

# DATA SCIENCE (資料科學) SPRING, 2023

## LECTURE 4: FEW-SHOT LEARNING & META LEARNING

Shuai, Hong-Han (帥宏翰)

Associate Professor

Department of Electronics and Electrical Engineering

National Yang Ming Chiao Tung University



PC: Midjourney

Online tutorial:

Thanks to the slides from Spyros Gidaris in CVPR 2020 Tutorial

Week	Contents
1	Introduction to Data Science
2	Data Crawling [HW1: Crawling Releases]
3	Peace Memorial Day
4	Model Compression [HW2: Model Compression]
5	Reinforcement Learning
6	Few-shot Learning [HW3: FSL]
7	Meta Learning



# Motivations

# Machines vs. Humans

- Machines: specialists.
- Humans: generalists.
- Machines ~100 billion parameters
- Humans ~100 trillion connections  
x1000 factor and closing fast.
- Machines have super-human performance on many specific tasks that are well defined and formulated.

# Human Brain Connectome

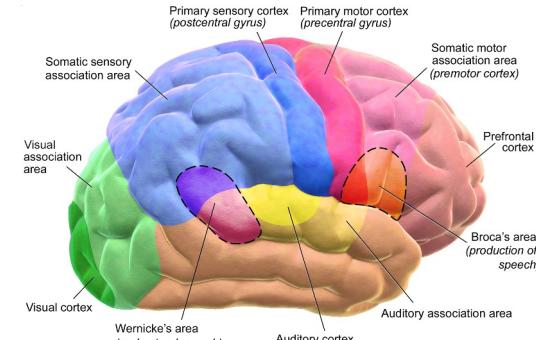
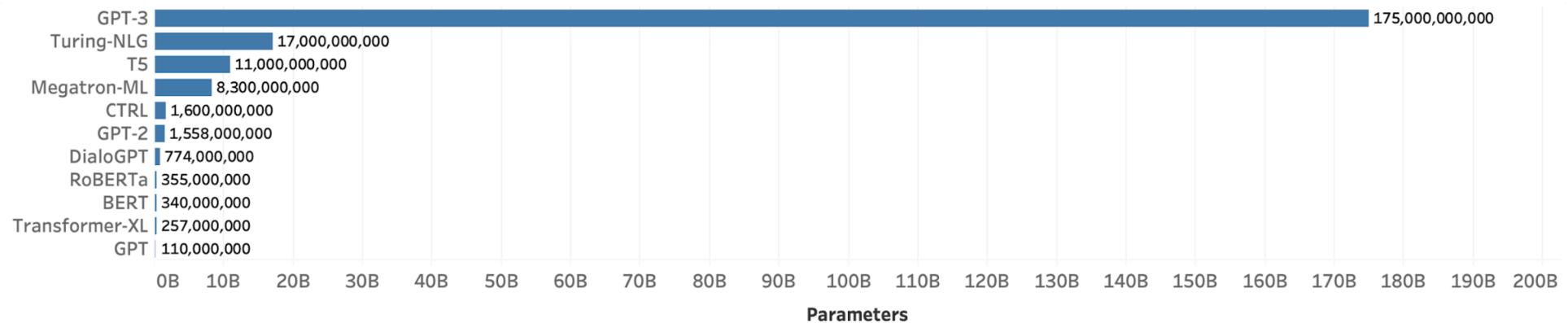


Image Source: Wikipedia

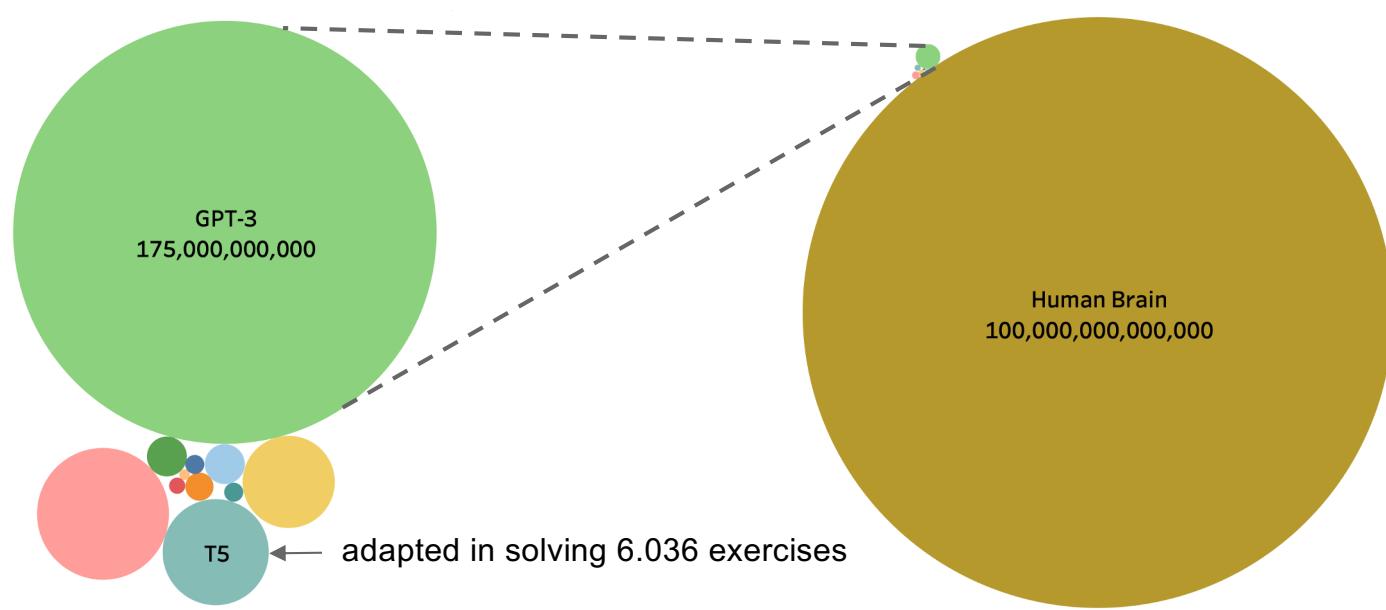
- 100 Billion neurons (1 Billion neurons in cat brain)
- 100 Trillion connections: each neuron connects to 5k-200k
- 10k different types of neurons
- 1k new neurons per day our entire life

# Transformer Parameters



- 3 orders of magnitude less parameters than number of connections in human brain

# Number of Connections or Parameters



- Transformers have 3 orders of magnitude less parameters than number of connections in human brain

# Super-Human ML Systems: AlphaX

- AlphaZero: board games
- AlphaStar: multiplayer online games
- AlphaFold: protein structure prediction
- AlphaD3M: automated machine learning
- AlphaStock: stock trading
- ..
- AlphaDogfight: fighter pilot

- Self driving grand challenge 2 decades ago: competitive.

Recent collaborative efforts

- Data-Driven Discovery of Models (D3M): AutoML
- Learning with Less Labels (LwLL): few shot learning
- Lifelong Learning Machines (L2M): online learning
- Machine Common Sense (MCS)

Automated machine learning, few shot learning, online learning, learn to read, natural language understanding

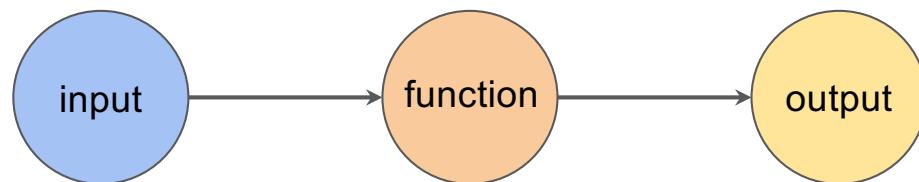
# **Meta Learning Definitions**

# Definitions

- Supervised learning
  - Transfer learning
  - Meta learning
  - Automated machine learning
- 
- Adaptation
  - Multi-task learning
  - Few-shot learning
  - Online learning

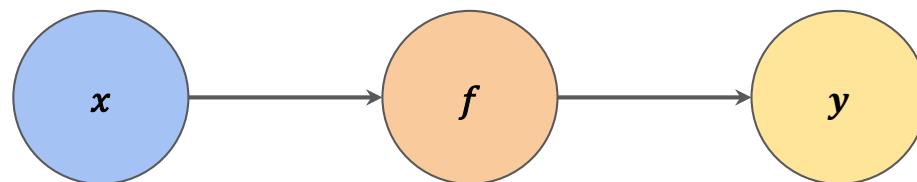
# Observation

- Input  $\mathbf{x}_{d \times 1}$
- Function  $f$
- Output  $y_{1 \times 1}$



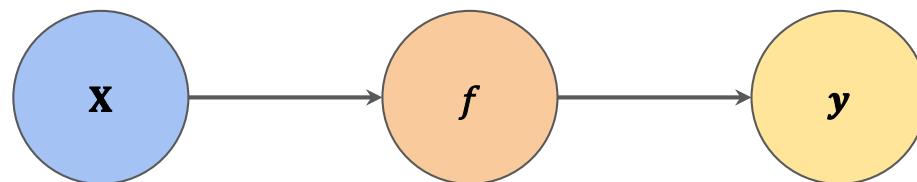
# Observation

- Input  $\mathbf{x}_{d \times 1}$
- Function  $f$
- Output  $y_{1 \times 1}$



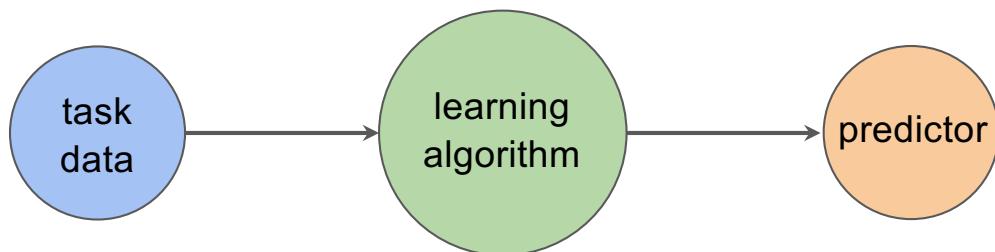
# Observations

- Input  $\mathbf{X}_{d \times m}$
- Function  $f$
- Output  $\mathbf{y}_{m \times 1}$

