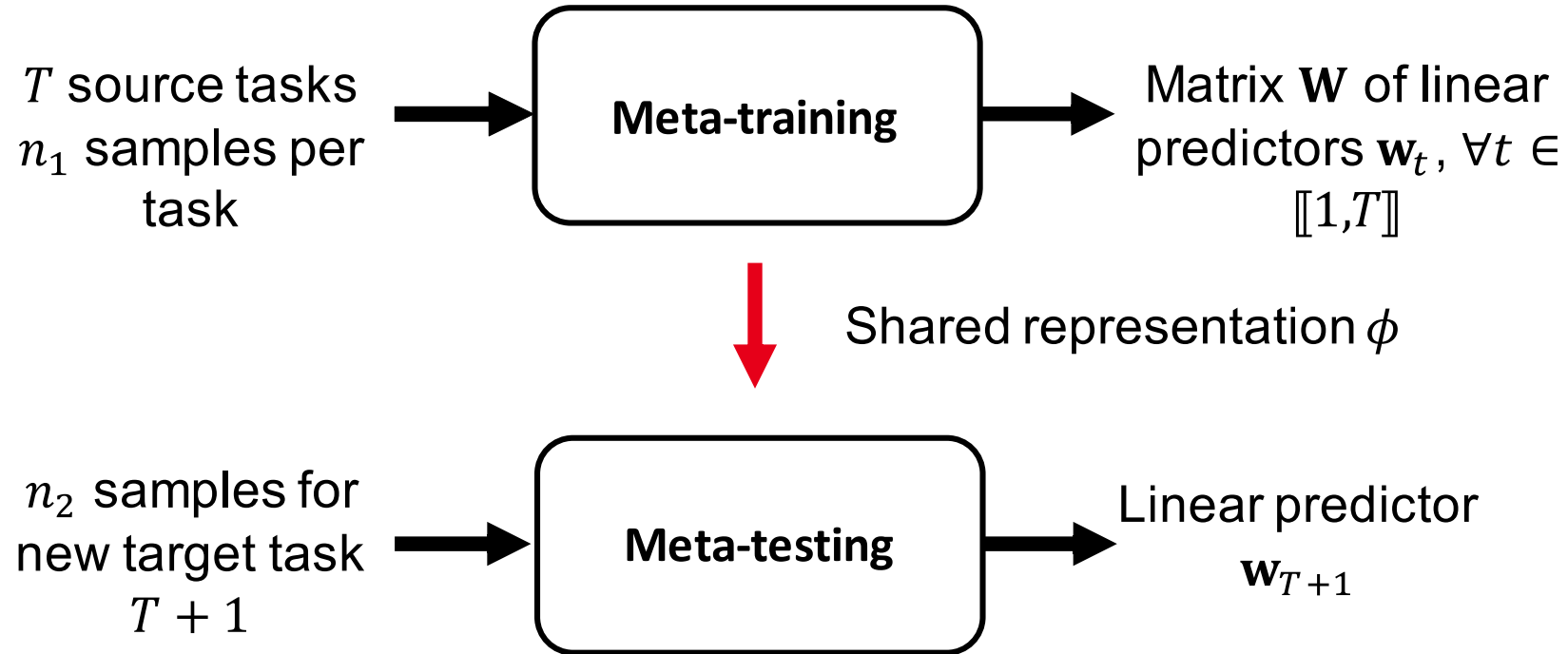


Goal: Minimize *excess risk* $ER = \mathcal{L}(\hat{\phi}, \hat{\mathbf{w}}_{T+1}) - \mathcal{L}(\phi^*, \mathbf{w}_{T+1}^*)$

► True risk \mathcal{L}

► Optimal weights ϕ^*

► \mathbf{w}_{T+1}^* ideal target linear predictor



Goal: Minimize *excess risk* $ER = \mathcal{L}(\hat{\phi}, \hat{\mathbf{w}}_{T+1}) - \mathcal{L}(\phi^*, \mathbf{w}_{T+1}^*)$

► True risk \mathcal{L}

► Optimal weights ϕ^*

► \mathbf{w}_{T+1}^* ideal target linear predictor

- **Multi-task training \neq Episodic training**
 - ▶ Mismatch in problem formulation and objectives

Wang H. et al. (2021), *Bridging Multi-Task Learning and Meta-Learning: Towards Efficient Training and Effective Adaptation* In ICML 2021.

- **Multi-task training \neq Episodic training**
 - ▶ Mismatch in problem formulation and objectives

- **But shared optimization formulation, with some simplification**
 - ▶ The differences are empirically negligible.

Wang H. et al. (2021), *Bridging Multi-Task Learning and Meta-Learning: Towards Efficient Training and Effective Adaptation* In ICML 2021.

- Traditional PAC-bounds

$$\mathbb{E}R(\phi, \mathbf{w}_{T+1}) \leq O\left(\frac{1}{n_1} + \frac{1}{T}\right)$$

- × Requires n_1 and T to tend to infinity.
- × Doesn't explain the success in **few data regime**.

Du S. et al. (2020), *Few-Shot Learning via Learning the Representation, Provably*. In ICRL 2021
Tripuraneni N. et al. (2020). *Provable Meta-Learning of Linear Representations*. In arXiv 2020.

- Assumption 1: Diversity of the source tasks
 - ▶ Optimal predictors $\mathbf{W}^* = [\mathbf{w}_1^*, \dots, \mathbf{w}_T^*]$ cover all the directions evenly

▶ Condition Number $\kappa(\mathbf{W}^*) = \frac{\sigma_{max}(\mathbf{W}^*)}{\sigma_{min}(\mathbf{W}^*)}$ should not increase with T

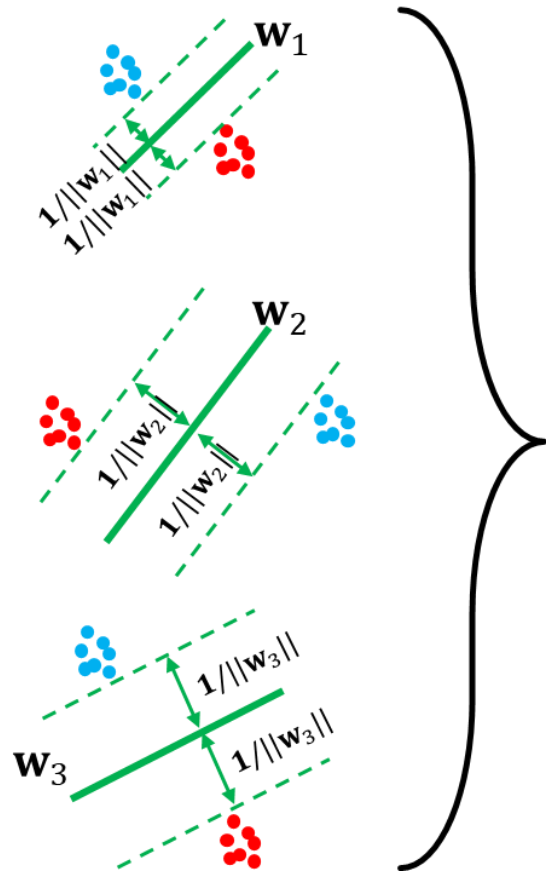
- Assumption 1: Diversity of the source tasks
 - ▶ Optimal predictors $\mathbf{W}^* = [\mathbf{w}_1^*, \dots, \mathbf{w}_T^*]$ cover all the directions evenly

▶ Condition Number $\kappa(\mathbf{W}^*) = \frac{\sigma_{max}(\mathbf{W}^*)}{\sigma_{min}(\mathbf{W}^*)}$ should not increase with T

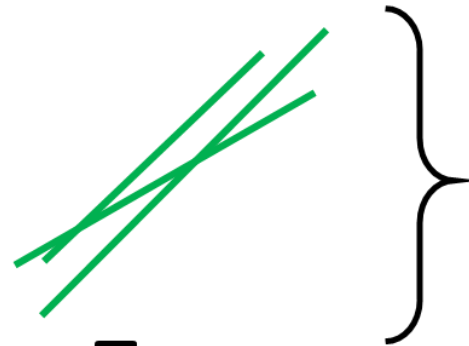
- Assumption 2: Constant classification margin

