# CAp 2021
# Conférence sur l'Apprentissage Automatique

**VERS UNE MEILLEURE COMPREHENSION DES MÉTHODES DE MÉTA-APPRENTISSAGE À TRAVERS LA THÉORIE DE L'APPRENTISSAGE DE REPRESENTATIONS MULTI-TÂCHES**
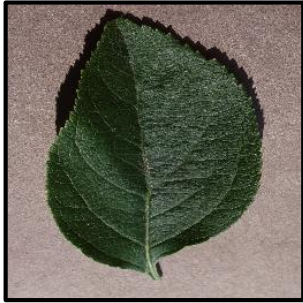
**Quentin BOUNIOT**
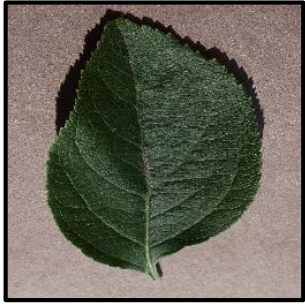
Apple

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

Quentin Bouniot, Ievgen Redko, Romaric Audigier, Angélique Loesch | CAp2021 | 14/06/2021

# META-LEARNING 101

- **What is Meta-learning ?**

  ► A meta-learner trained on multiple tasks.

  ► For each task, the meta-learner trains a learner.

  ► The meta-learner is evaluated on new unseen tasks.

# META-LEARNING 101

- **What is Meta-learning ?**

  ► A meta-learner trained on multiple tasks.

  ► For each task, the meta-learner trains a learner.

  ► The meta-learner is evaluated on new unseen tasks.

► **Meta-Learning can be used for a lot of problems (classification, regression, RL, ...)**

- **What is Meta-learning ?**

  ▶ A meta-learner trained on multiple tasks.

  ▶ For each task, the meta-learner trains a learner.

  ▶ The meta-learner is evaluated on new unseen tasks.

▶ **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

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

Support Set

Query Set

Set of classes

Model

Learning predictors

Evaluating

Support Set

Learning predictors

Model

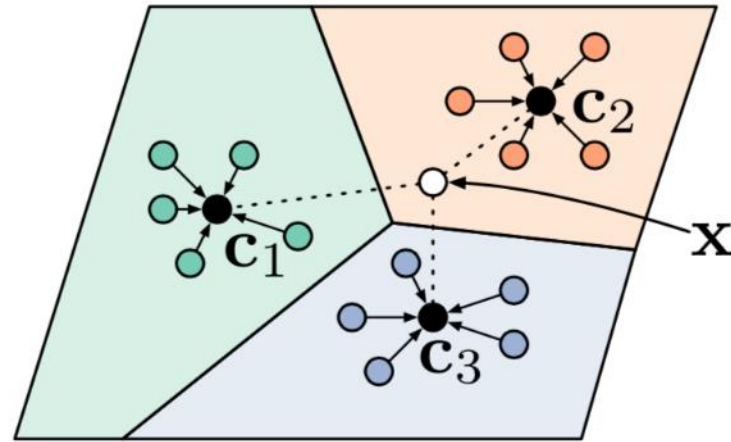Query Set

Evaluating
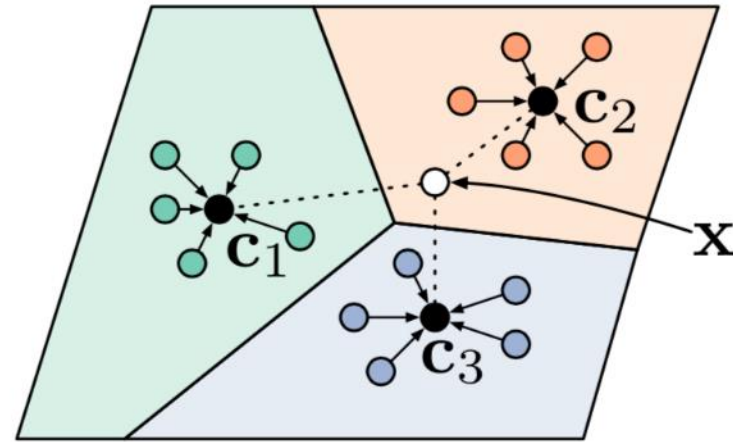
Set of classes

- **Disjoint** sets of classes between **meta-training** and **meta-testing** classes
- Construction of **episodes** from dataset
- **Non-overlapping** class labels between episodes