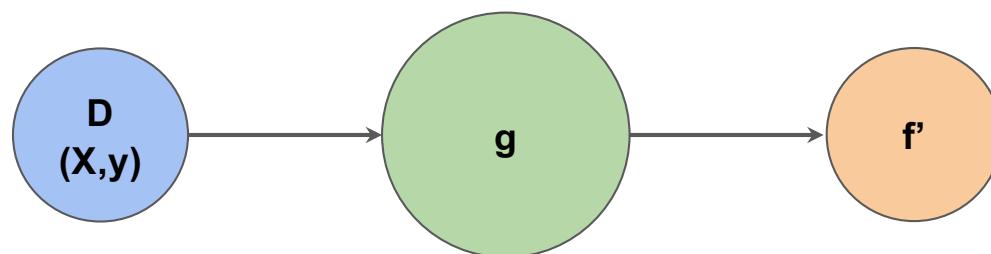


$$\mathbf{y} = f(\mathbf{X})$$

# Supervised Learning

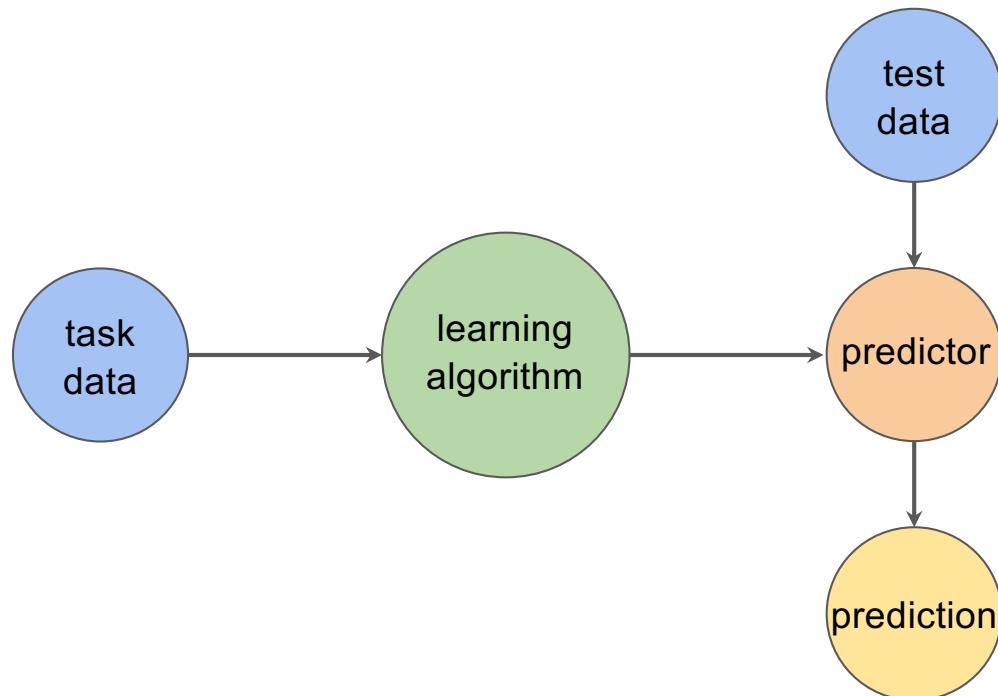


# Supervised Learning

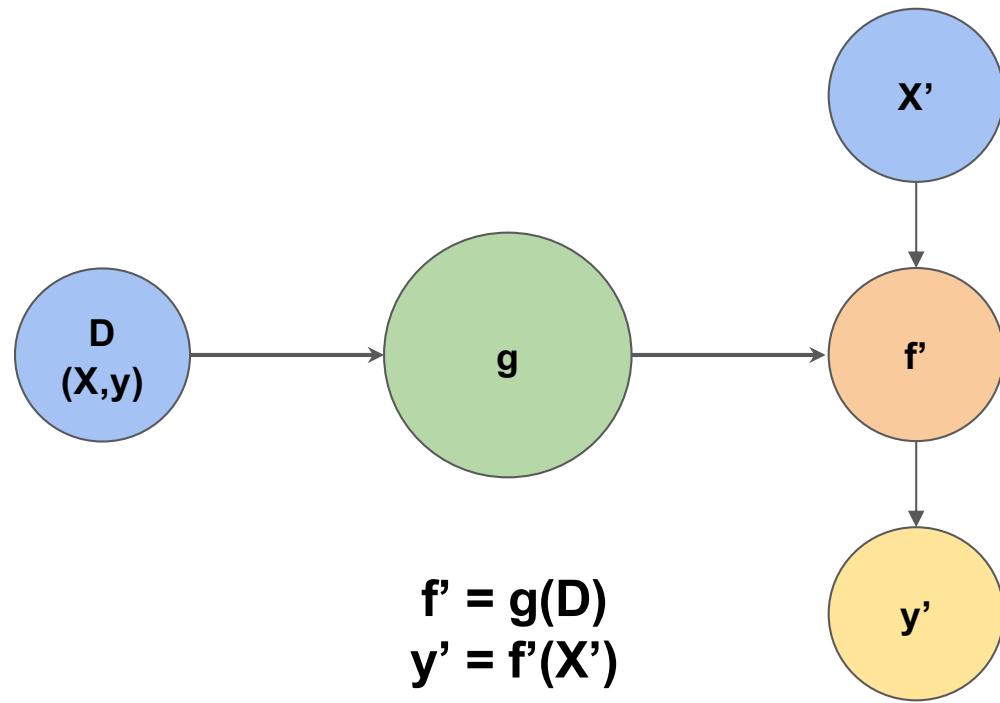


$$f' = g(D)$$

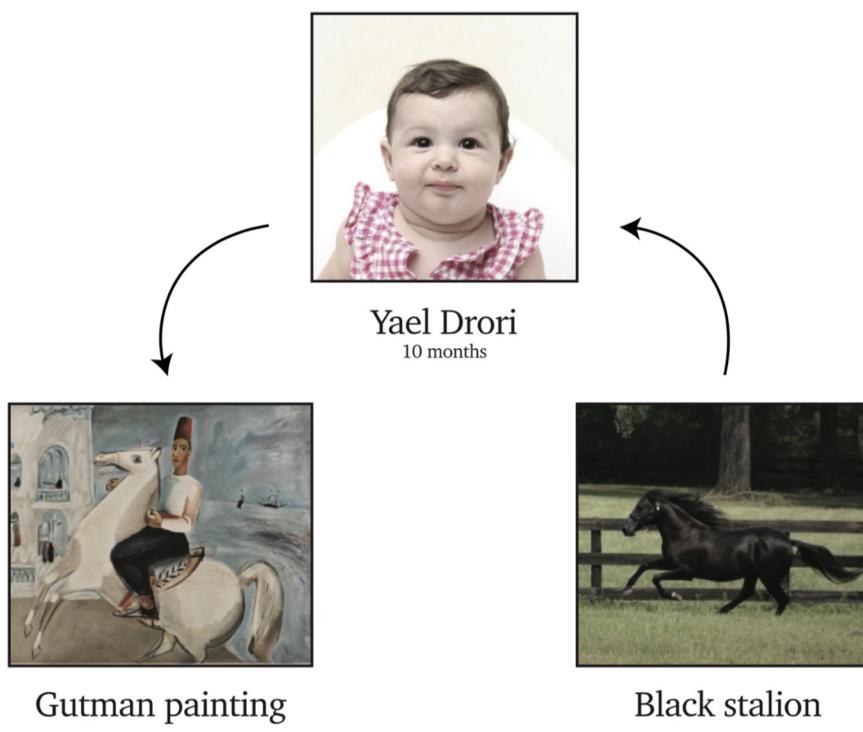
# Supervised Learning



# Supervised Learning

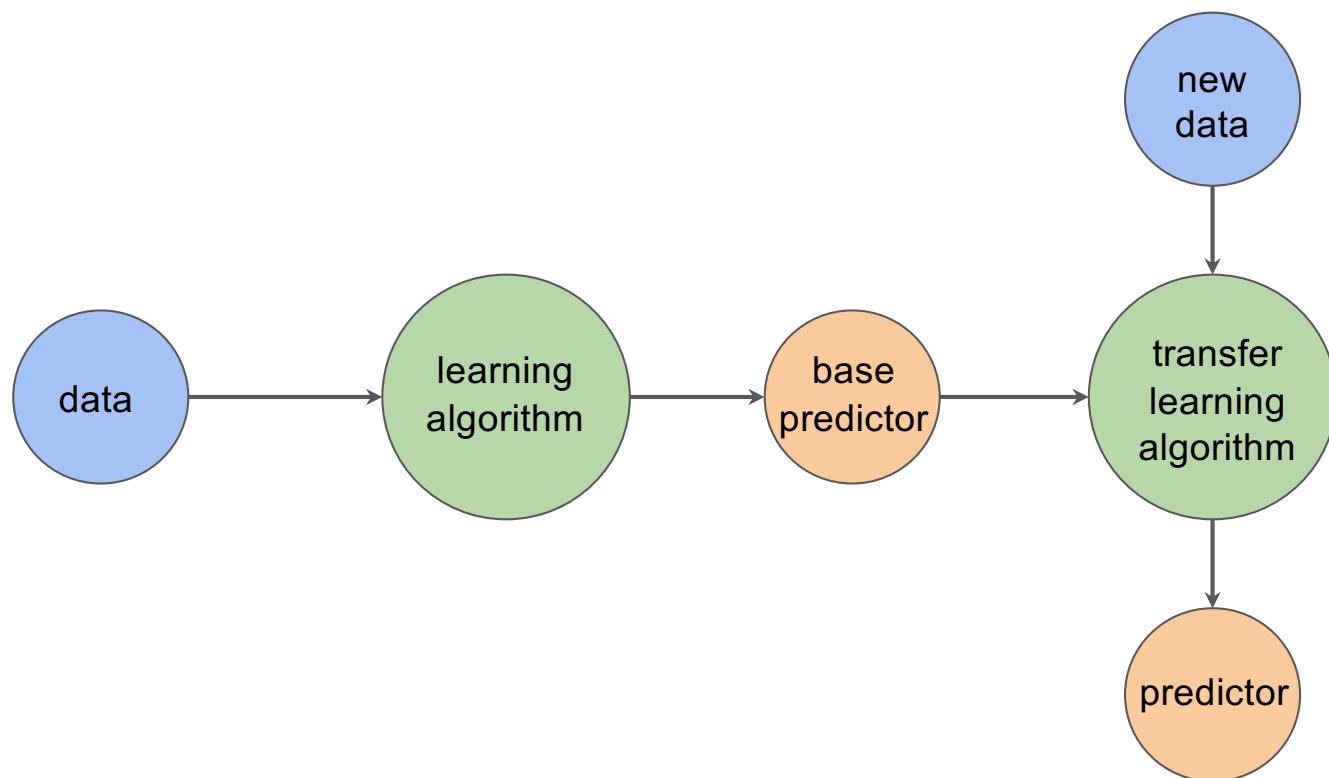


# Transfer Learning

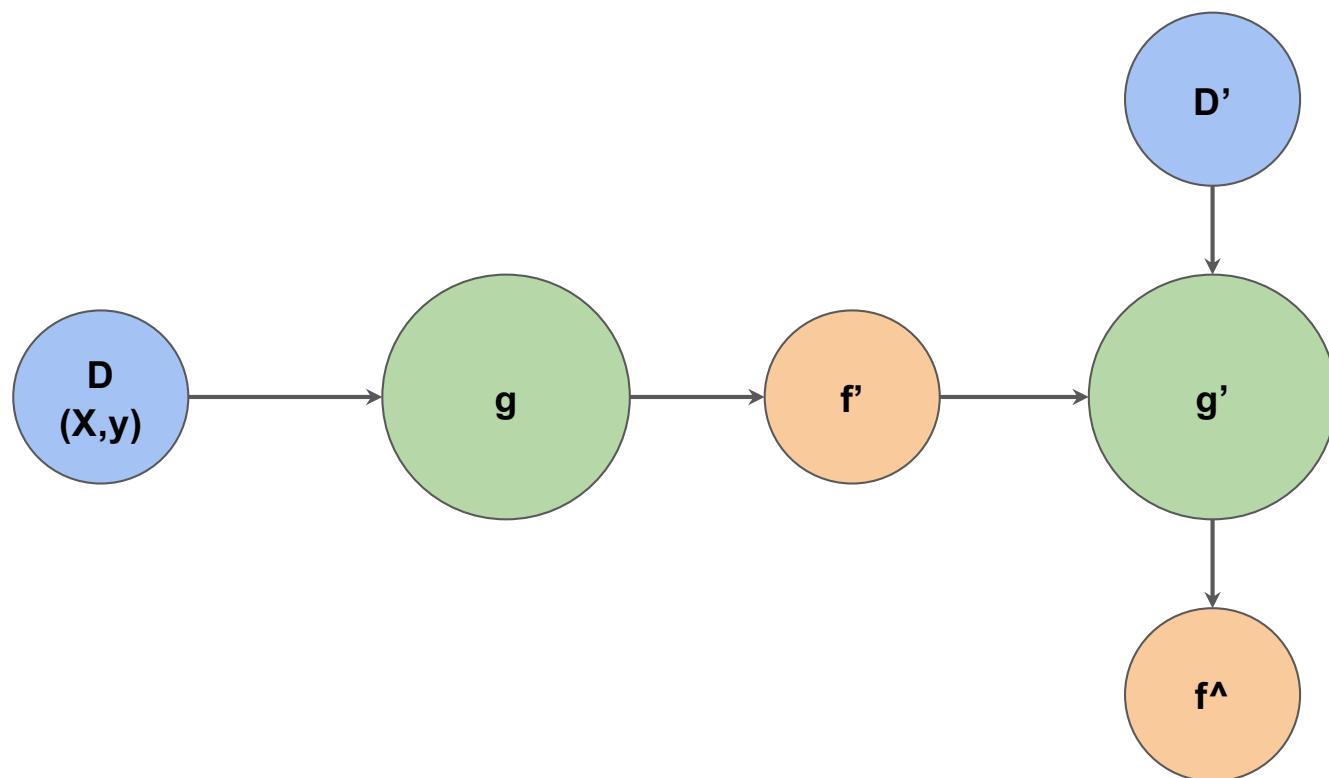