▶ Norm of the predictors $\|\mathbf{w}_t^*\|_{t \in [1, T]}$ should not increase with T

Source tasks

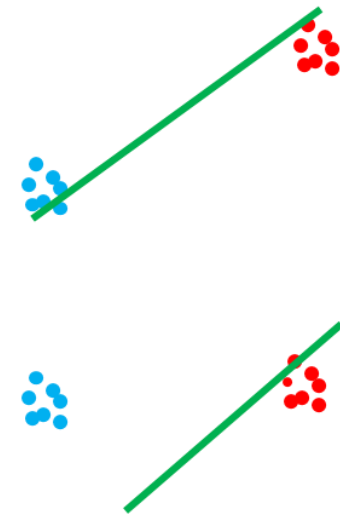


$$W = [w_1, w_2, w_3]$$



$$\frac{\sigma_{max}}{\sigma_{min}} = \kappa(W) \gg 1$$

Target tasks

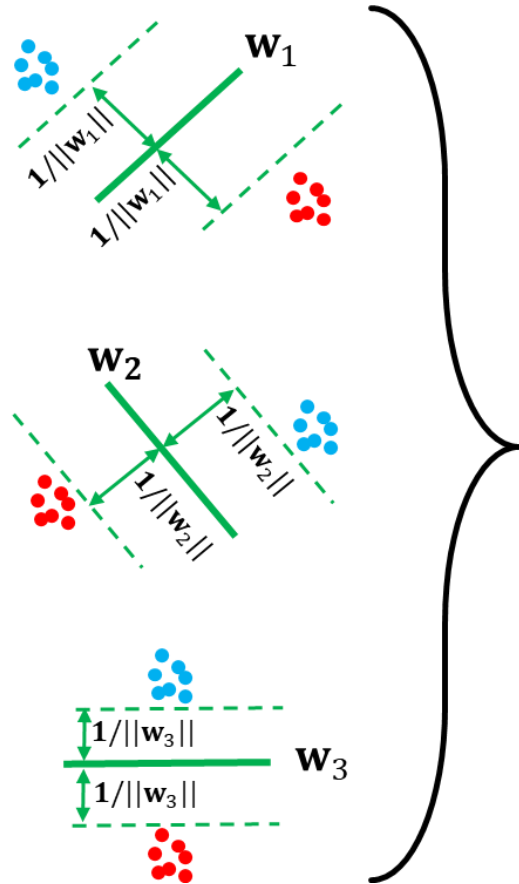


× Linear predictors cover **only part** of the space or **over-specialize** to the tasks

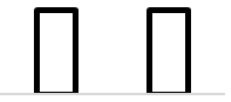
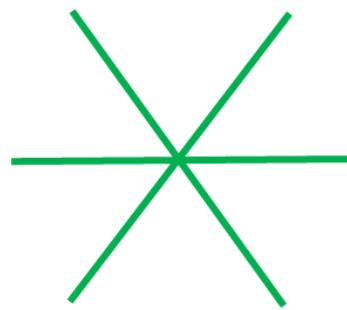
ILLUSTRATION

SATISFIED ASSUMPTIONS

Source tasks



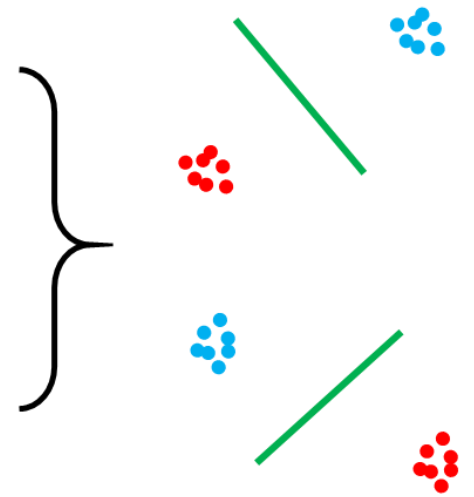
$$W = [w_1, w_2, w_3]$$



$$\sigma_{max} \quad \sigma_{min}$$

$$\kappa(W) \approx 1$$

Target tasks



- ✓ Satisfying assumption 1 makes sure that linear predictors are **complementary**
- ✓ Satisfying assumption 2 avoids **under- or over-specialization** to the tasks

- Few-shot generalization bound

If assumptions are satisfied,

$$\text{ER}(\phi, \mathbf{w}_{T+1}) \leq O\left(\frac{1}{n_1 T} + \frac{1}{n_2}\right)$$

- ✓ All source and target data are useful to decrease the bound of *excess risk*



FROM THEORY TO PRACTICE CONTRIBUTIONS



CAN WE FORCE THE ASSUMPTIONS ?

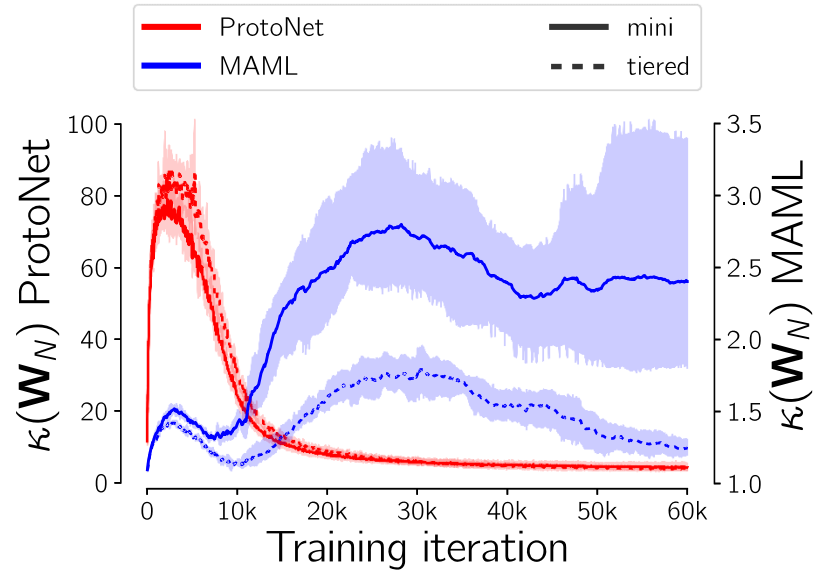
Given \mathbf{W}^* such that $\kappa(\mathbf{W}^*) \gg 1$, can we learn $\hat{\mathbf{W}}$ with $\kappa(\hat{\mathbf{W}}) \approx 1$ while solving the underlying classification problems equally well ?

Given \mathbf{W}^* such that $\kappa(\mathbf{W}^*) \gg 1$, can we learn $\hat{\mathbf{W}}$ with $\kappa(\hat{\mathbf{W}}) \approx 1$ while solving the underlying classification problems equally well ?

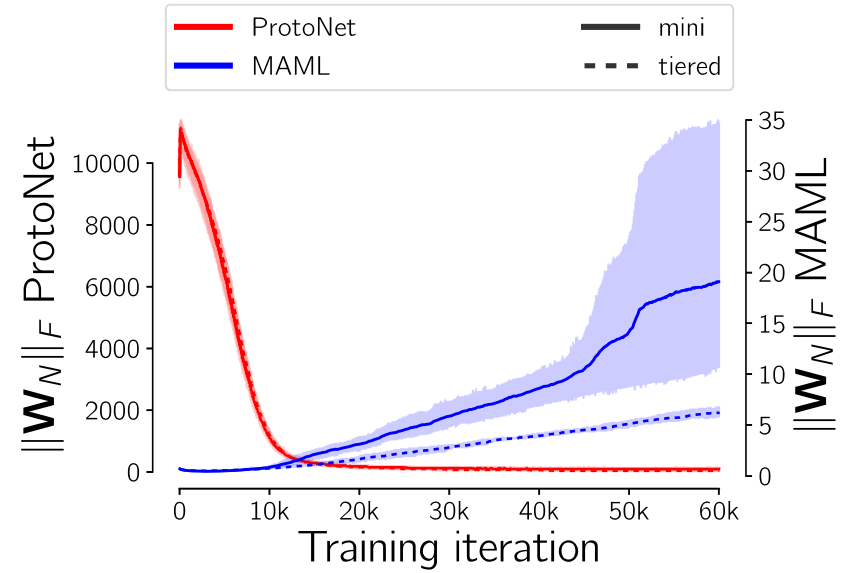
- ✓ Even when \mathbf{W}^* does not satisfy the assumptions, it is **possible to learn** $\hat{\phi}$ to respect them

- **Idea:**
 - ▶ Verify Assumptions 1 and 2 for meta-learning algorithms
- **How ?**
 - ▶ Monitor Condition Number and Norm of the predictors

WHAT HAPPENS IN PRACTICE ?