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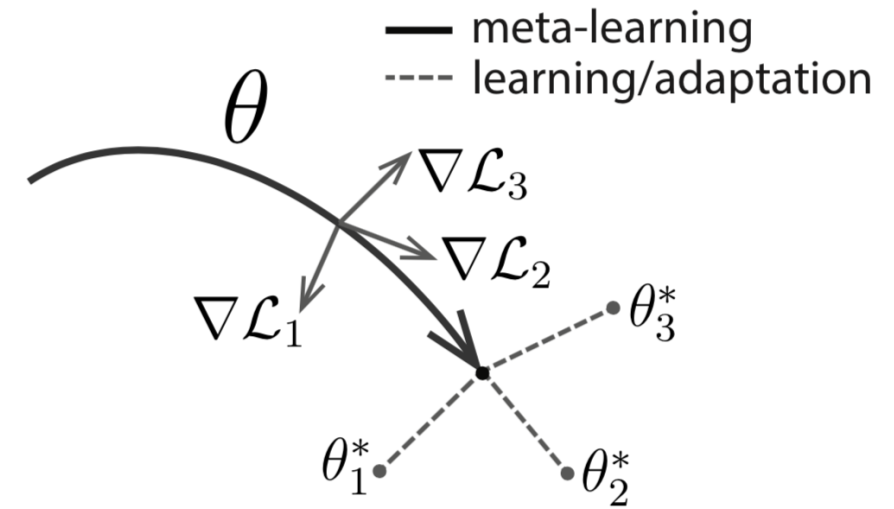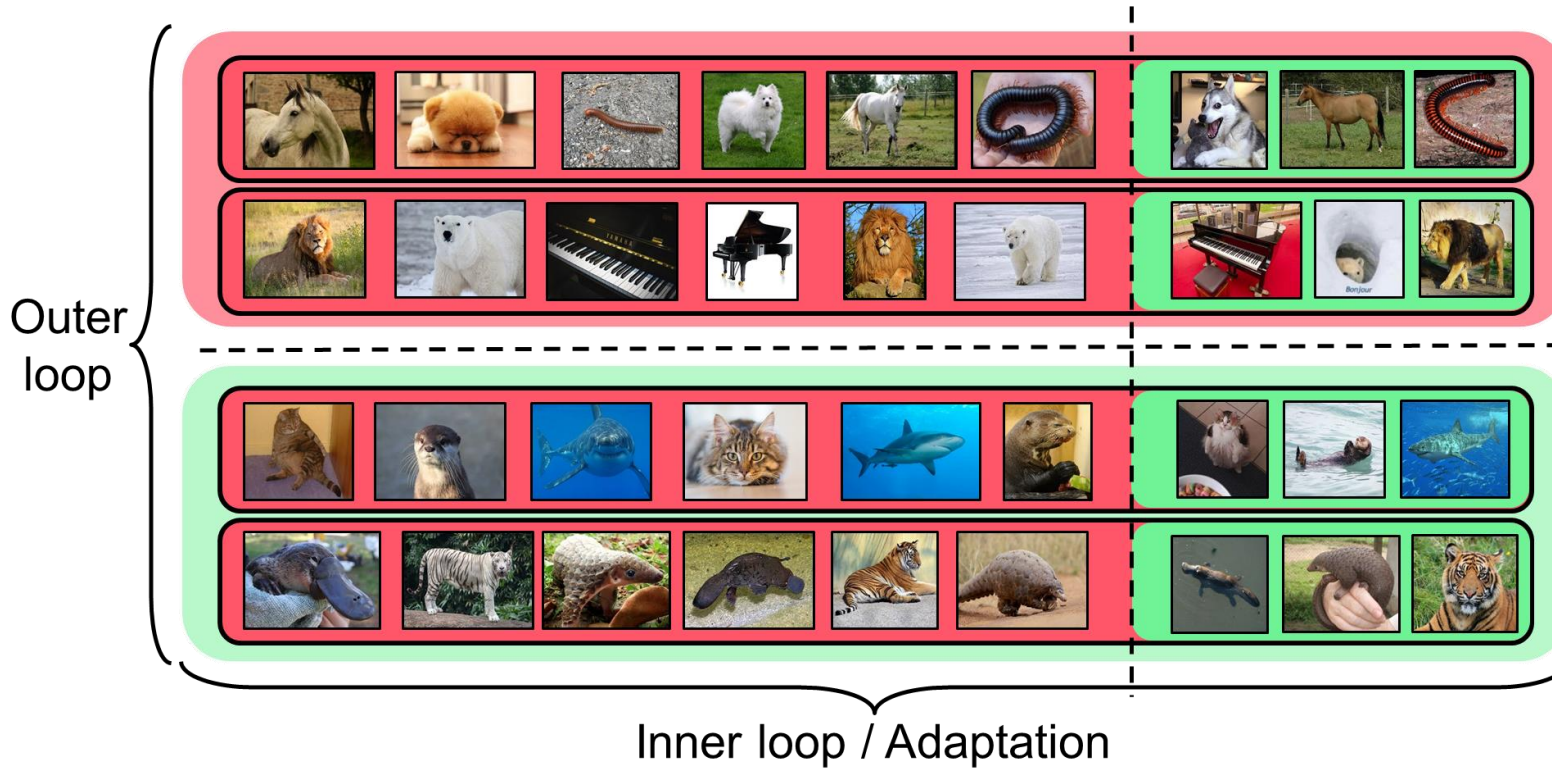
- **Embedding function** to encode query and support samples.
- Support samples fused into **prototypes $c_i$** for each class