Source: Deep Learning course 2017, Iddo Drori

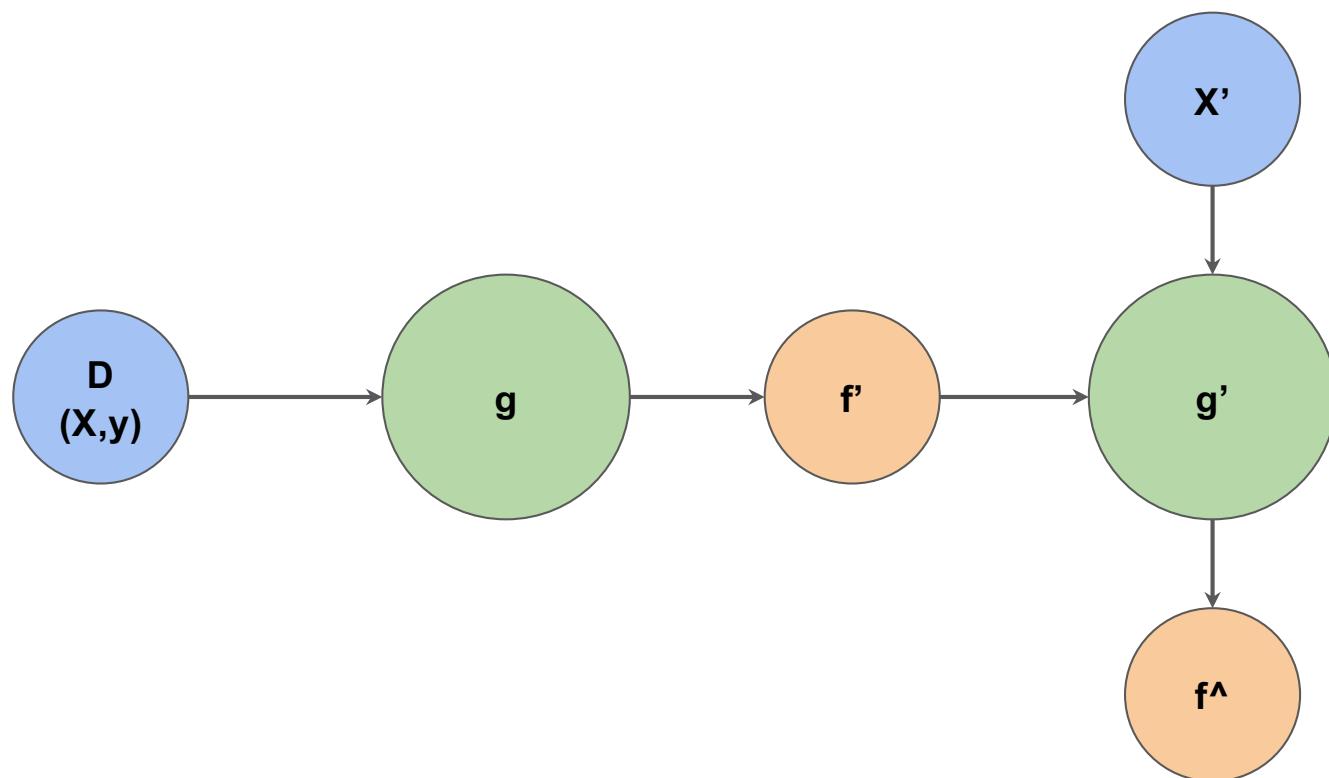
# Transfer Learning



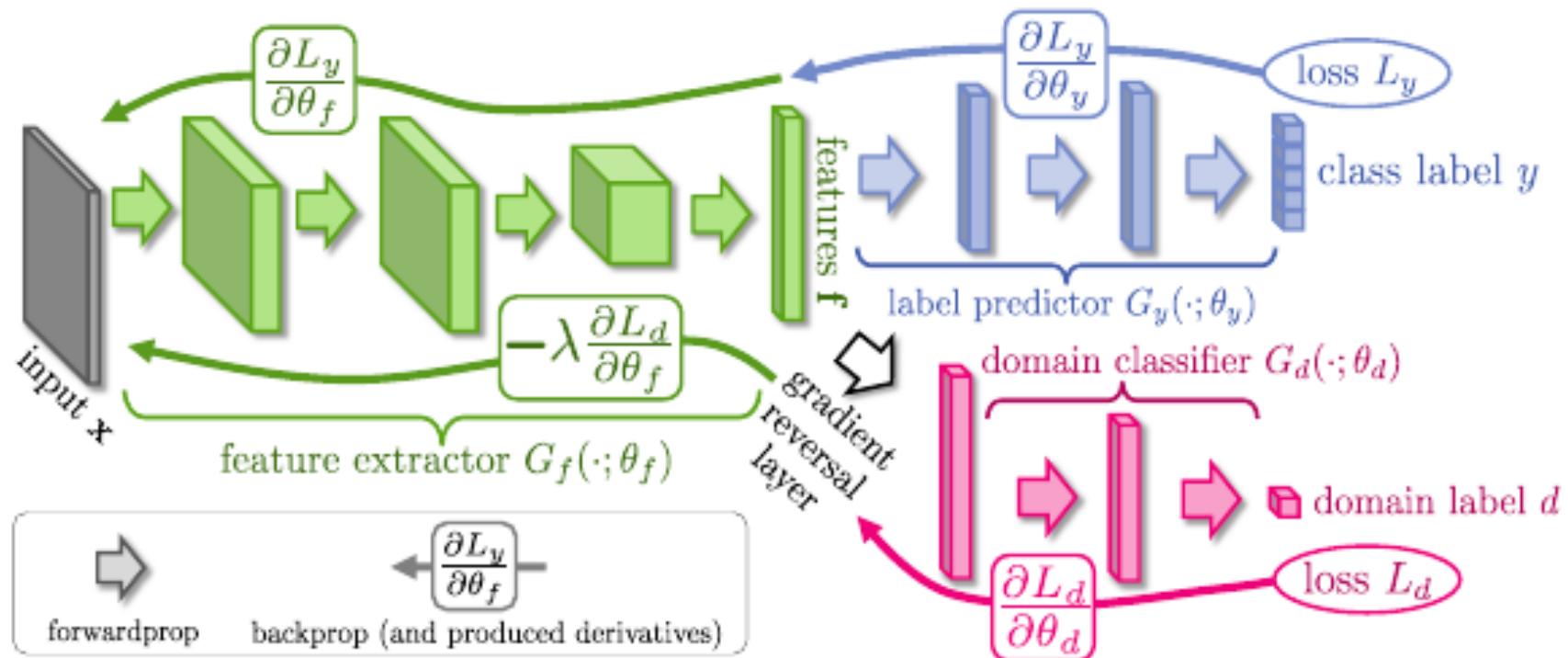
# Transfer Learning



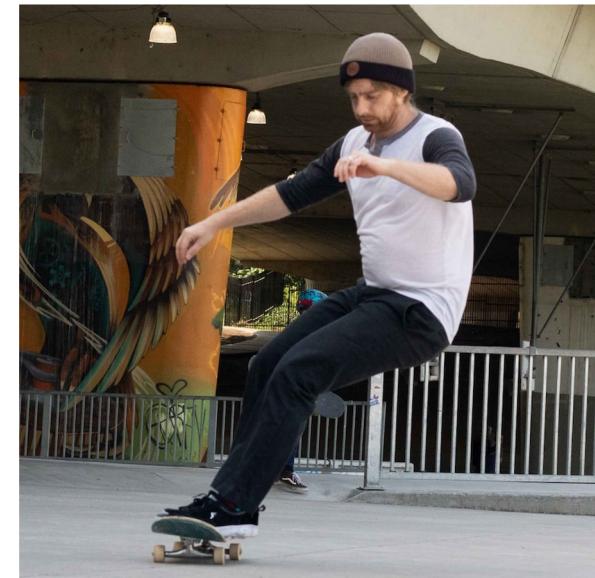
# Adaptation (Unsupervised Transfer Learning)



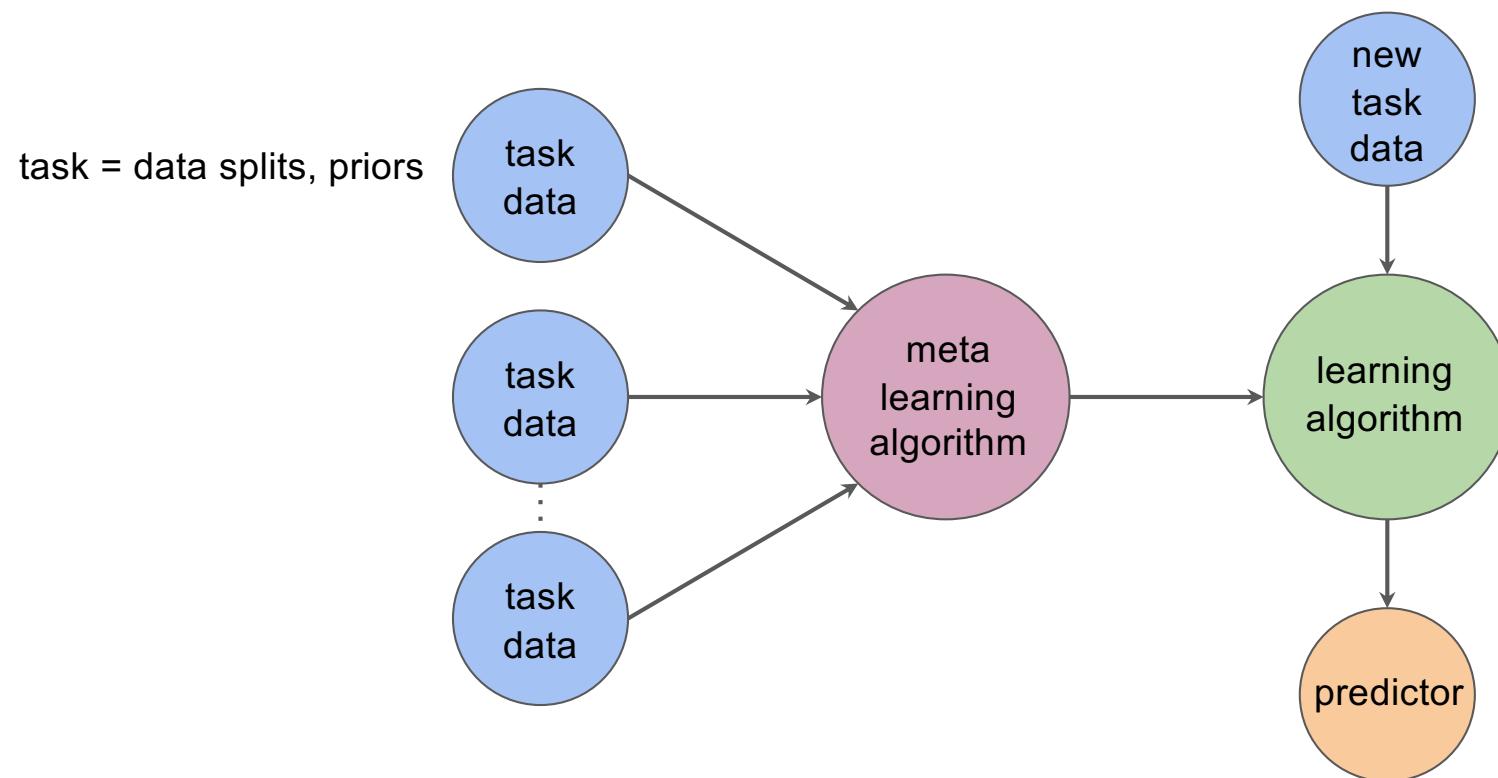
# Unsupervised Domain Adaptation by Backpropagation