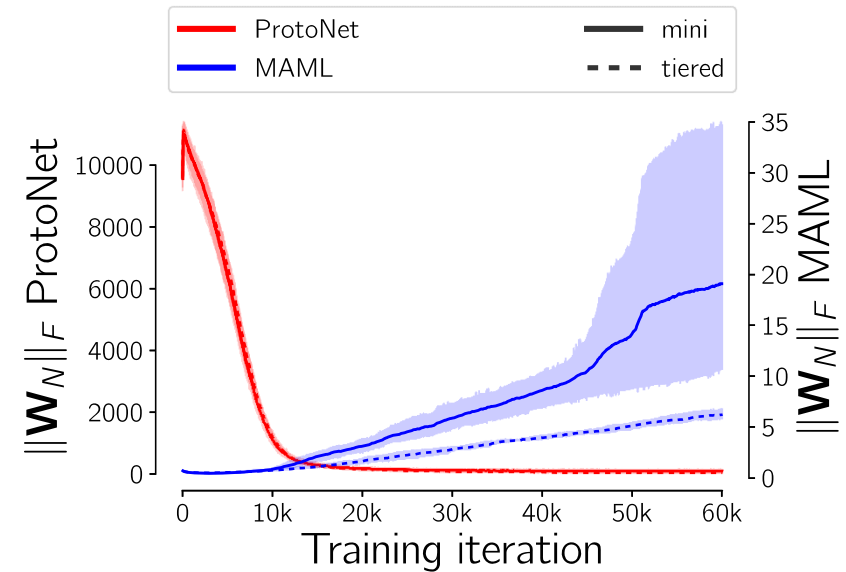
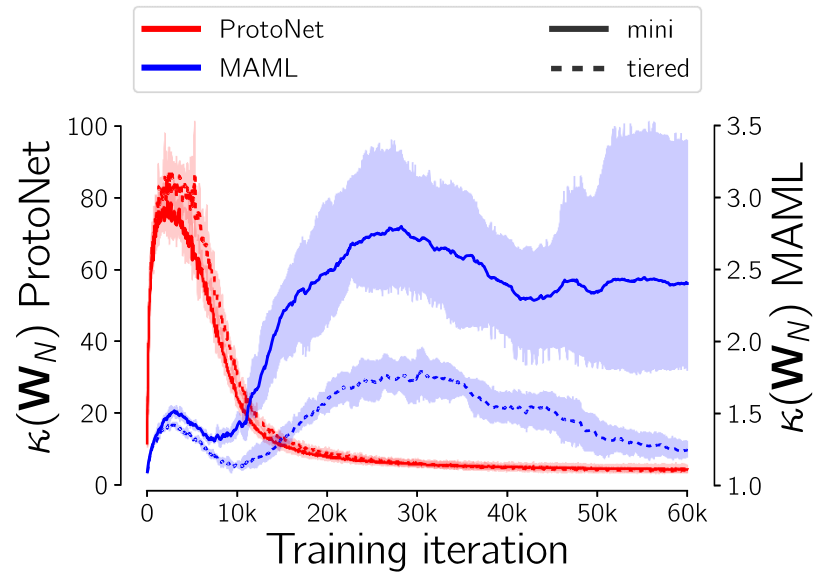


Monitoring the Condition Number



Monitoring the Norm

WHAT HAPPENS IN PRACTICE ?



- ✓ ProtoNet naturally verifies the assumptions
- ✗ MAML does not verify the assumptions



WHY DOES IT HAPPEN ?

- Case of ProtoNet:
 - Theorem (informal)

If all prototypes are normalized,
then all ProtoNet encoders verify Assumption 1.

- ✓ Norm minimization is enough to obtain well-behaved condition number for ProtoNet.

- Case of MAML:
 - Theorem (informal)

At iteration i , if $\sigma_{min} = 0$ for last two tasks,
then $\kappa(\widehat{W}_2^{i+1}) \geq \kappa(\widehat{W}_2^i)$

× The condition number for MAML can increase between iterations.

- Ensuring Assumption 1: Spectral or entropic regularization

- Ensuring Assumption 1: Spectral or entropic regularization

$$\kappa(\mathbf{W}_N) = \frac{\sigma_{\max}(\mathbf{W}_N)}{\sigma_{\min}(\mathbf{W}_N)} \quad \text{or} \quad H_\sigma(\mathbf{W}_N) = \sum_{i=1}^N \text{softmax}(\sigma(\mathbf{W}_N))_i \cdot \log \text{softmax}(\sigma(\mathbf{W}_N))_i$$

- Ensuring Assumption 1: Spectral or entropic regularization

$$\kappa(\mathbf{W}_N) = \frac{\sigma_{max}(\mathbf{W}_N)}{\sigma_{min}(\mathbf{W}_N)} \quad \text{or} \quad H_\sigma(\mathbf{W}_N) = \sum_{i=1}^N \text{softmax}(\sigma(\mathbf{W}_N))_i \cdot \log \text{softmax}(\sigma(\mathbf{W}_N))_i$$

- ✓ Regularizing with $\kappa(\mathbf{W}_N)$ or $H_\sigma(\mathbf{W}_N)$ leads to a **better coverage** of the searched space

- Ensuring Assumption 1: Spectral or entropic regularization

$$\kappa(\mathbf{W}_N) = \frac{\sigma_{max}(\mathbf{W}_N)}{\sigma_{min}(\mathbf{W}_N)} \quad \text{or} \quad H_\sigma(\mathbf{W}_N) = \sum_{i=1}^N \text{softmax}(\sigma(\mathbf{W}_N))_i \cdot \log \text{softmax}(\sigma(\mathbf{W}_N))_i$$

✓ Regularizing with $\kappa(\mathbf{W}_N)$ or $H_\sigma(\mathbf{W}_N)$ leads to a **better coverage** of the searched space

- Ensuring Assumption 2: Norm regularization or normalization for linear predictors