- Probability distribution using **inverse of distances** to prototypes.
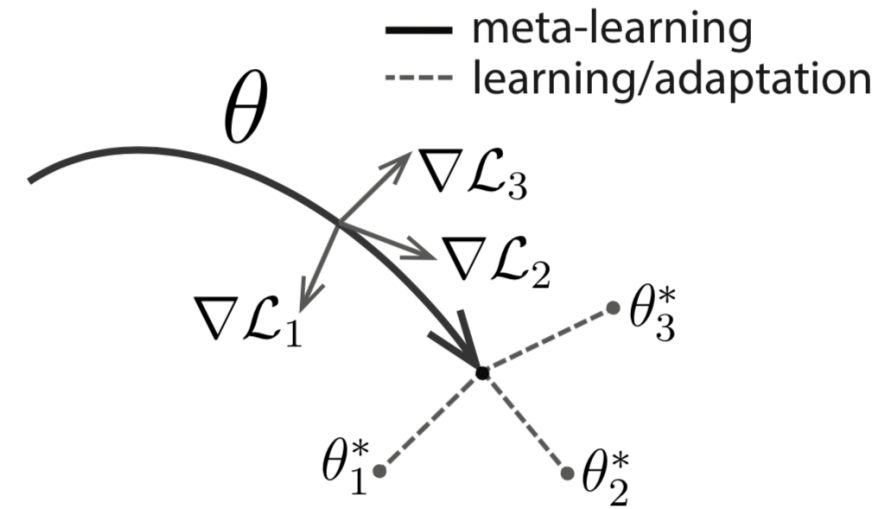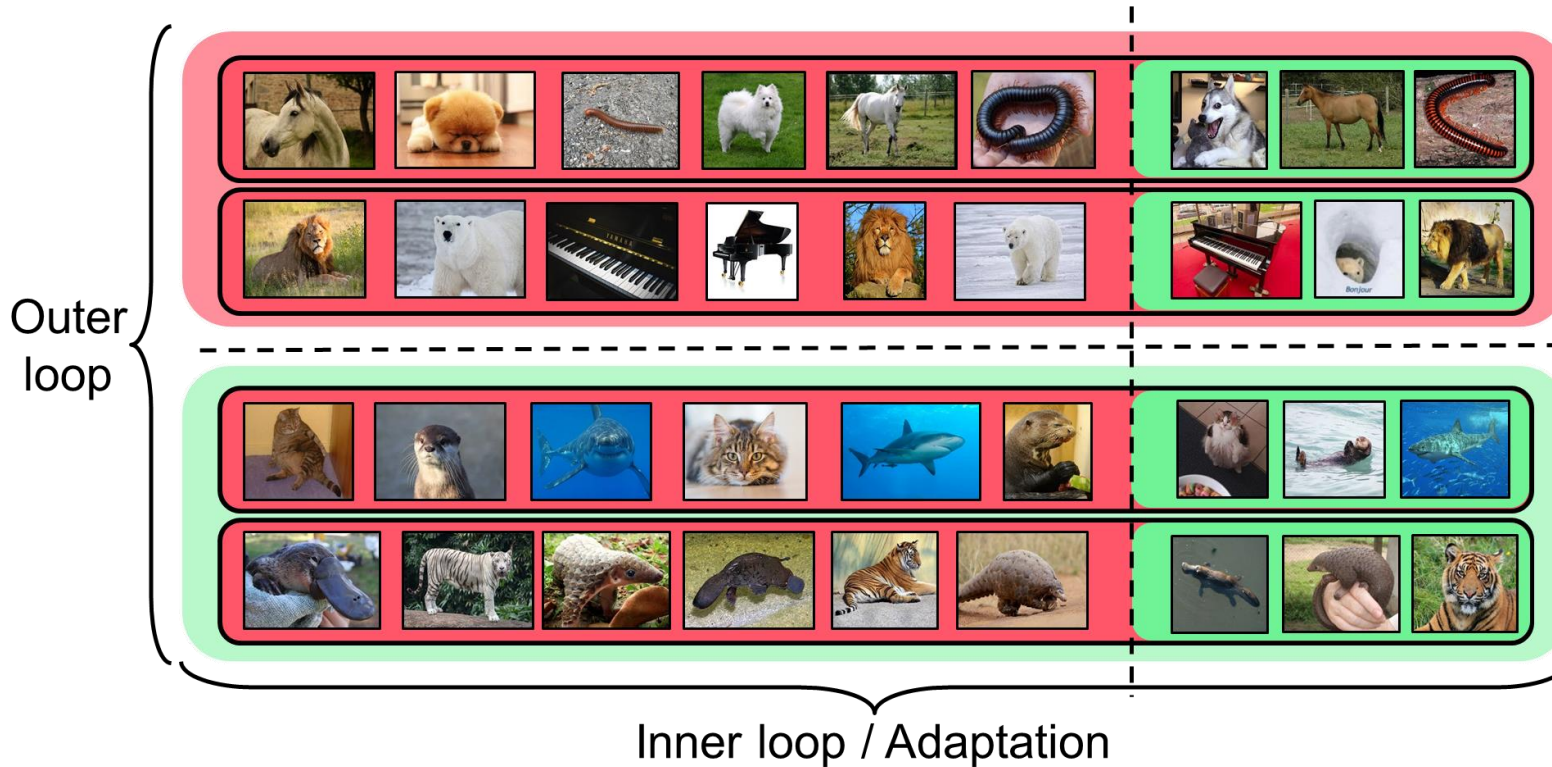- **Contrastive loss** according to distance function.