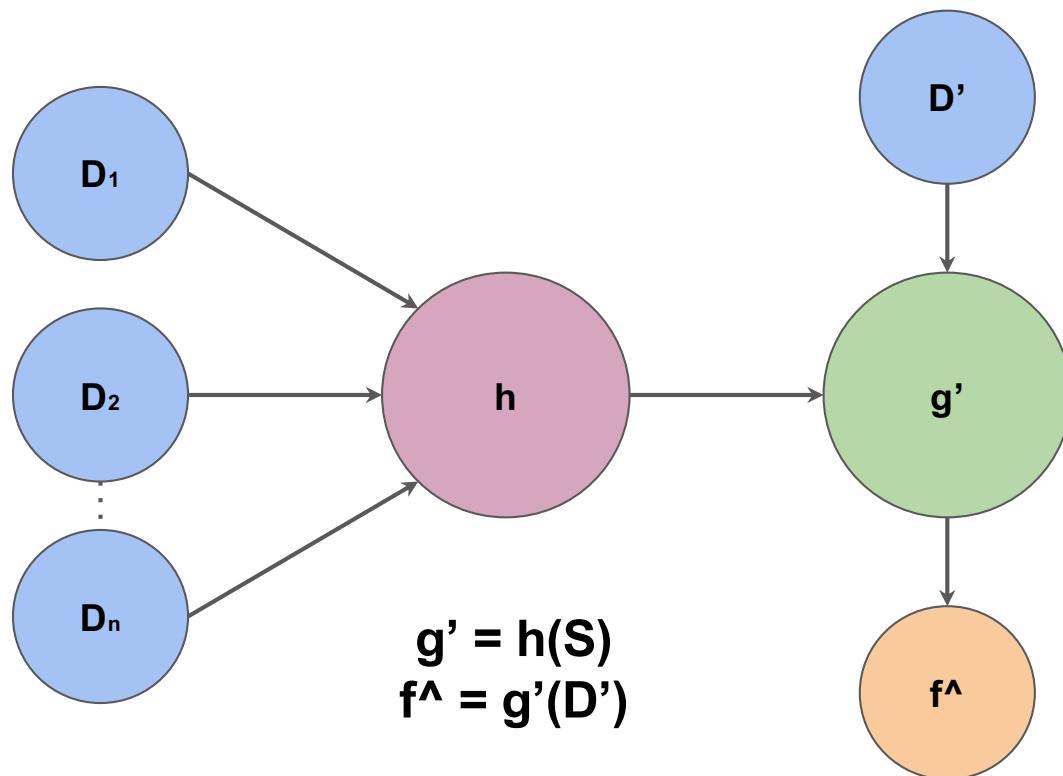
# Human Example



# Meta Learning

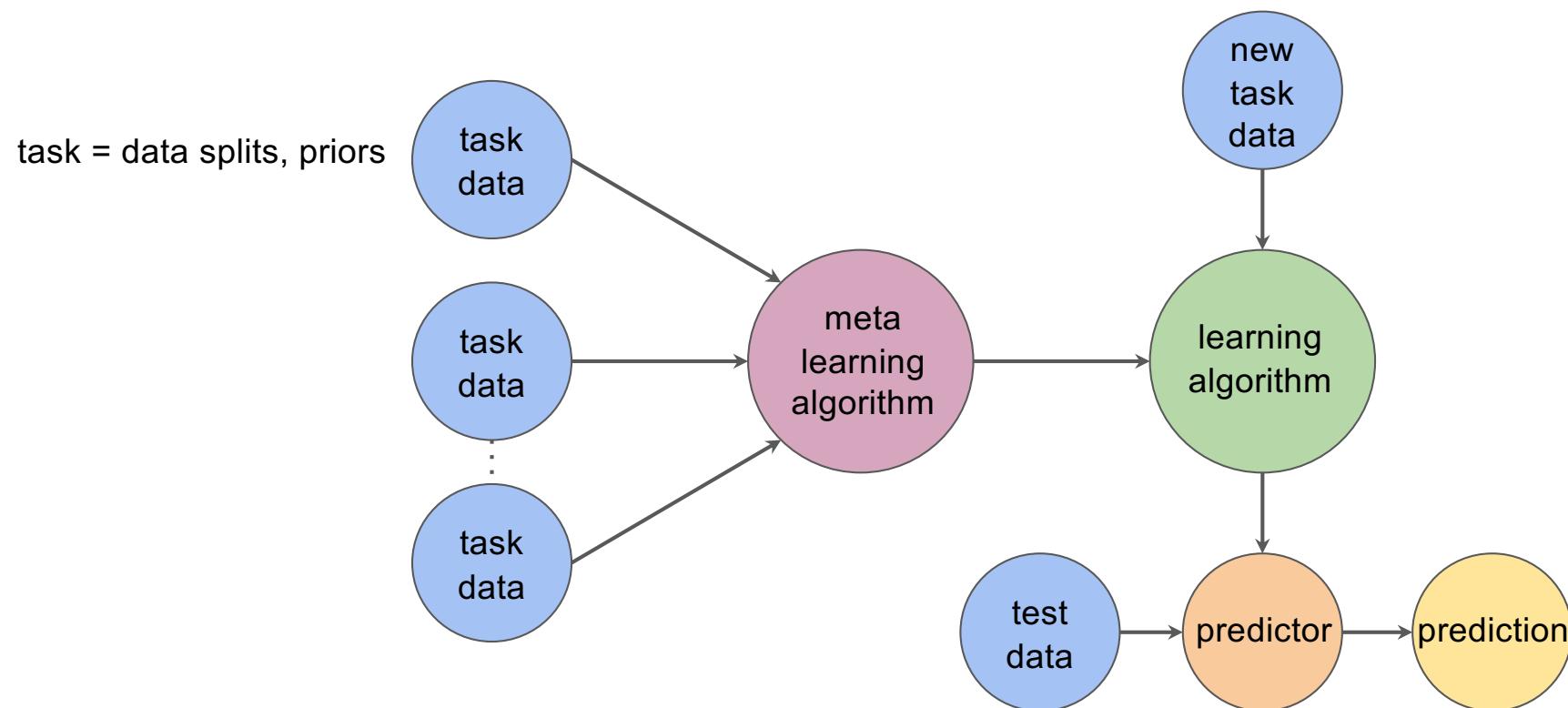


# Meta Learning

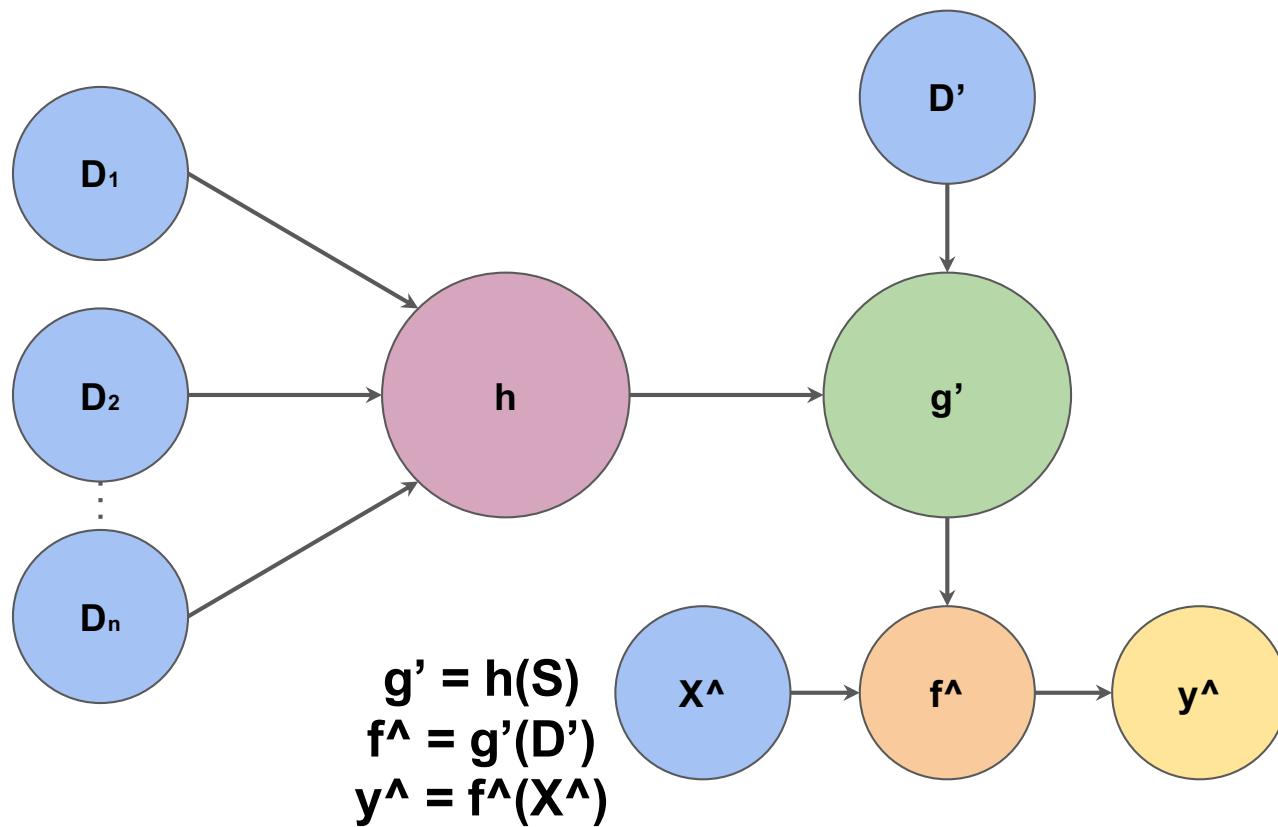


$$g' = h(S)$$
$$f^\wedge = g'(D')$$

# Meta Learning

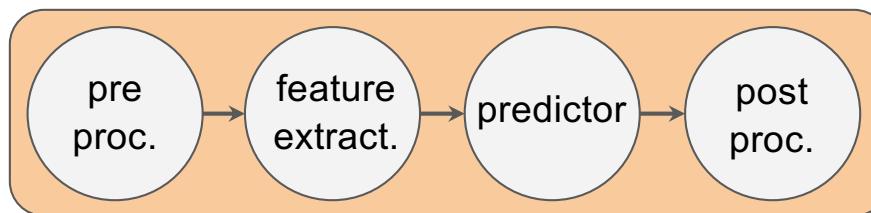


# Meta Learning

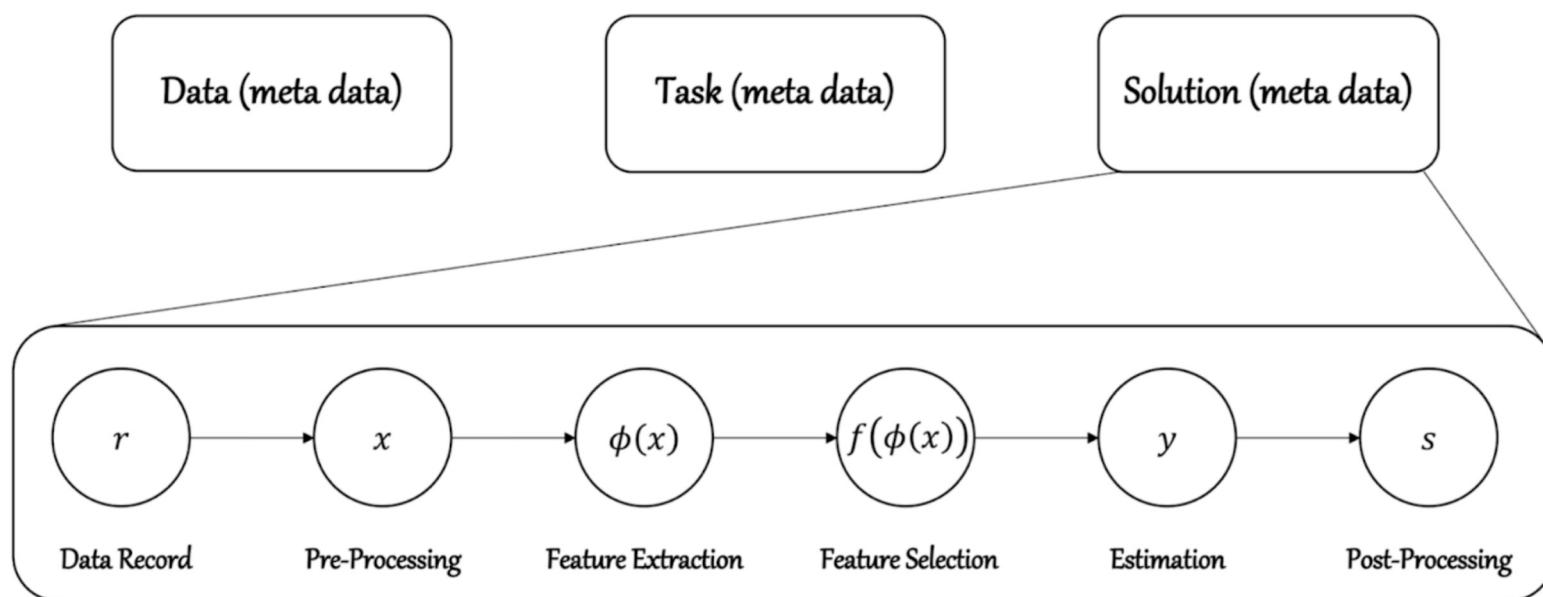


# Machine Learning System

- Predictor is part of a machine learning system
- Built from data science / machine learning primitives
- Machine learning primitive = {PCA, SVM, NN,...}
- Example machine learning pipeline:

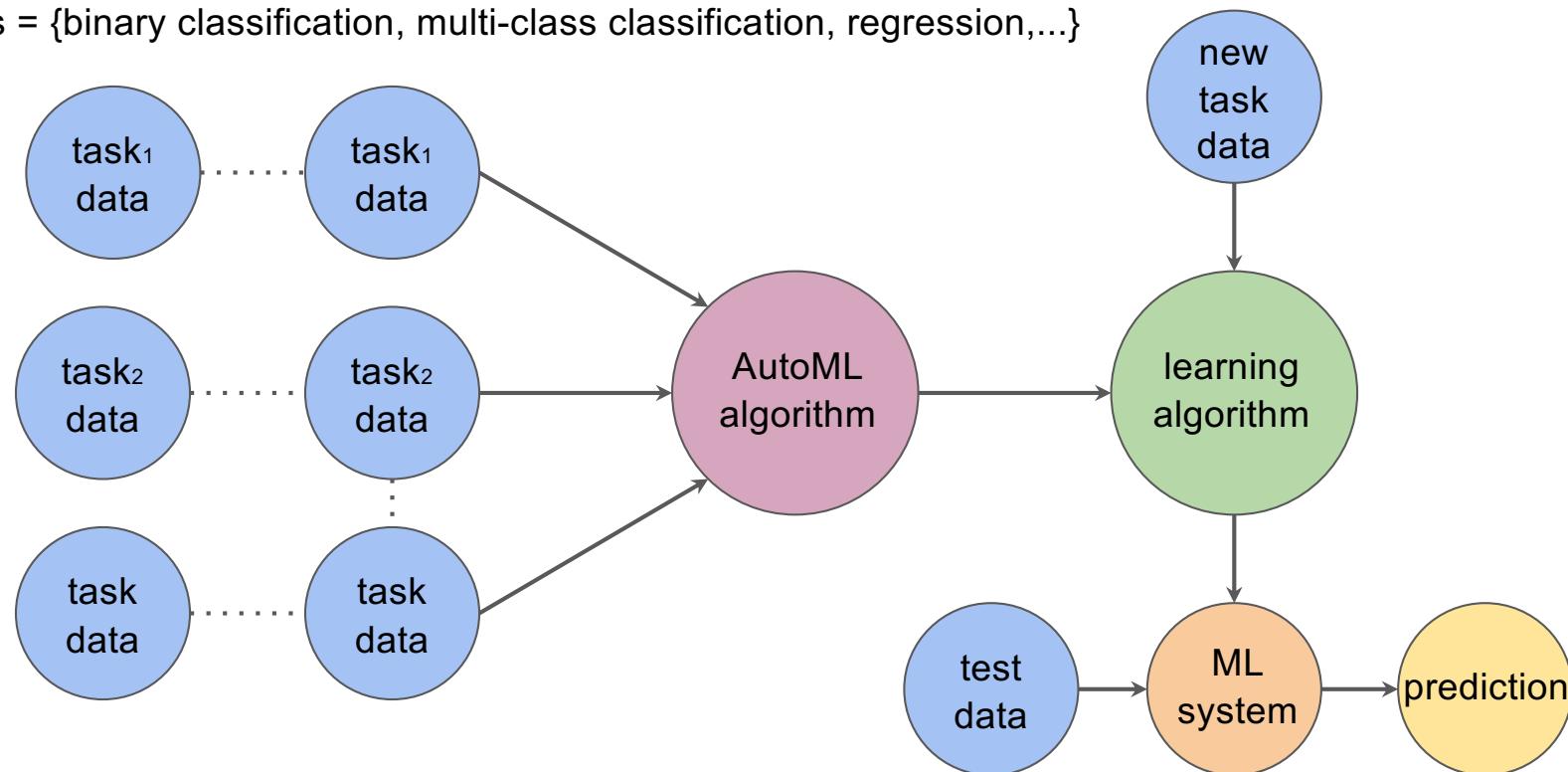


# Machine Learning Pipeline

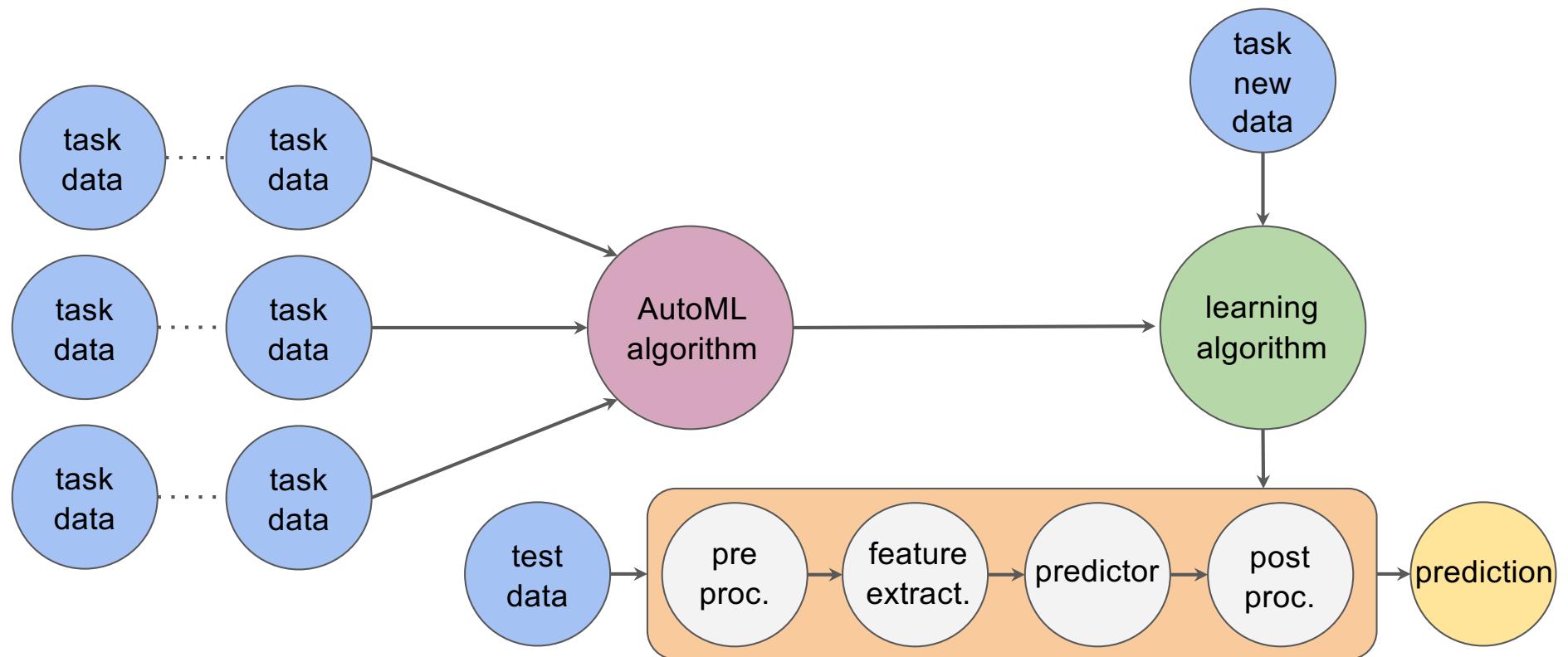


# Automated Machine Learning

tasks = {binary classification, multi-class classification, regression,...}



# Automated Machine Learning (AutoML)

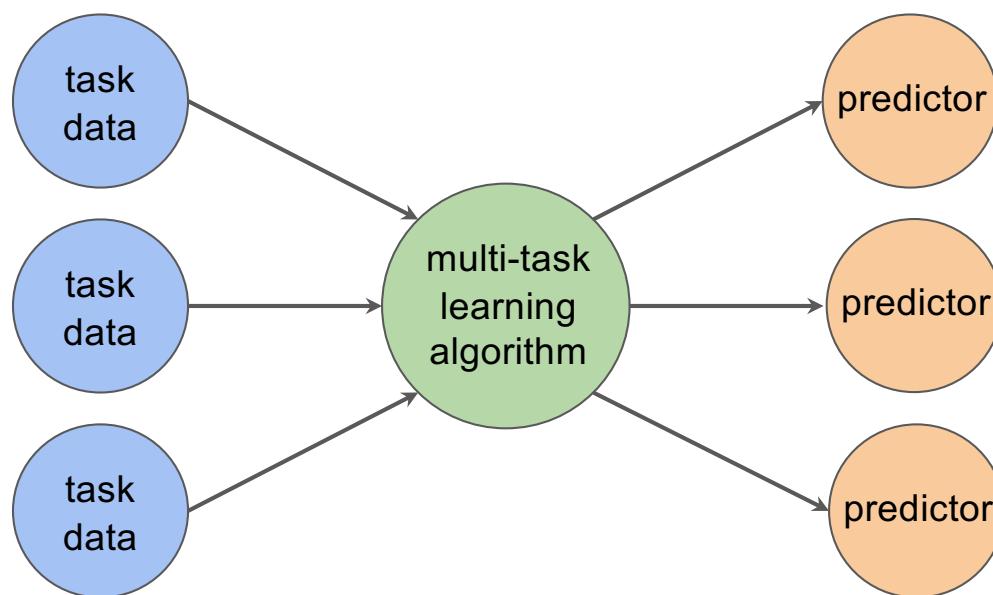


# Learning to Learn

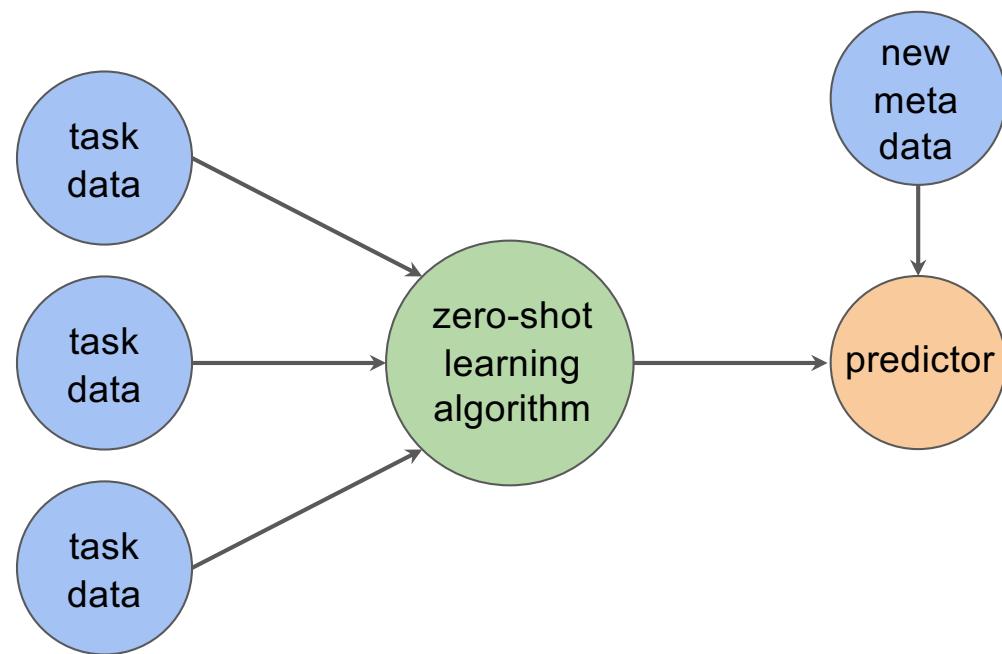
- Machine Learning: learn parameters of **M**
- Learning to learn: learn **M** and parameters

where **M** is a classifier or machine learning pipeline or machine learning algorithm or reinforcement learning method, ect.

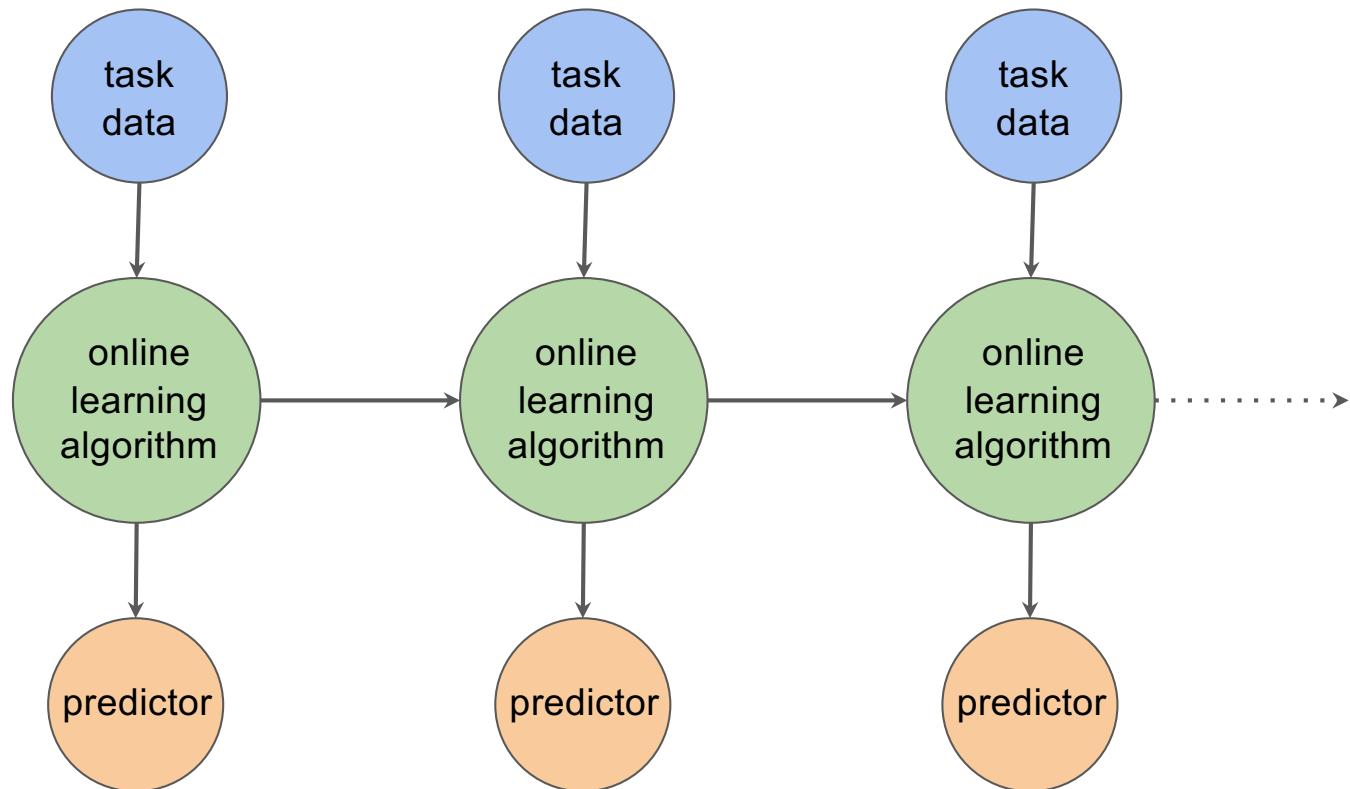
# Multi-Task Learning



# Zero-Shot Learning



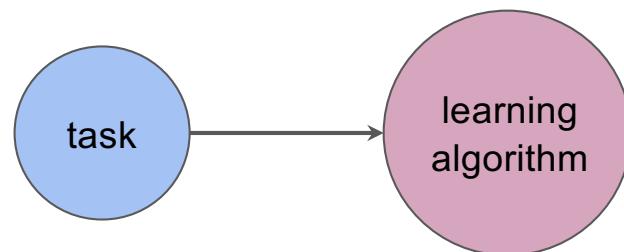
# Online Learning (Sequential, Lifelong learning)