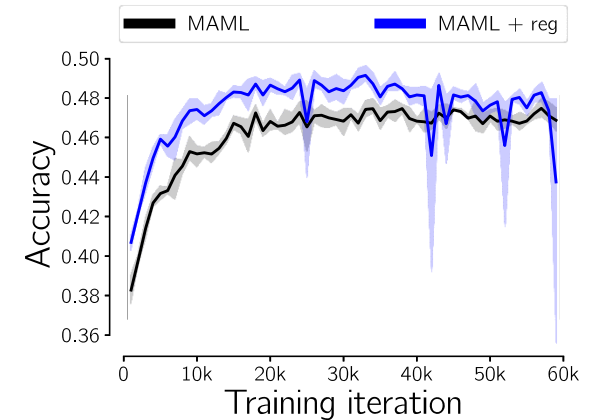
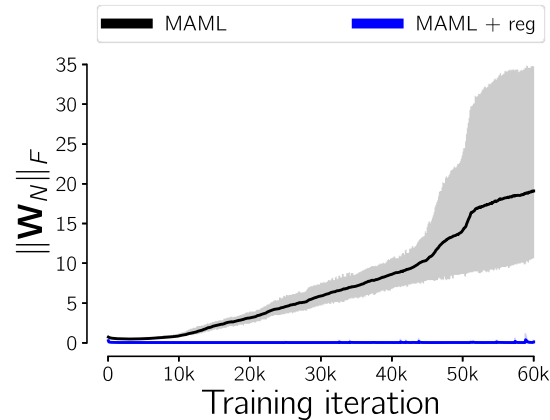
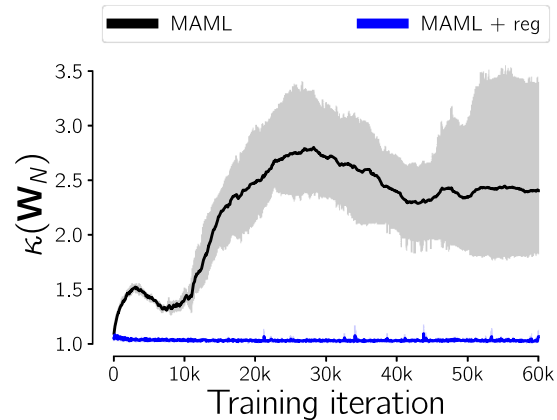
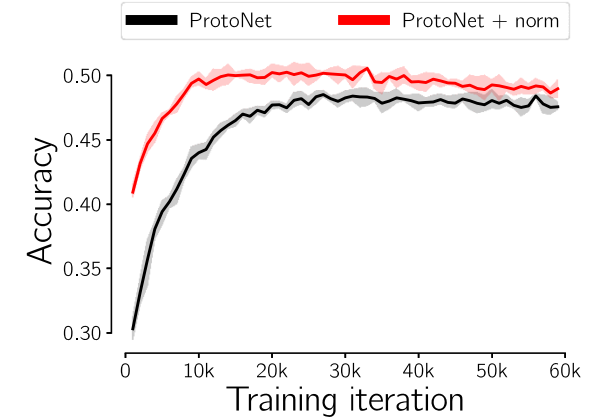
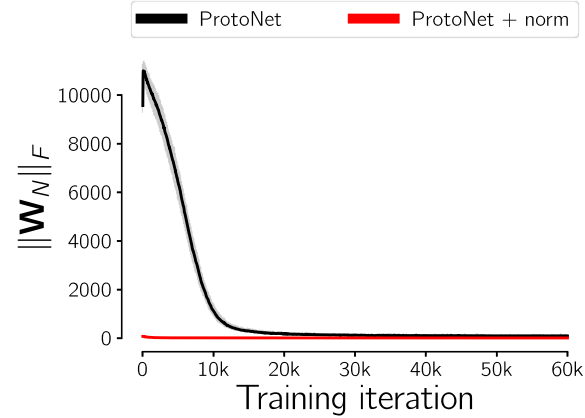
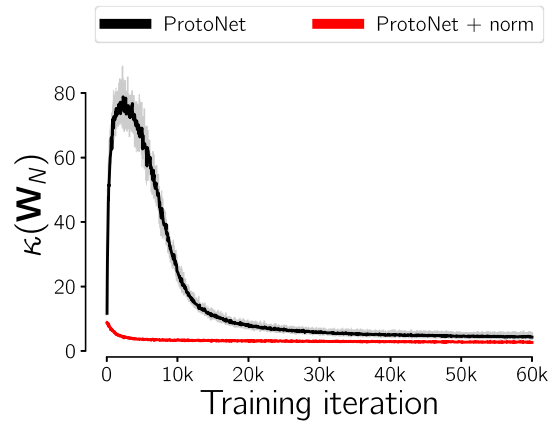
- Ensuring Assumption 1: Spectral or entropic regularization

$$\kappa(\mathbf{W}_N) = \frac{\sigma_{max}(\mathbf{W}_N)}{\sigma_{min}(\mathbf{W}_N)} \quad \text{or} \quad H_\sigma(\mathbf{W}_N) = \sum_{i=1}^N \text{softmax}(\sigma(\mathbf{W}_N))_i \cdot \log \text{softmax}(\sigma(\mathbf{W}_N))_i$$

- ✓ Regularizing with $\kappa(\mathbf{W}_N)$ or $H_\sigma(\mathbf{W}_N)$ leads to a **better coverage** of the searched space
- Ensuring Assumption 2: Norm regularization or normalization for linear predictors
 - ✓ Normalizing predictors ensures **constant margin** that does not change with T

EXPERIMENTAL RESULTS

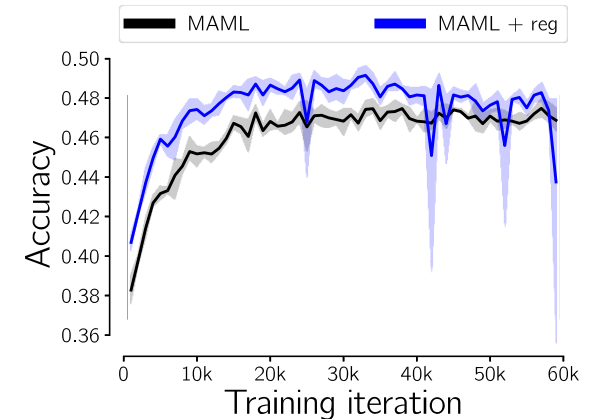
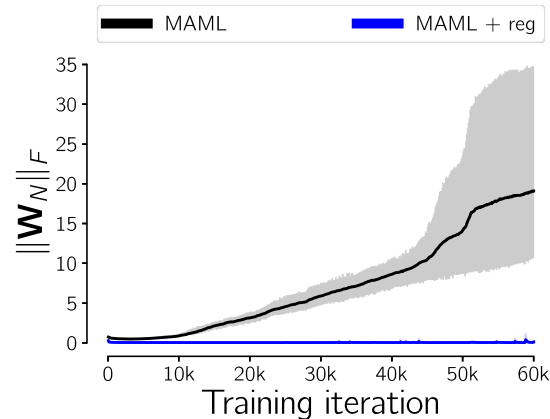
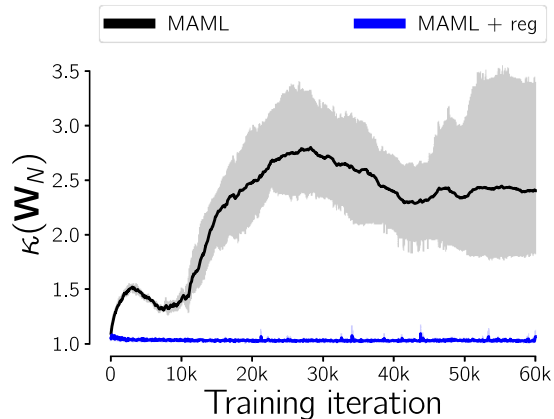
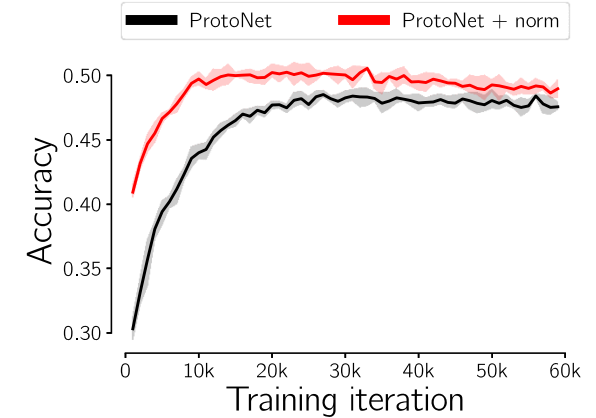
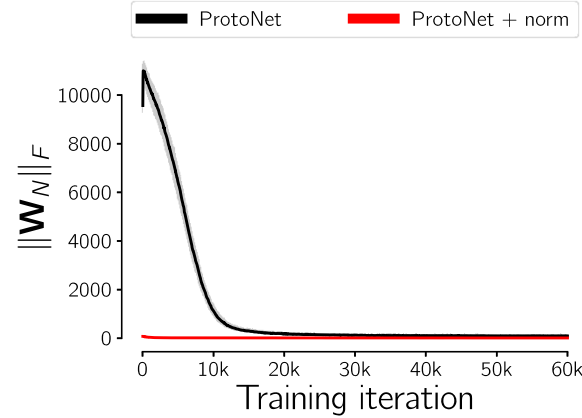
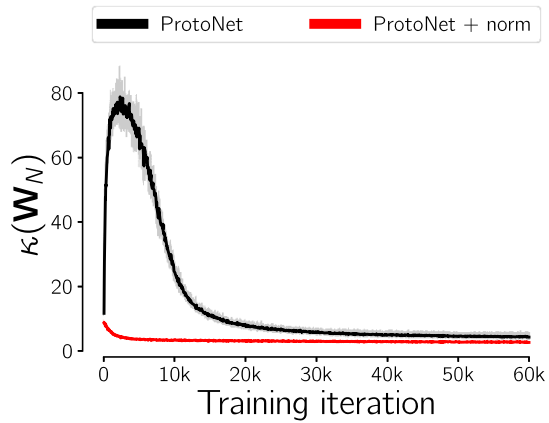
MONITORING THE CONDITION NUMBER AND THE NORM