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.

Outer loop

Inner loop / Adaptation

— meta-learning
--- learning/adaptation

$\theta$

$\nabla\mathcal{L}_3$

$\nabla\mathcal{L}_2$

$\nabla\mathcal{L}_1$

$\theta_3^*$

$\theta_1^*$

$\theta_2^*$

Finn C. et al. (2017), *Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks*. In ICML 2017
Park E. et Oliva J.B. (2019). *Meta-curvature*. In NeurIPS 2019

- **Inner Loop:**
  - ► Performs a few gradient updates over the *k* labelled examples (the support set) of **current episode/task**.
- **Outer Loop:**
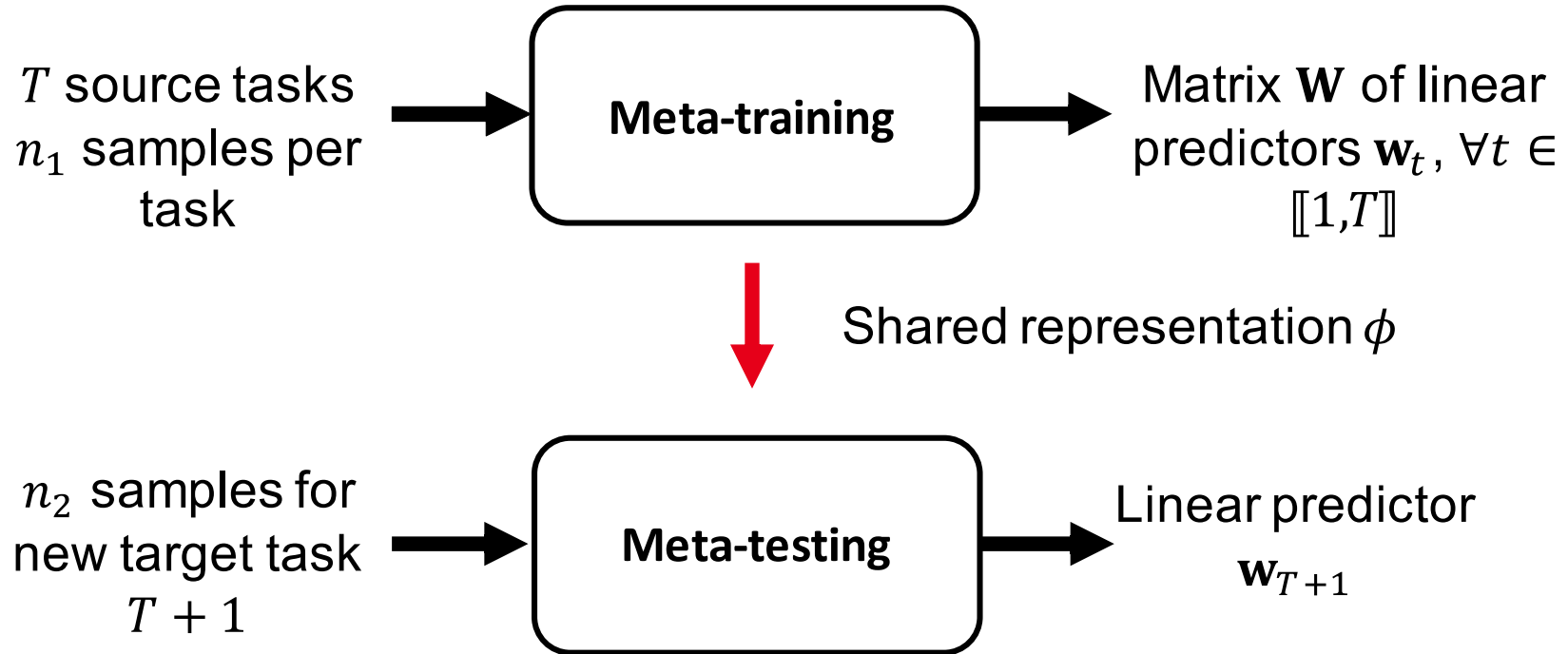  - ► Updates the **initialization** of the parameters (often called the *meta-initialization*).
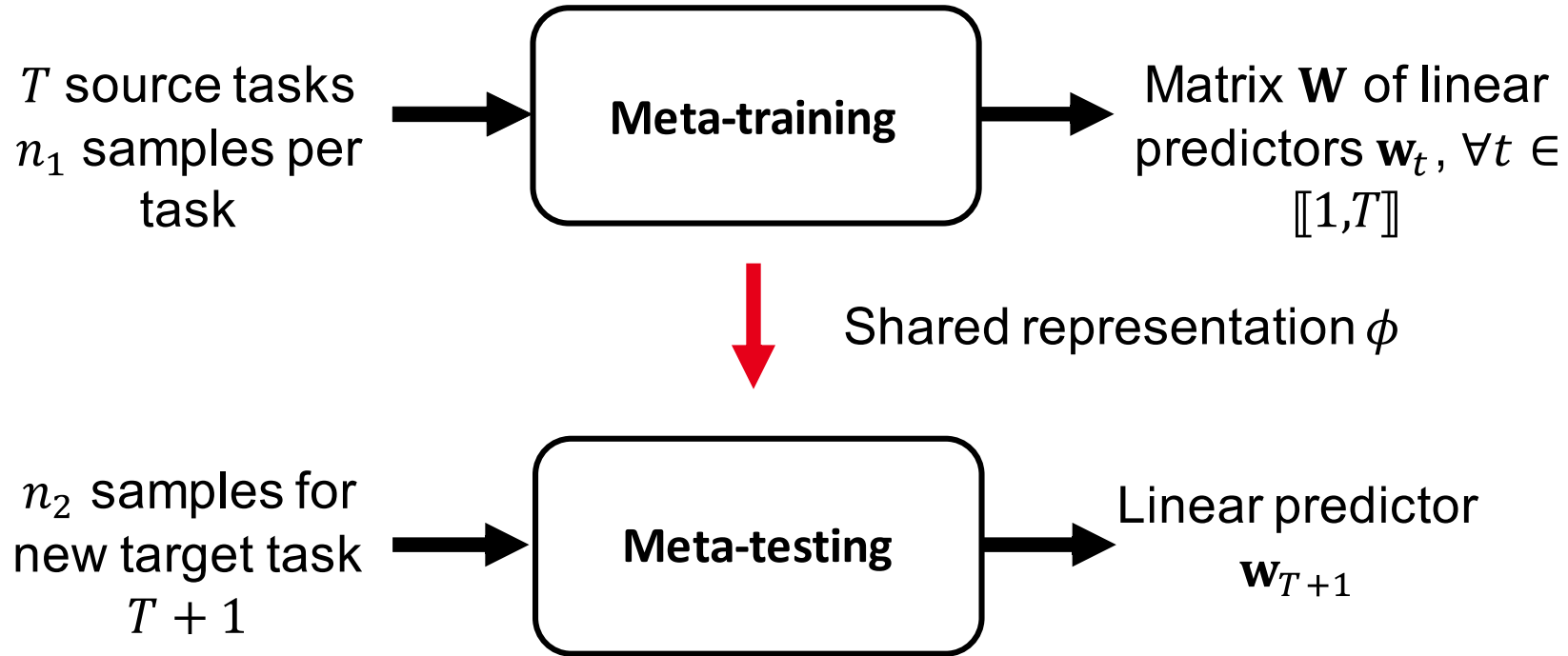
Finn C. et al. (2017), *Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks.* In ICML 2017
Park E. et Oliva J.B. (2019). *Meta-curvature.* In NeurIPS 2019