# Learning Tasks

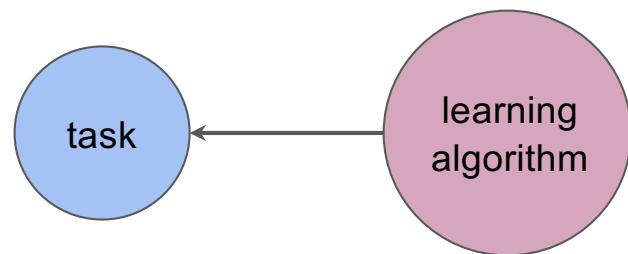
# Learn to Solve a given Task

- Optimization
- Neural networks, relatively easy

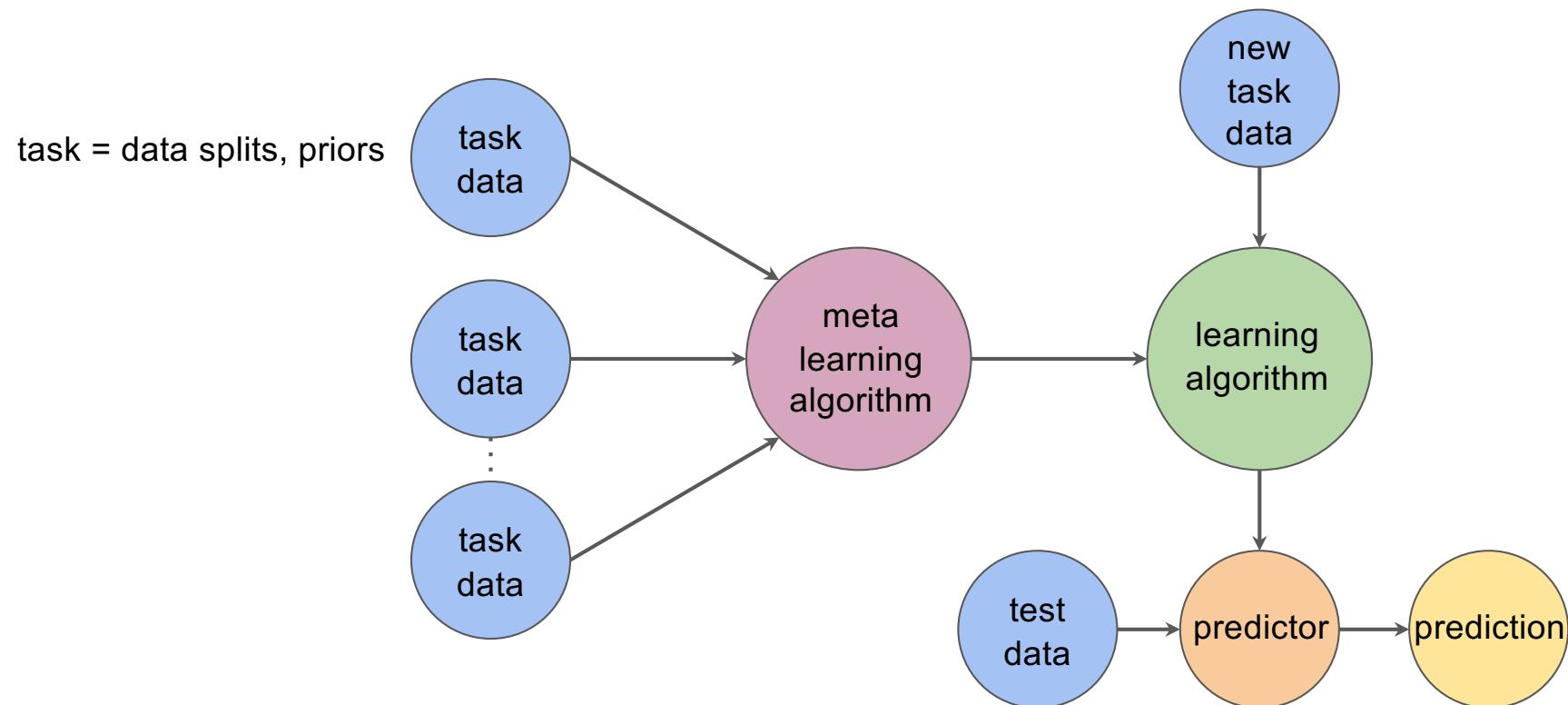


# Learning the Task

- Reverse direction

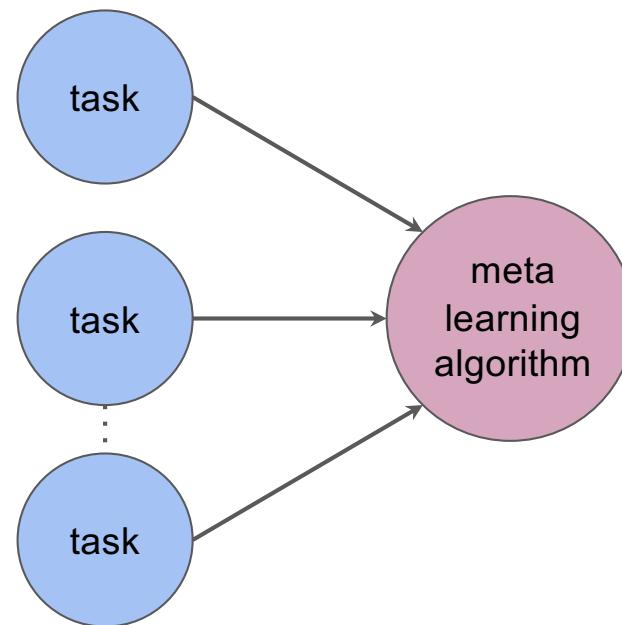


# Meta Learning



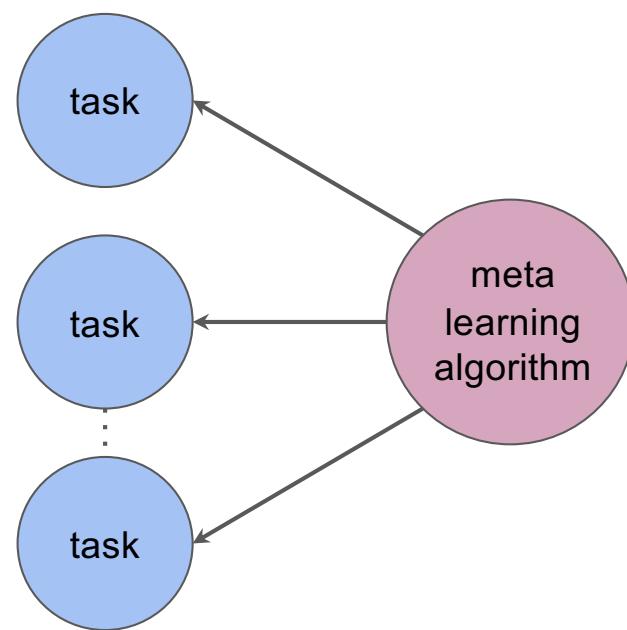
# Meta Learning

- Optimization

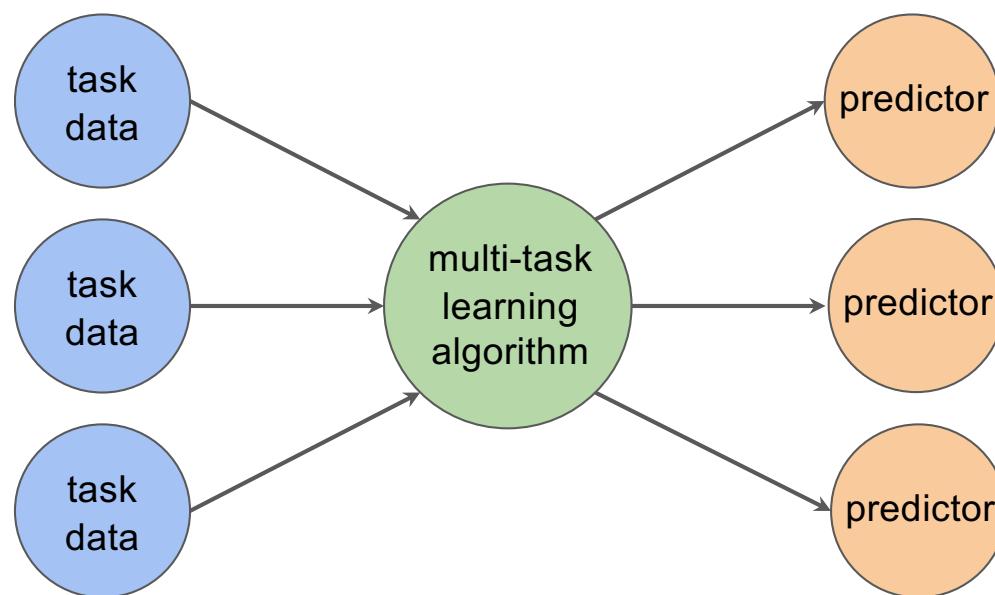


# Learning the Tasks

- Reverse direction

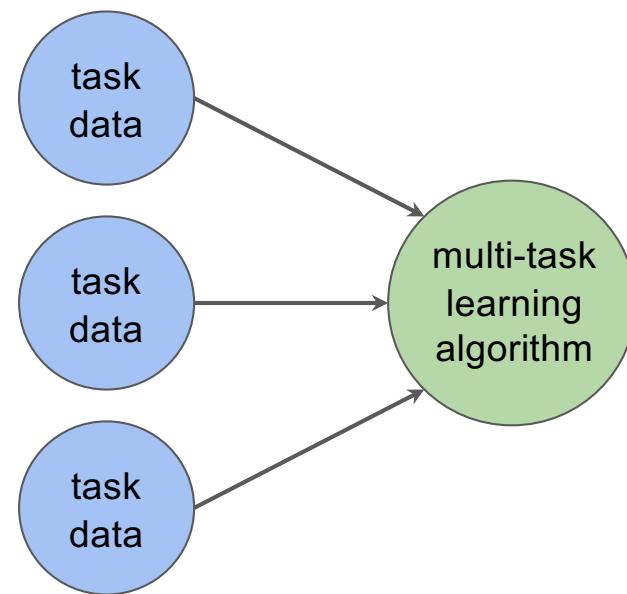


# Multi-Task Learning



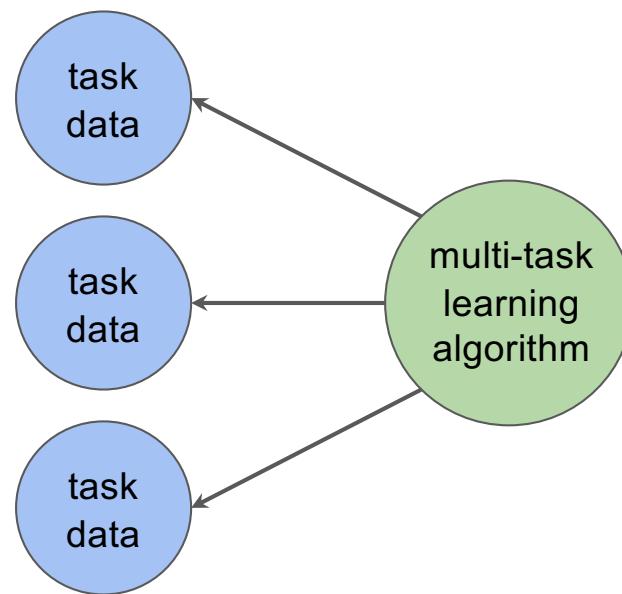
# Multi-Task Learning

- Optimization



# Learn the Tasks

- Reverse direction



# Curriculum Learning

- Append loss on top of any existing loss
- Downweight contributions of hard samples, with a large loss

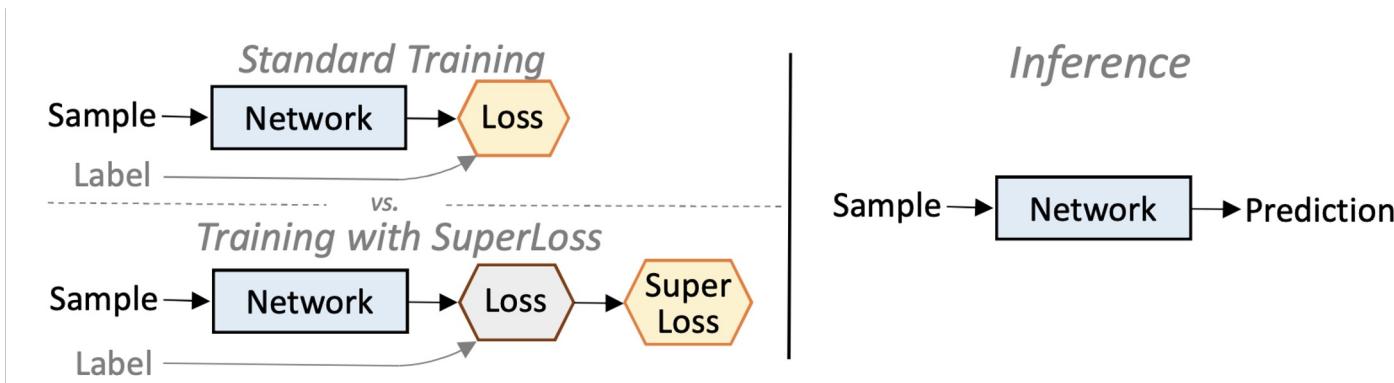


Figure source: SuperLoss: A generic loss for robust curriculum learning, Castells et al, 2020

# Self-Supervised Learning

- No negative examples
- Better performance than contrastive methods

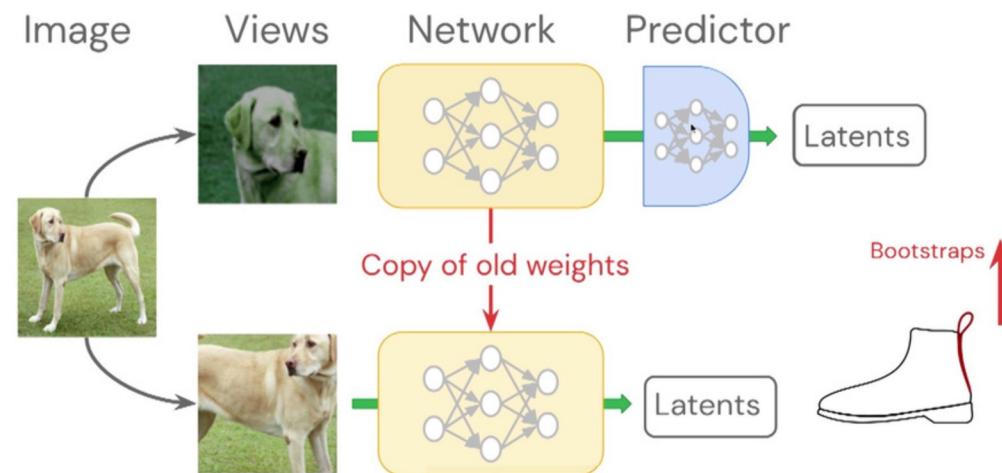


Figure source: Bootstrap your own latent: A new approach to self-supervised learning, Grill et al, 2020



# Agenda

- Introduction
- Main types of few-shot algorithms
- Few-shot learning without forgetting

# Agenda

## Introduction

- **Few-shot learning**
- Meta-learning paradigm
- How to evaluate

## Main types of few-shot learning

## Few-shot learning without forgetting

## Support Set

Armadillo



Pangolin



## Support Set

Armadillo



Pangolin



Query



Armadillo or Pangolin?