Experiments on mini-ImageNet 5-way 1-shot

EXPERIMENTAL RESULTS

MONITORING THE CONDITION NUMBER AND THE NORM

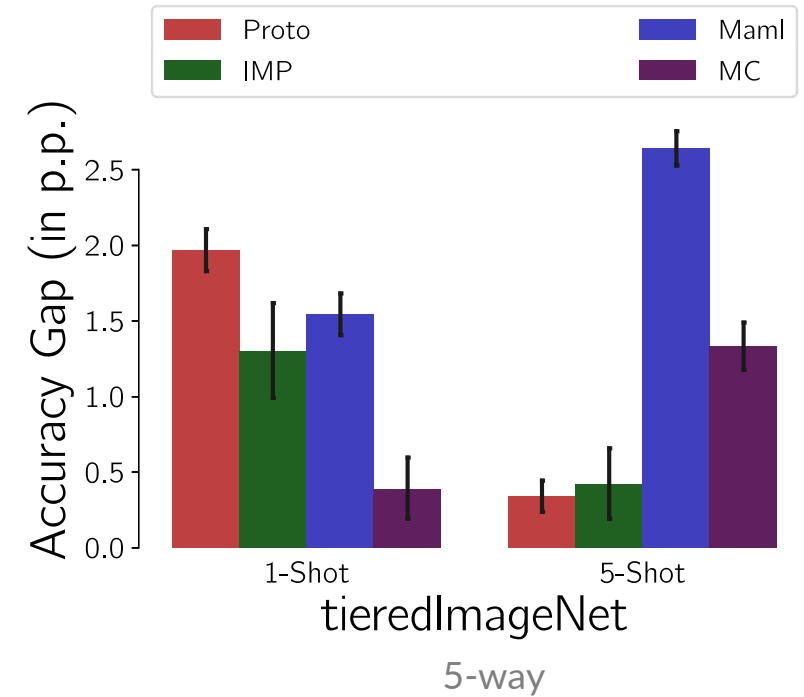
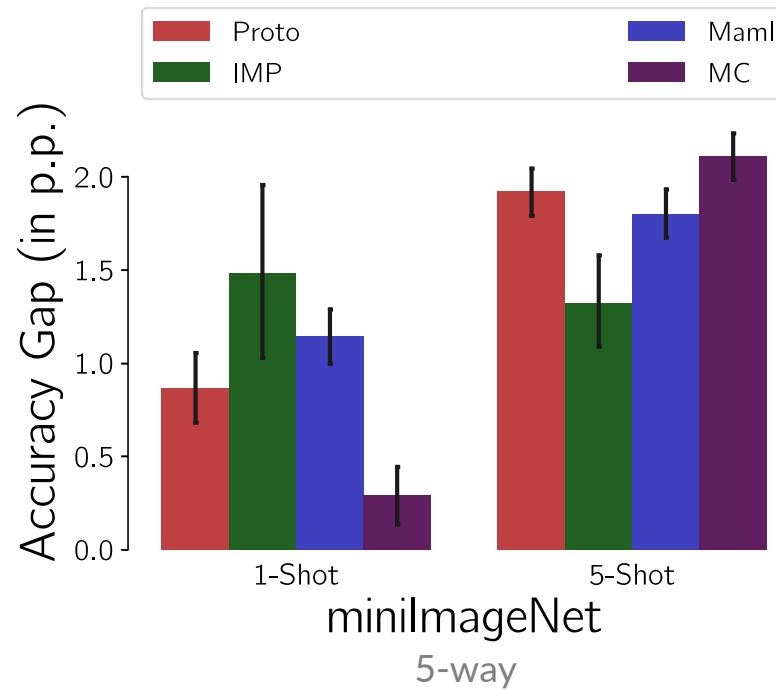
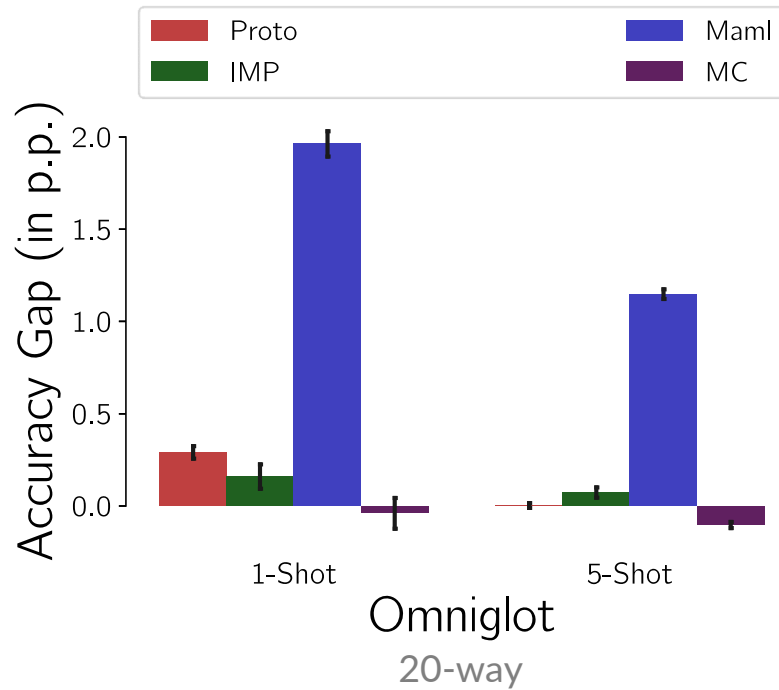