# MULTI-TASK REPRESENTATION LEARNING THEORY

Quentin Bouniot, Ievgen Redko, Romaric Audigier, Angélique Loesch | CAp2021 | 14/06/2021

$T$ source tasks $n_1$ samples per task $\rightarrow$ **Meta-training** $\rightarrow$ Matrix $\mathbf{W}$ of linear predictors $\mathbf{w}_t$, $\forall t \in [\![1,T]\!]$

$T$ source tasks $n_1$ samples per task

**Meta-training**

Matrix $\mathbf{W}$ of linear predictors $\mathbf{w}_t$, $\forall t \in [\![1,T]\!]$

Shared representation $\phi$

$n_2$ samples for new target task $T+1$

**Meta-testing**

Linear predictor $\mathbf{w}_{T+1}$

$T$ source tasks $n_1$ samples per task → **Meta-training** → Matrix $\mathbf{W}$ of linear predictors $\mathbf{w}_t$, $\forall t \in [\![1,T]\!]$

Shared representation $\phi$

$n_2$ samples for new target task $T+1$ → **Meta-testing** → Linear predictor $\mathbf{w}_{T+1}$

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

▶ True risk $\mathcal{L}$ ▶ Optimal weights $\phi^*$ ▶ $\mathrm{w}_{T+1}^*$ ideal target linear predictor

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 $\mathbf{W}^* = [\mathbf{w}_1^*, \ldots, \mathbf{w}_T^*]$ cover all the directions evenly

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

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 $\mathbf{W}^* = [\mathbf{w}_1^*, \ldots, \mathbf{w}_T^*]$ cover all the directions evenly

  > ► Condition Number $\kappa(\mathbf{W}^*) = \dfrac{\sigma_{max}(\mathbf{W}^*)}{\sigma_{min}(\mathbf{W}^*)}$ 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

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 2:        **Constant classification margin**

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

  > If satisfied, $\mathrm{ER}(\phi, \mathbf{w}_{T+1}) \leq O(\frac{1}{n_1 T} + \frac{1}{n_2})$

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

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.

Source tasks

$\mathbf{w}_1$

$1/\|\mathbf{w}_1\|$
$1/\|\mathbf{w}_1\|$

$\mathbf{w}_2$

$1/\|\mathbf{w}_2\|$
$1/\|\mathbf{w}_2\|$

$1/\|\mathbf{w}_3\|$
$1/\|\mathbf{w}_3\|$
$\mathbf{w}_3$

$\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3]$

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

$\kappa(\mathbf{W}) \gg 1$

Target tasks

**Source tasks**

$W = [w_1, w_2, w_3]$

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

$\kappa(W) \gg 1$

**Target tasks**

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

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

# FROM THEORY TO PRACTICE
## CONTRIBUTIONS

Quentin Bouniot, Ievgen Redko, Romaric Audigier, Angélique Loesch | CAp2021 | 14/06/2021

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

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

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

Monitoring the Condition Number

Monitoring the Norm

Monitoring the Condition Number

Monitoring the Norm

▶ $\mathbf{W}_N$ restriction to $N$ last predictors

✓ ProtoNet **naturally verifies** the assumptions

✗ MAML **does not** verify the assumptions

# WHY DOES IT HAPPEN ?

- **Theorem** (Normalized ProtoNet):

$$\text{if } \forall i \; \|prototype_i\| = 1, \text{then } \exists \phi \in \arg\min loss \text{ such that } \kappa(\mathbf{W}^*) = 1$$

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

- **Theorem** (Normalized ProtoNet):

$$\text{if } \forall i \, \|prototype_i\| = 1, \text{then } \exists \phi \in \arg\min loss \text{ such that } \kappa(\mathbf{W}^*) = 1$$

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

- **Proposition** (Condition Number for MAML):

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

    ✕  The condition number for MAML can **increase** between iterations.