## Support Set

Husky



Malamute



## Support Set

Husky



Malamute



Query



Husky or Malamute?

# Training Set

Husky



⋮



Elephant



⋮



Tiger



⋮



Macaw



⋮



Car



⋮



**Are they the same kind of animal?**



**Are they the same kind of animal?**

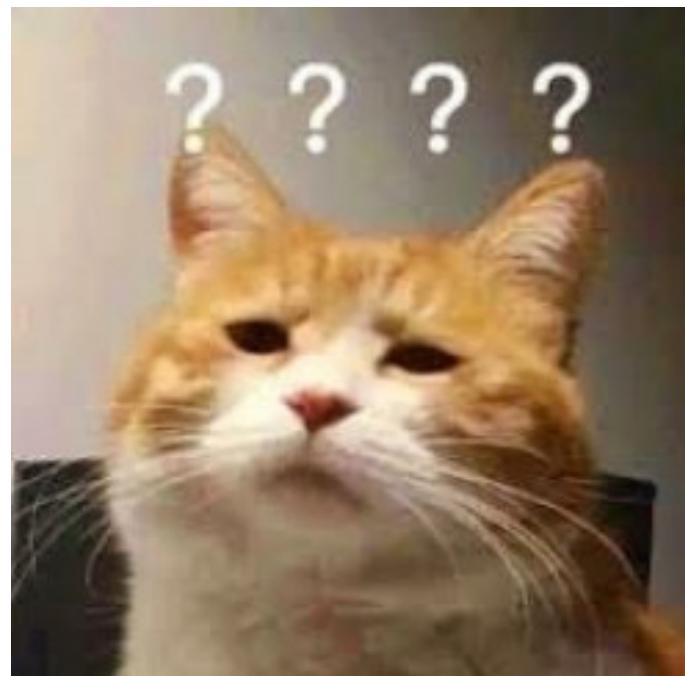


**Are they the same kind of animal?**



# Few-Shot Learning

**Query:**



# Few-Shot Learning

Query:



Support Set:

Fox



Squirrel



Rabbit



Hamster



Otter



Beaver



# Meta Learning

Reference:

Fei-Fei, Fergus, & Perona. One-shot learning of object categories. *IEEE Transactions on PAMI*, 2006.

- Few-shot learning is a kind of meta learning.
- Meta learning: learn to learn.

# Few-Shot Learning



# Few-Shot Learning



Give the kid some cards:

Fox



Squirrel



Rabbit



Hamster



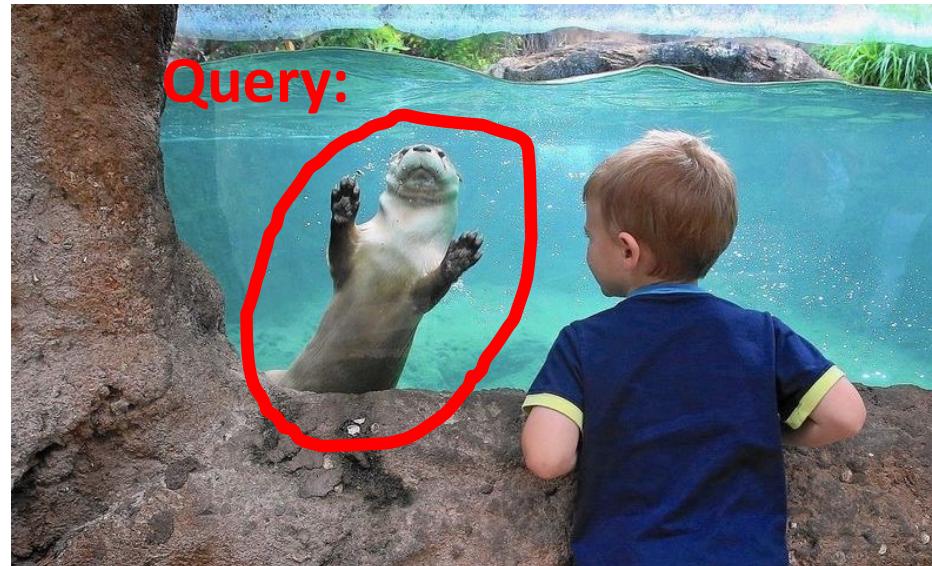
Otter



Beaver



# Meta Learning



Support set:

Fox



Squirrel



Rabbit



Hamster



Otter



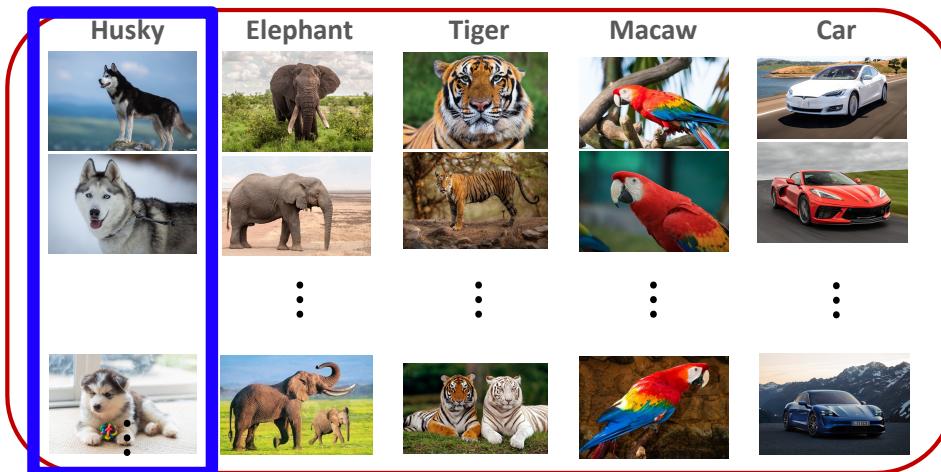
Beaver



# Supervised Learning vs. Few-Shot Learning

- Traditional supervised learning:
  - Test samples are **never seen before**.
  - Test samples are from **known classes**.

Training Set



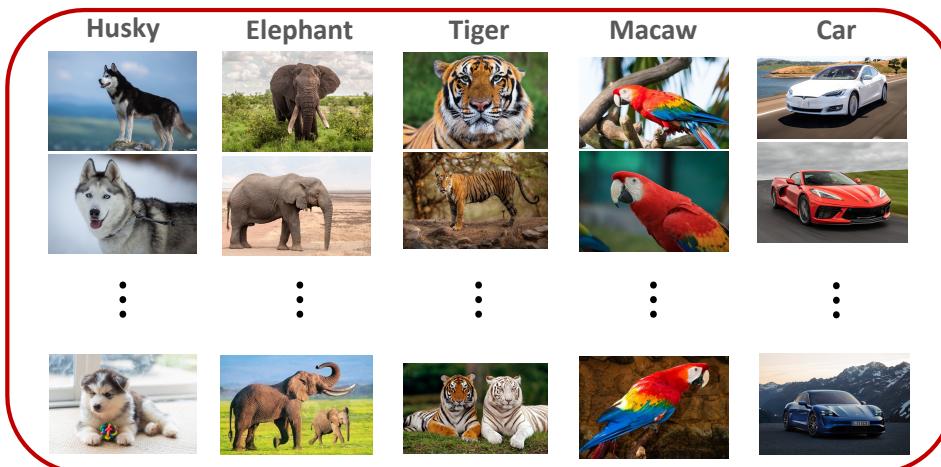
Test Sample



# Supervised Learning vs. Few-Shot Learning

- Few-shot learning:
  - Query samples are never seen before.
  - Query samples are from unknown classes.

Training Set

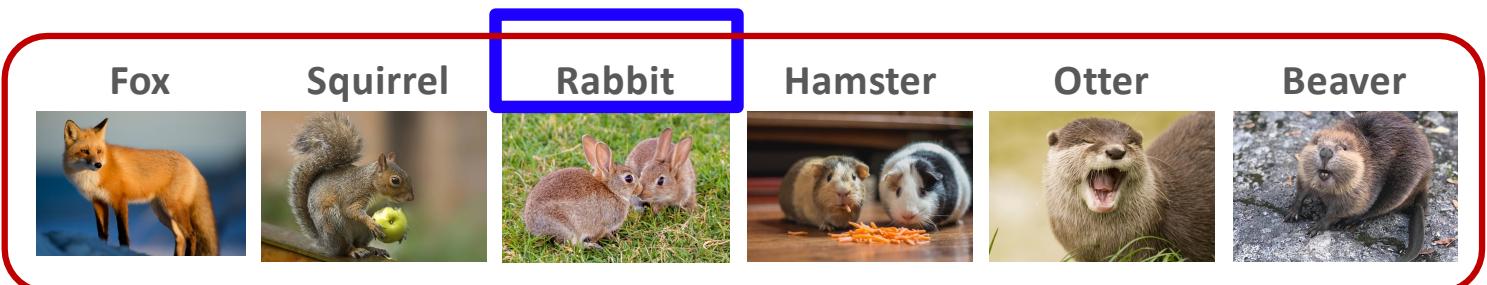


Query Sample

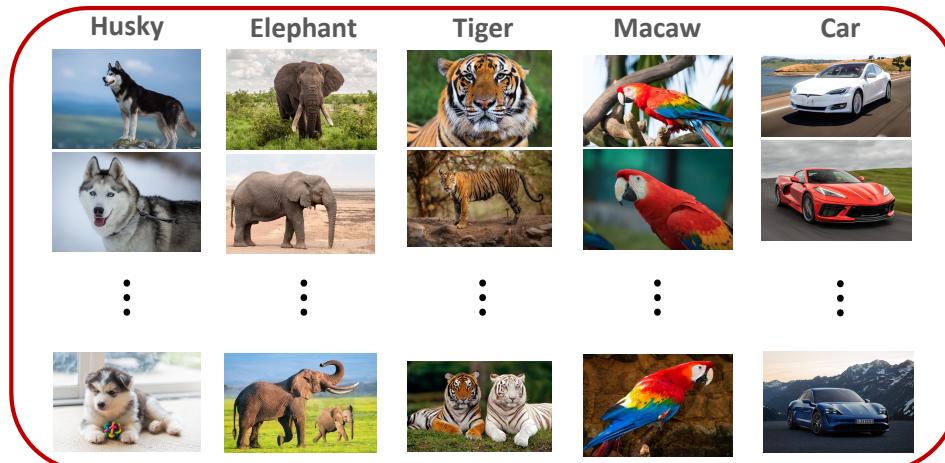


# Training Set, Support Set, and Query

Support Set:



Training Set



Query Sample



# $k$ -way $n$ -shot Support Set

Support Set:



- $k$ -way: the support set has  $k$  classes.
- $n$ -shot: every class has  $n$  samples.