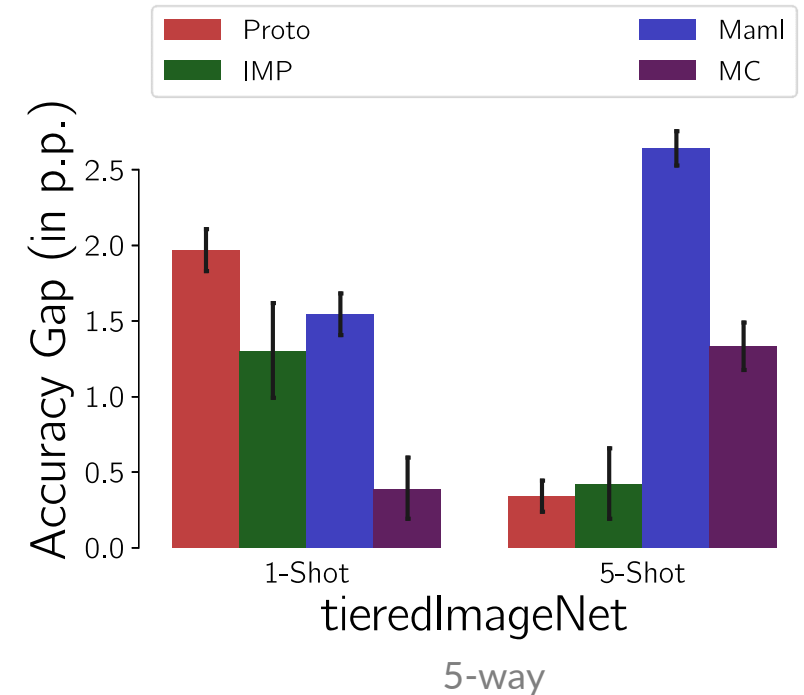
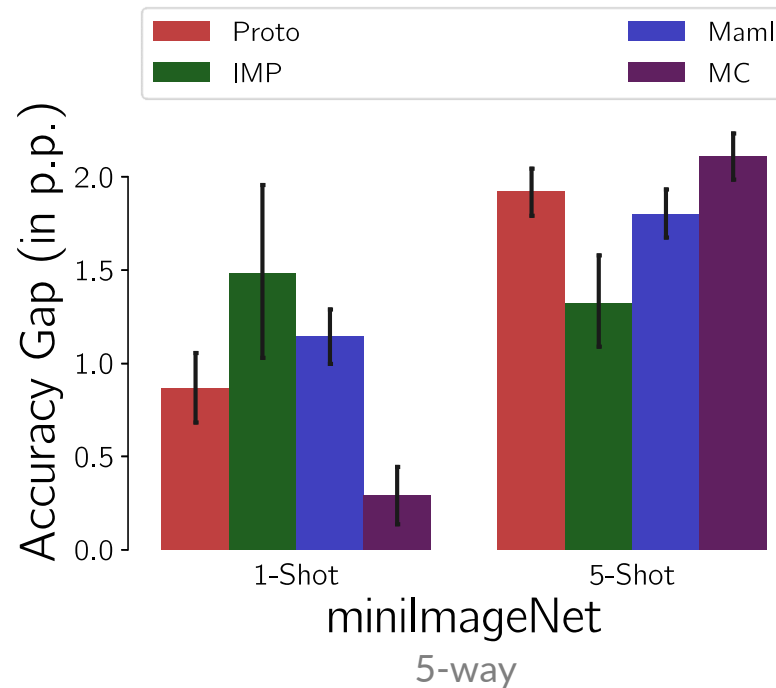
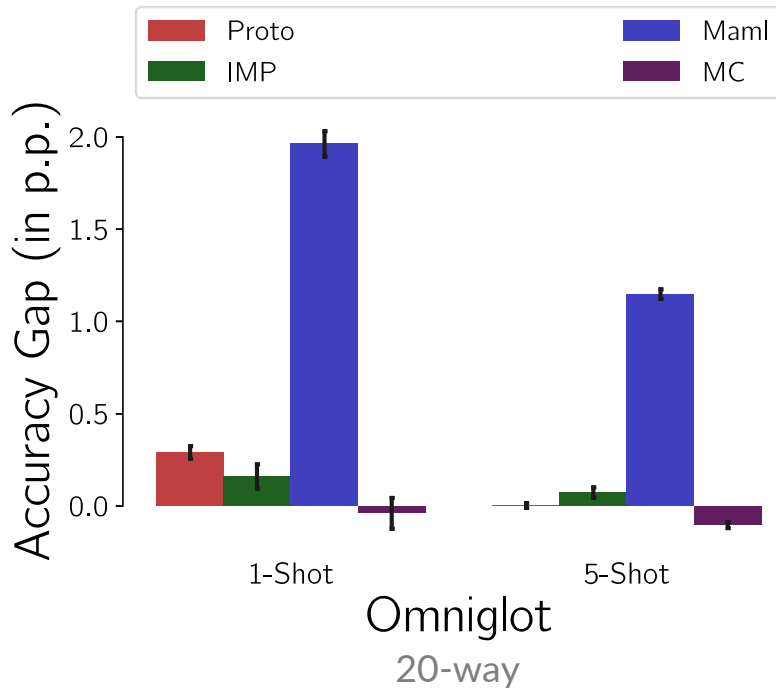
Experiments on mini-ImageNet 5-way 1-shot

✓ Our regularization and normalization have the intended effects.

EXPERIMENTAL RESULTS

ACCURACY GAPS





- ✓ Statistically significant improvement with our regularization and normalization.
- ✓ Better generalization when the assumptions are not verified naturally.

EXPERIMENTAL RESULTS: CROSS-DOMAIN

Source Domain:



ImageNet:
Perspective
Natural Images
Color

Target Domains:
(Disjoint Label Spaces)

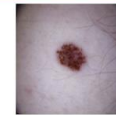
Decreasing Similarity to ImageNet



CropDisease:
Perspective
Natural Images
Color



EuroSAT:
No Perspective
Natural Images
Color



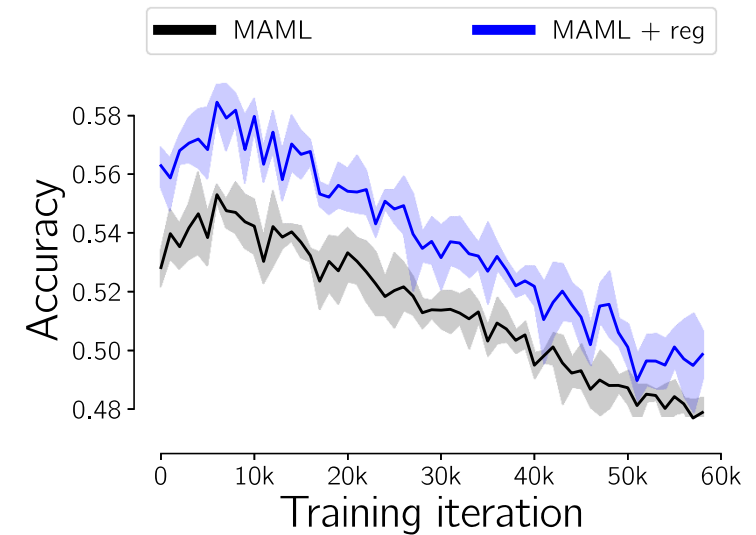
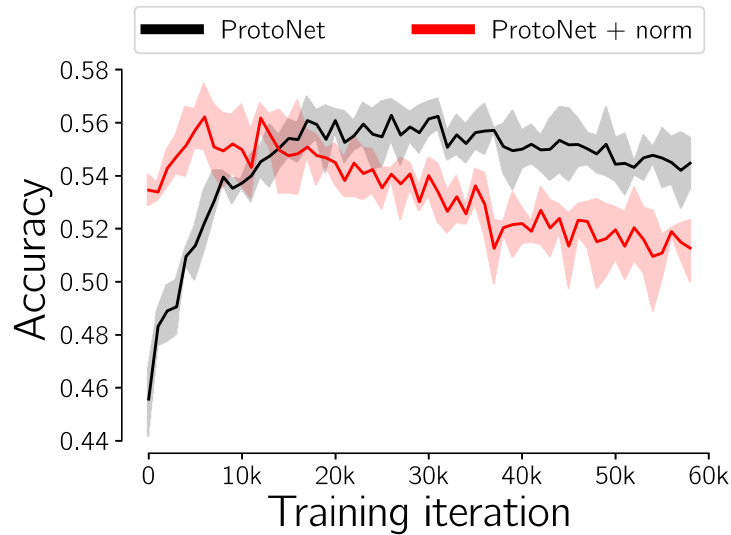
ISIC:
No Perspective
Medical Images
Color



ChestX:
No Perspective
Medical Images
Grayscale

5-way
1-shot

Guo et al. 2020.
A Broader Study of Cross-Domain Few-Shot Learning.
In ECCV 2020



EXPERIMENTAL RESULTS: CROSS-DOMAIN

Source Domain:



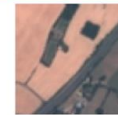
ImageNet:
Perspective
Natural Images
Color

Target Domains:
(Disjoint Label Spaces)