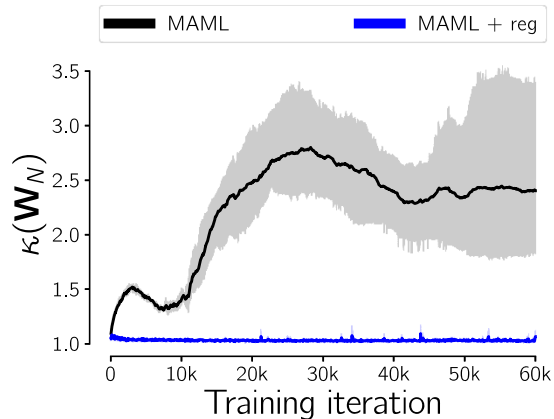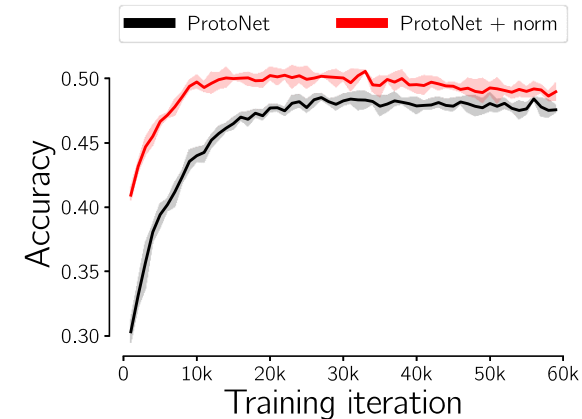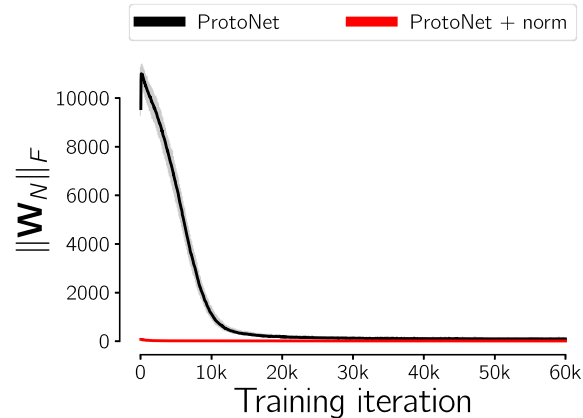- **Ensuring Assumption 1:** **Spectral** or **entropic** regularization
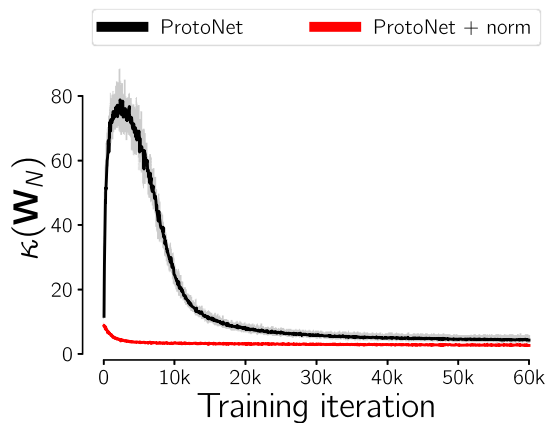
- **Ensuring Assumption 1:**     **Spectral** or **entropic** regularization

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

- **Ensuring Assumption 1:**     **Spectral** or **entropic** regularization

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

     ✓  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

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

✓ 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

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

✓ 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 $T$
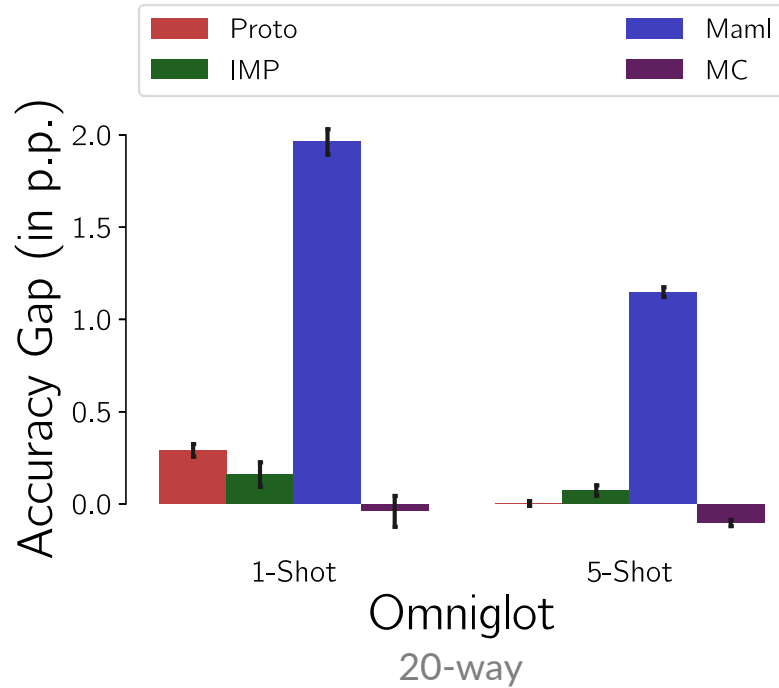
Experiments on mini-ImageNet 5-way 1-shot

Experiments on mini-ImageNet 5-way 1-shot

✓ Our **regularization** and **normalization** have the intended effects.

- ✓ **Statistically significant** improvement with our regularization and normalization.
- ✓ Enforcing the assumptions leads to **better generalization** when not verified naturally.

# EXPERIMENTAL RESULTS:
## CROSS-DOMAIN

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



5-way
1-shot

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



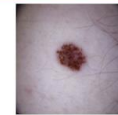Source Domain:

Target Domains:
(Disjoint Label Spaces)

Decreasing Similarity to ImageNet

**ImageNet:**
Perspective
Natural Images
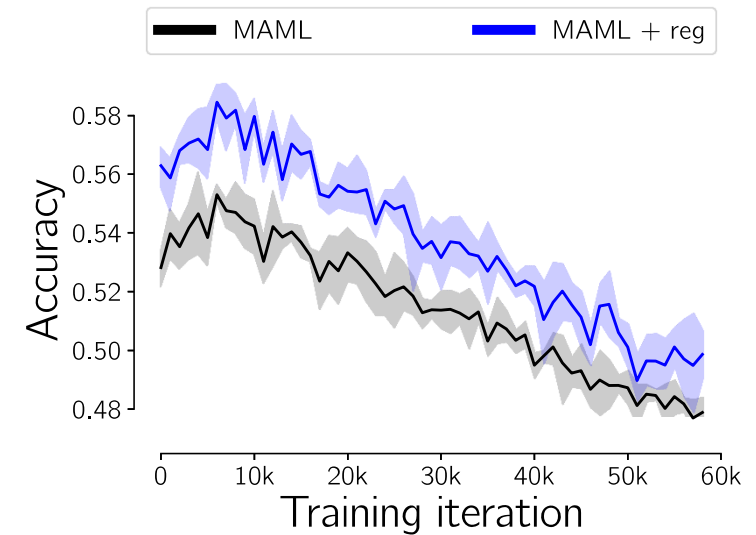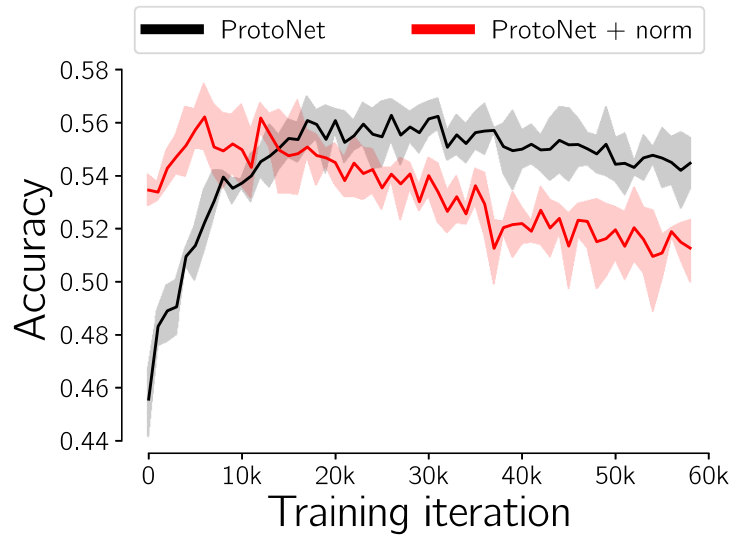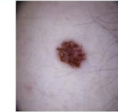Color

**CropDisease:**
Perspective
Natural Images
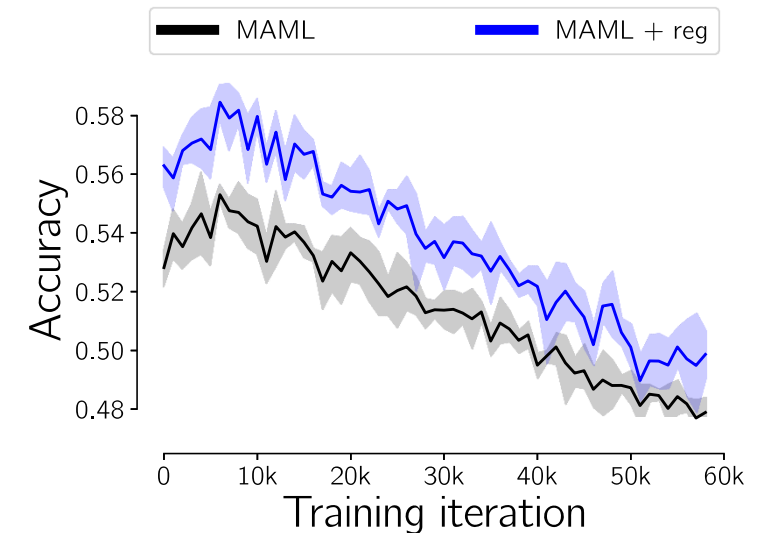Color
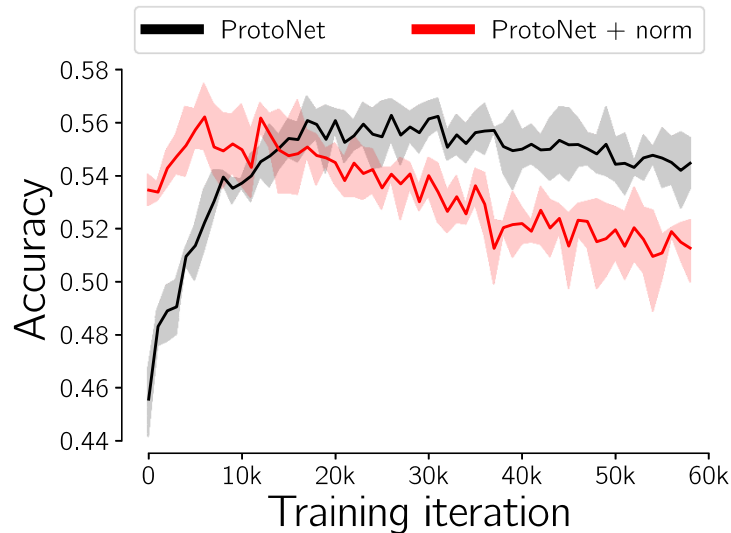
**EuroSAT:**
No Perspective
Natural Images
Color

**ISIC:**
No Perspective
Medical Images
Color

**ChestX:**
No Perspective
Medical Images
Grayscale

5-way
1-shot





× Improvement does **not** translate to cross-domain for *metric-based methods*.