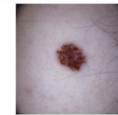
Decreasing Similarity to ImageNet



CropDisease:
Perspective
Natural Images
Color



EuroSAT:
No Perspective
Natural Images
Color



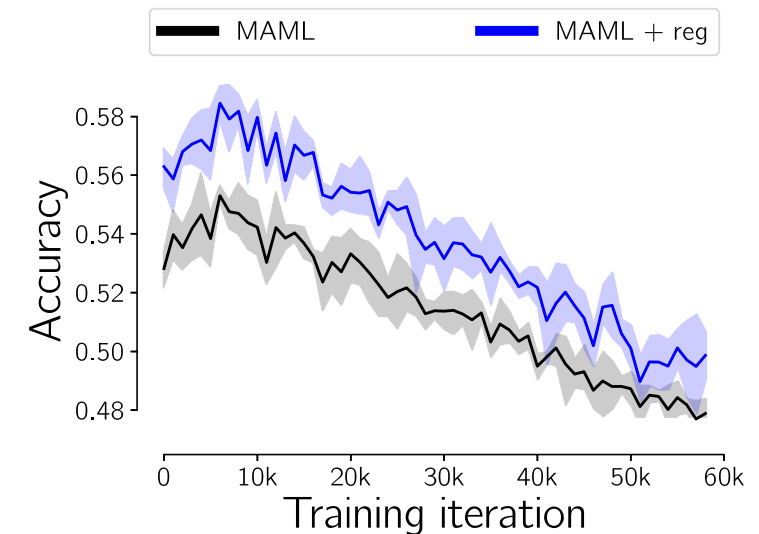
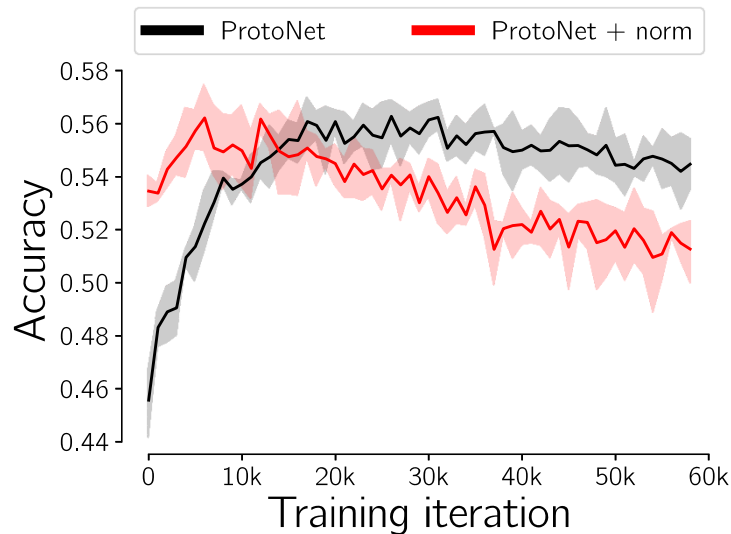
ISIC:
No Perspective
Medical Images
Color



ChestX:
No Perspective
Medical Images
Grayscale

5-way
1-shot

Guo et al. 2020.
A Broader Study of Cross-Domain Few-Shot Learning.
In ECCV 2020



- ✗ Improvement does **not** translate to cross-domain for *metric-based methods*.
- ✓ *Gradient-based methods* keep their accuracy gains.



TAKE HOME MESSAGE



- Improving Few-Shot Learning Through Multi-Task Representation Learning Theory
 - ✓ Connection between Meta-Learning and Multi-Task Representation Learning Theory
 - ✓ Explaining why some meta-learning methods **naturally fulfill** theoretical assumptions of the best learning bounds.
 - ✓ **Practical ways** to enforce the assumptions which leads to **significant** performance improvements.

More details in
arXiv paper:



Contact:

 quentin.bouniot@cea.fr

 @QBouniot

 <https://qbouniot.github.io>

Thank you for listening !

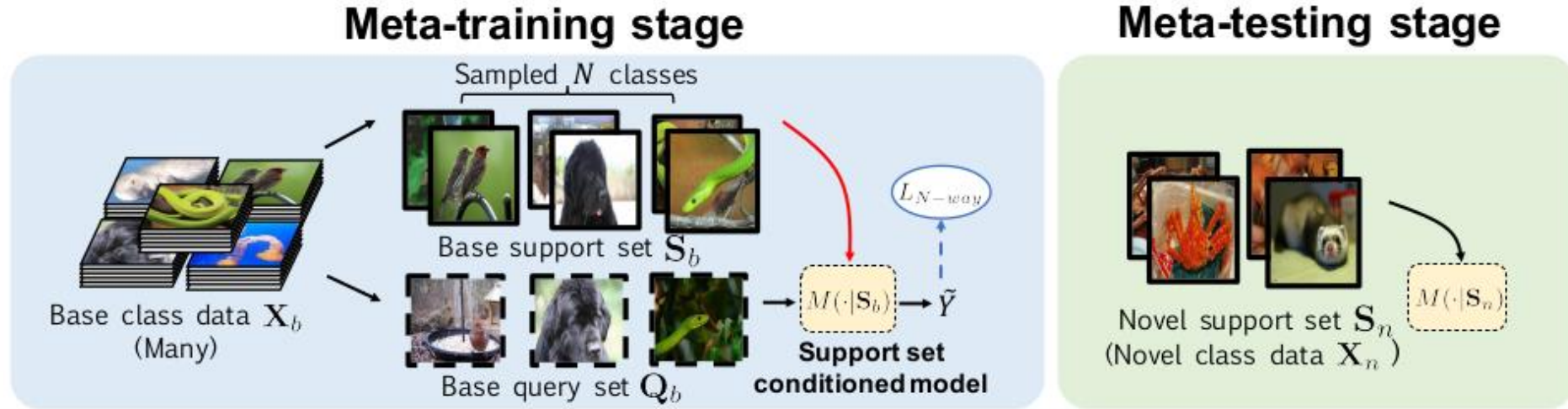


Commissariat à l'énergie atomique et aux énergies alternatives
Institut List | CEA SACLAY NANO-INNOV | BAT. 861 - PC142
91191 Gif-sur-Yvette Cedex - FRANCE
www-list.cea.fr

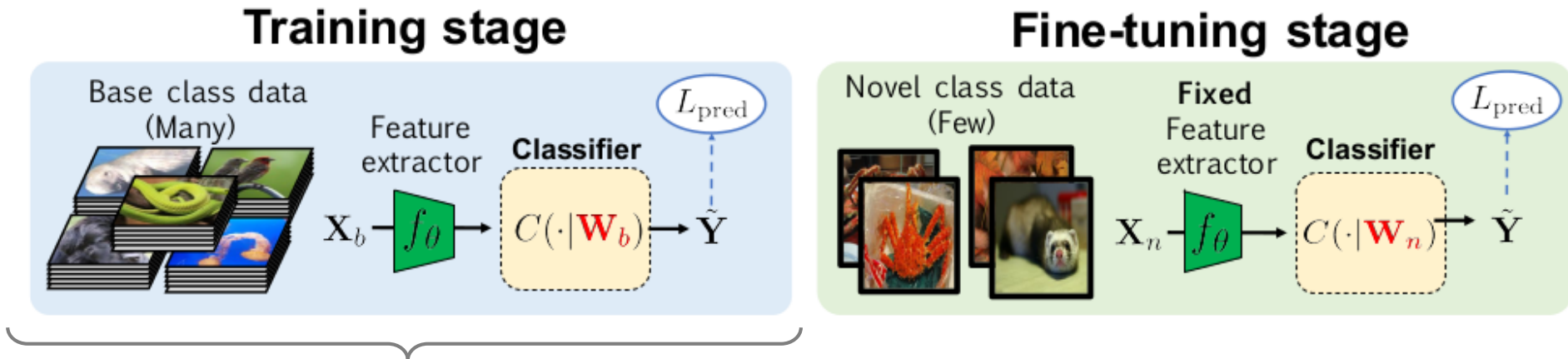
Établissement public à caractère industriel et commercial | RCS Paris B 775 685 019