✓ *Gradient-based methods* keep their accuracy gains.

Quentin Bouniot, Ievgen Redko, Romaric Audigier, Angélique Loesch | CAp2021 | 14/06/2021

- **Improving Few-Shot Learning Through Multi-Task Representation Learning Theory**

  - ✓ Connection between Meta-Learning and Multi-Task Representation Learning Theory

  - ✓ Explanations of 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 !

Episodic Training



Regular Training



Learning a single episode

Chen W.-Y. et al., *A Closer Look at Few-Shot Classification.* In ICLR 2019

Regular Training



**Training stage**

Base class data (Many)

Feature extractor

Classifier

$X_b - f_\theta \to C(\cdot | W_b) \to \hat{Y}$

$L_{pred}$

**Fine-tuning stage**

Novel class data (Few)

**Fixed** Feature extractor

Classifier

$X_n - f_\theta \to C(\cdot | W_n) \to \tilde{Y}$

$L_{pred}$

**Classifier** $C(\cdot | W)$

**Baseline**

$f_\theta(x_i) \to$ Linear layer $\to$ Softmax $\sigma \to \tilde{y}_i$

$W \in \mathbb{R}^{d \times c}$    $\tilde{y}_i = \sigma(W^\top f_\theta(x_i))$

**Baseline++**

$f_\theta(x_i) \to$ Cosine distance $\to$ Softmax $\sigma \to \tilde{y}_i$

$W = [w_1, w_2, ...w_c] \in \mathbb{R}^{d \times c}$
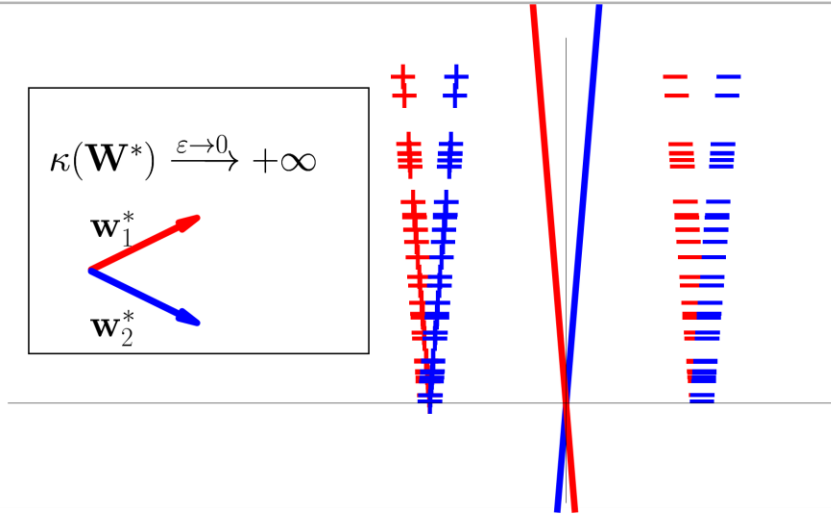
Adapted from [Chen19]

- *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
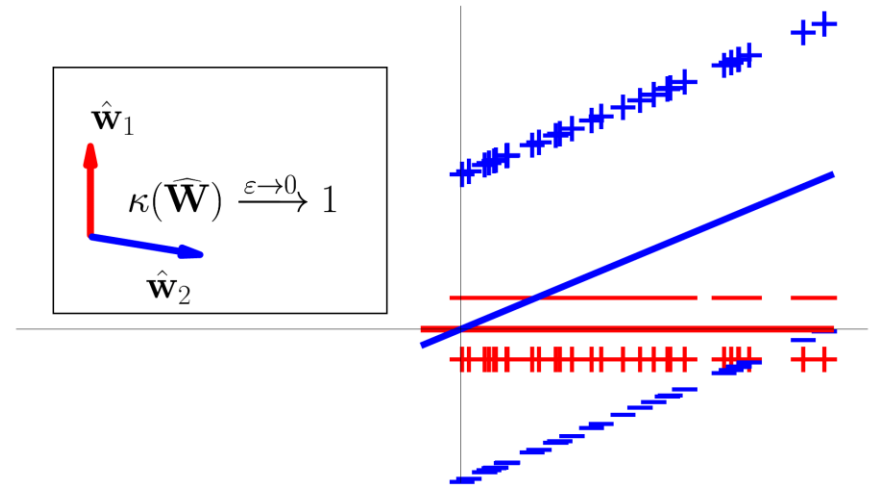
Gidaris S. et Komodakis N., *Dynamic Few-Shot Visual Learning Without Forgetting.* In CVPR 2018

*Given $\mathbf{W}^*$ such that $\kappa(\mathbf{W}^*) \gg 1$, can we learn $\widehat{\mathbf{W}}$ with $\kappa(\widehat{\mathbf{W}}) \approx 1$ 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

► $\kappa(\mathbf{W}_N)$ shows **dynamics** during training, but values are not comparable

► $\kappa(\mathbf{W})$ is **intractable to compute** during training.