APPENDIX: EPISODIC TRAINING VS REGULAR TRAINING

Episodic Training



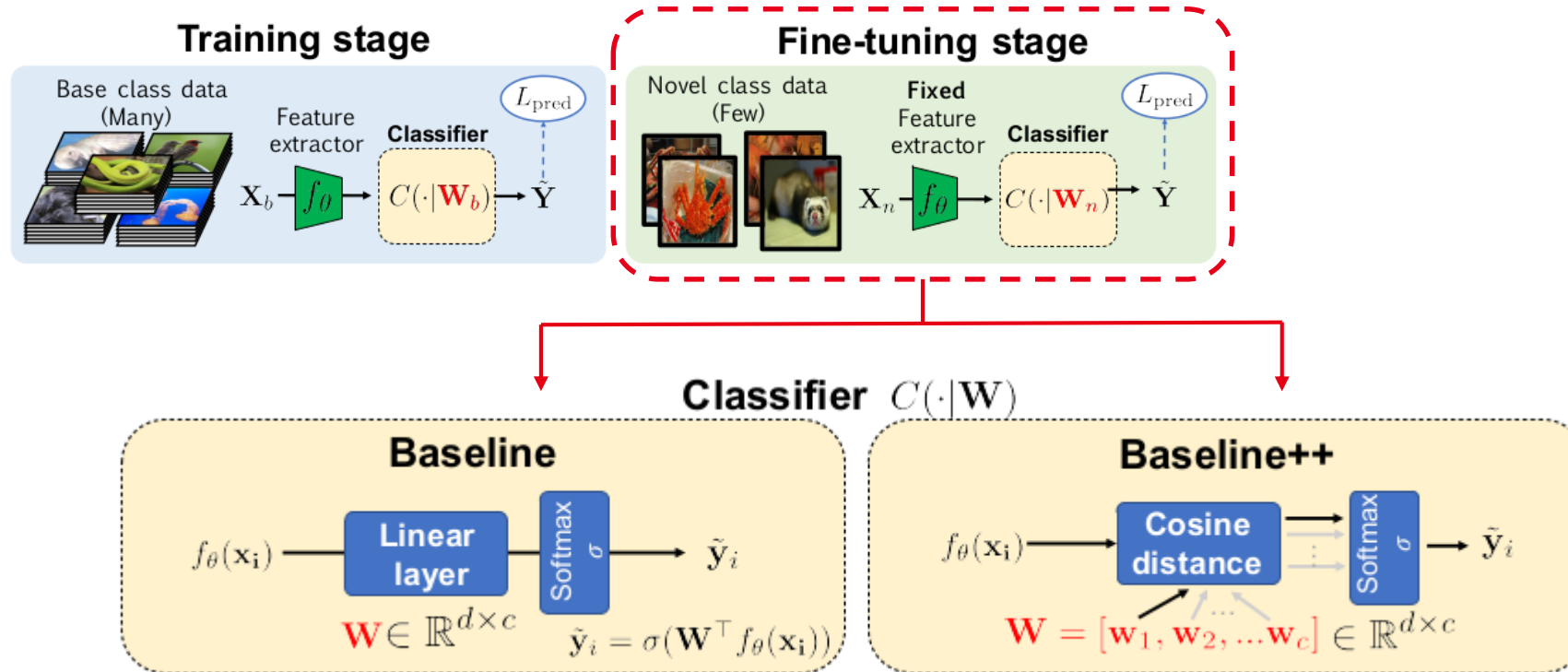
Regular Training



Learning a single episode

Appendix: Fine-tuning methods

Regular
Training



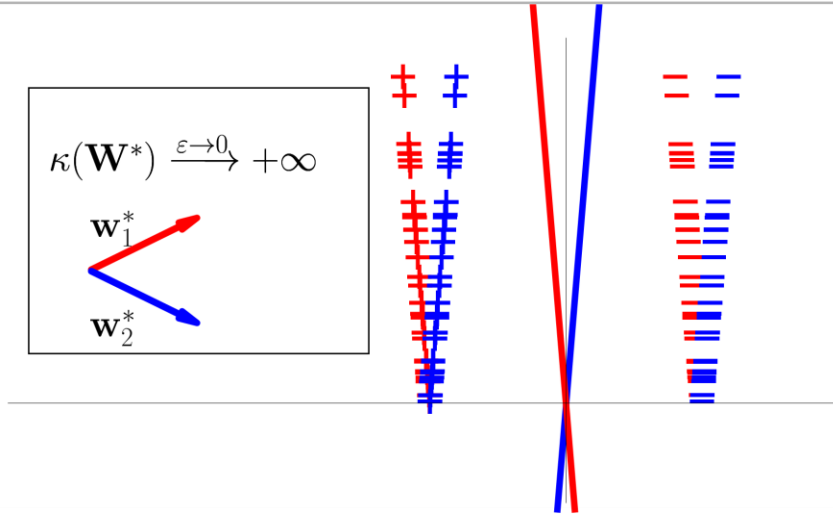
Adapted from [Chen19]

- **Baseline** uses a dot product in the classification layer followed by a softmax
- **Cosine classifier (or Baseline++)** uses a cosine similarity followed by a softmax

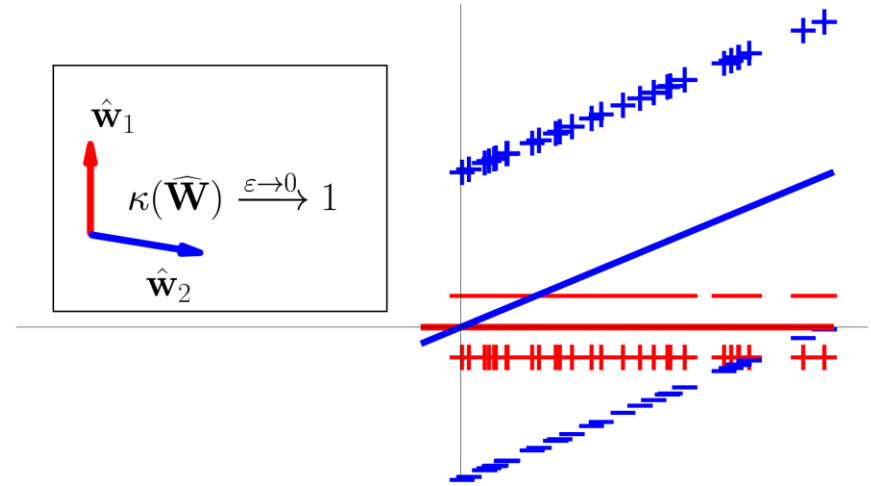
APPENDIX: CAN WE FORCE THE ASSUMPTIONS ?

Given \mathbf{W}^* such that $\kappa(\mathbf{W}^*) \gg 1$, can we learn $\hat{\mathbf{W}}$ with $\kappa(\hat{\mathbf{W}}) \approx 1$ while solving the underlying classification problems equally well ?

+ - Source task 1 in Φ^* space + - Source task 2 in Φ^* space

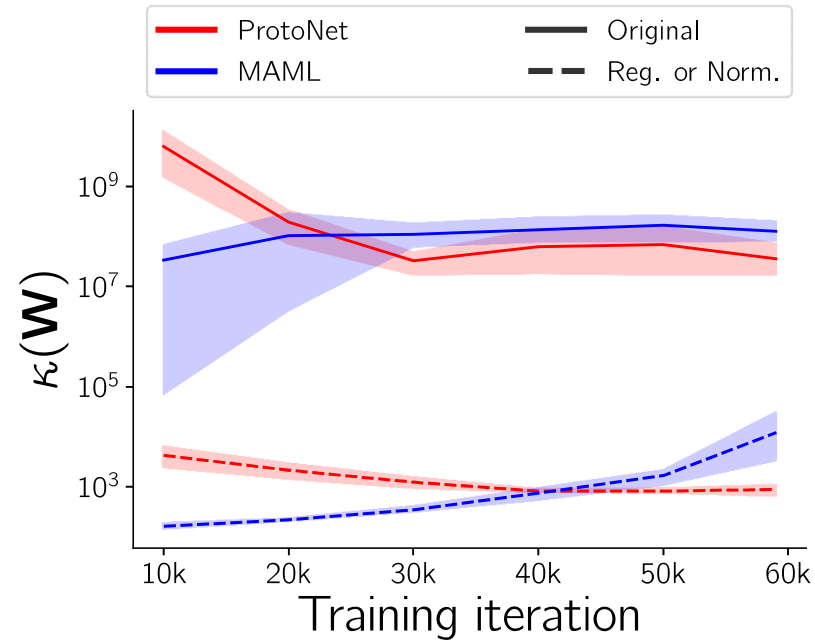


+ - Source task 1 in $\hat{\Phi}$ space + - Source task 2 in $\hat{\Phi}$ space



✓ Even when \mathbf{W}^* does not satisfy the assumptions, it is possible to learn $\hat{\phi}$ to respect them

APPENDIX: CONDITION NUMBER OF ALL PREDICTORS



- ▶ $\kappa(\mathbf{W}_N)$ shows dynamics during training, but values are not comparable
- ▶ $\kappa(\mathbf{W})$ is intractable to compute during training.