## *k-way* *n-shot* Support Set

Support Set:

Squirrel



Rabbit



Hamster



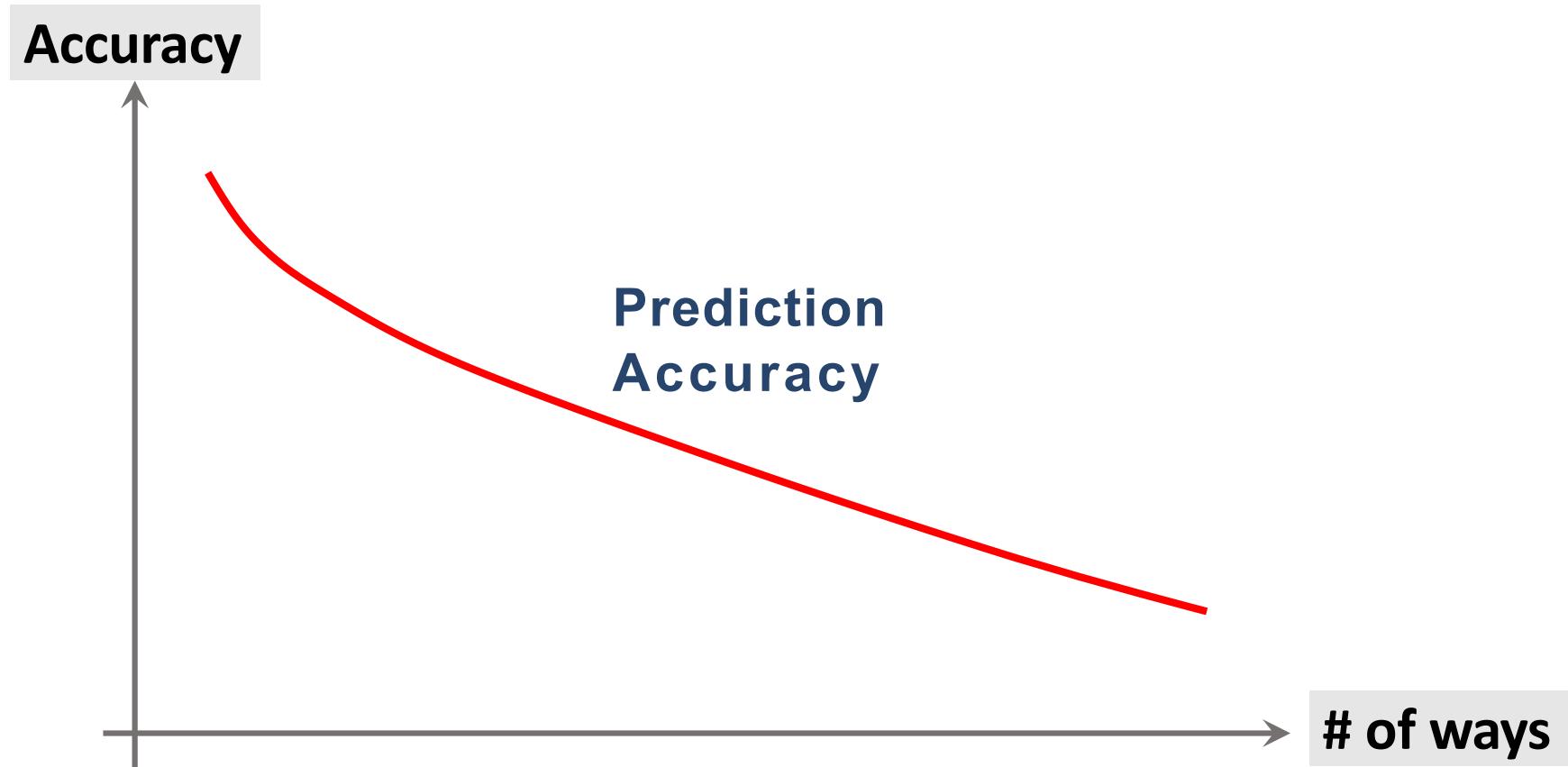
Otter



2-shot



4-way



Squirrel



Rabbit



Otter



**3-way is easier than 6-way**

Fox



Squirrel



Rabbit



Hamster



Otter



Beaver

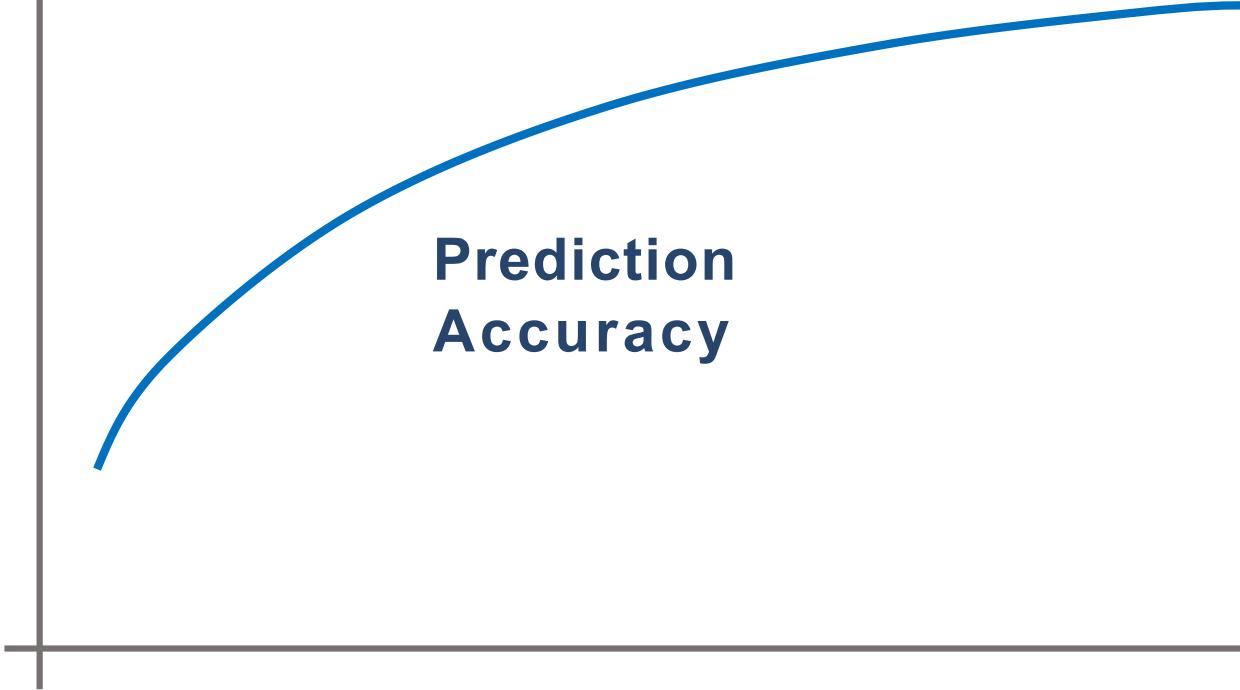


**Accuracy**



**Prediction  
Accuracy**

**# of shots**



Squirrel



Rabbit



Hamster



Otter



**2-shot is easier than 1-shot**

Squirrel



Rabbit



Hamster



Otter





**Idea: Learn a  
Similarity Function**

## Basic Idea

- Learn a similarity function:  $\text{sim}(\mathbf{x}, \mathbf{x}^*)$ .

## Basic Idea

- Learn a similarity function:  $\text{sim}(\mathbf{x}, \mathbf{x}^*)$ .
- Ideally,  $\text{sim}(\mathbf{x}_1, \mathbf{x}_2) = 1$ ,  $\text{sim}(\mathbf{x}_1, \mathbf{x}_3) = 0$ , and  $\text{sim}(\mathbf{x}_2, \mathbf{x}_3) = 0$ .

Bulldog



$\mathbf{x}_1$

Bulldog



$\mathbf{x}_2$

Fox



$\mathbf{x}_3$

- First, learn a similarity function from large-scale **training dataset**.



## Basic Idea

- First, learn a similarity function from large-scale training dataset.
- Then, apply the similarity function for prediction.
  - Compare the **query** with every sample in the **support set**.
  - Find the sample with the highest similarity score.

## Support Set:

Greyhound



Bulldog



Armadillo



Pangolin



Otter



Beaver



## Basic Idea

What is in the image?

Query:



Support Set:

Greyhound



Bulldog



Armadillo



Pangolin



Otter



Beaver



## Basic Idea

What is in the image?

Query:



sim = 0.2

Greyhound



Bulldog



Armadillo



Pangolin



Otter



Beaver



## Basic Idea

What is in the image?

Query:



sim = 0.2

sim = 0.1

Greyhound



Bulldog



Armadillo



Pangolin



Otter



Beaver



## Basic Idea

What is in the image?

Query:



sim = 0.2

sim = 0.1

sim = 0.03

Greyhound



Bulldog



Armadillo



Pangolin



Otter



Beaver



## Basic Idea

What is in the image?

Query:



sim = 0.2

sim = 0.1

sim = 0.03

sim = 0.05

sim = 0.7

sim = 0.5

Greyhound



Bulldog



Armadillo



Pangolin



Otter



Beaver



## Basic Idea

What is in the image?

Query:



sim = 0.2

sim = 0.1

sim = 0.03

sim = 0.05

sim = 0.7

sim = 0.5

Greyhound



Bulldog



Armadillo



Pangolin



Otter



Beaver



Restart and revisit

## Few-shot learning



- Have you seen before an **okapi**?
- Can you learn to recognize it from only this image?

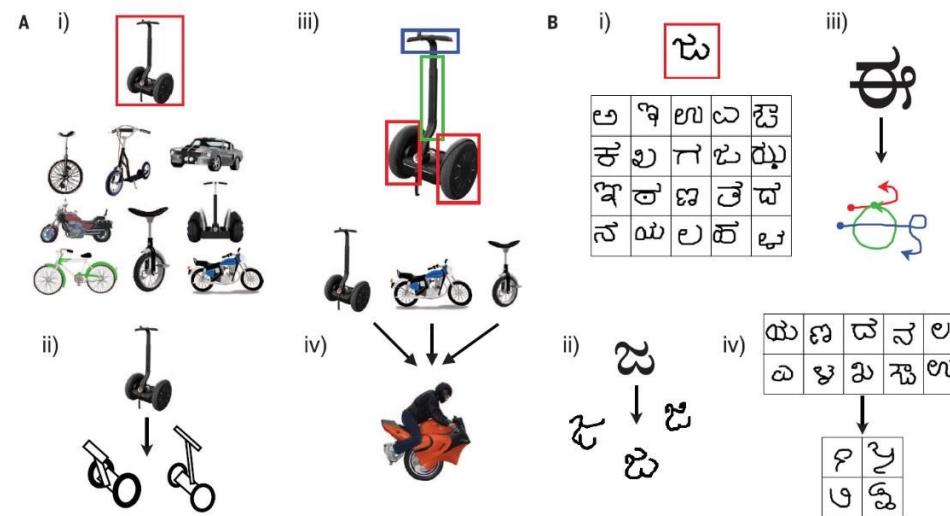
## Few-shot learning



- **Humans:** able to learn new concepts using few training examples
- **Goal of few-shot learning:** mimic this ability with machine learning methods

## Few-shot before the deep learning “revolution”

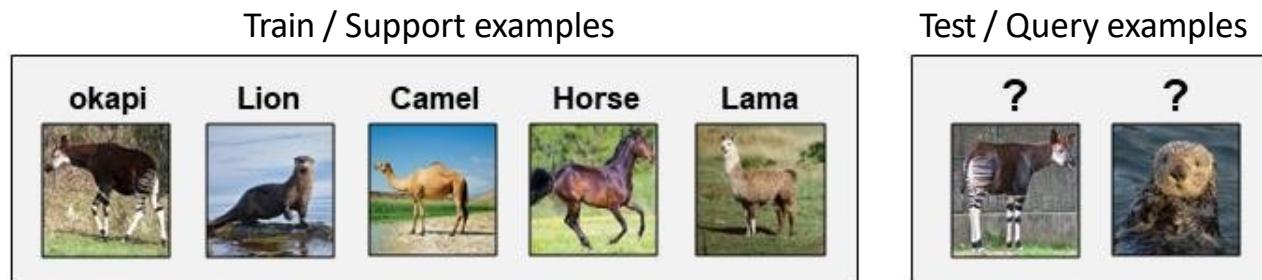
- “One-shot learning of simple visual concepts”, Lake et al. 11
- “One-Shot Learning with a Hierarchical Nonparametric Bayesian Model”, Salakhutdinov et al. 12
- “A Bayesian Approach to Unsupervised One-Shot Learning of Object Categories”, Fei Fei et al. 13
- “Human-level concept learning through probabilistic program induction”, Lake et al. 15



Here we will focus on deep learning based methods

## Formally: Learn N-way K-shot classification tasks

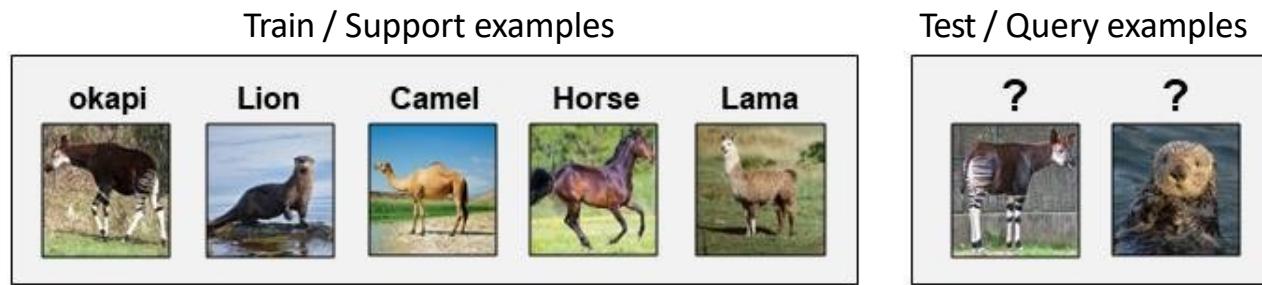
- **N** = number of classes
- **K** = training examples per class, **as small as 1 or 5!**



**Example: 5-way 1-shot classification task**

## Formally: Learn N-way K-shot classification tasks

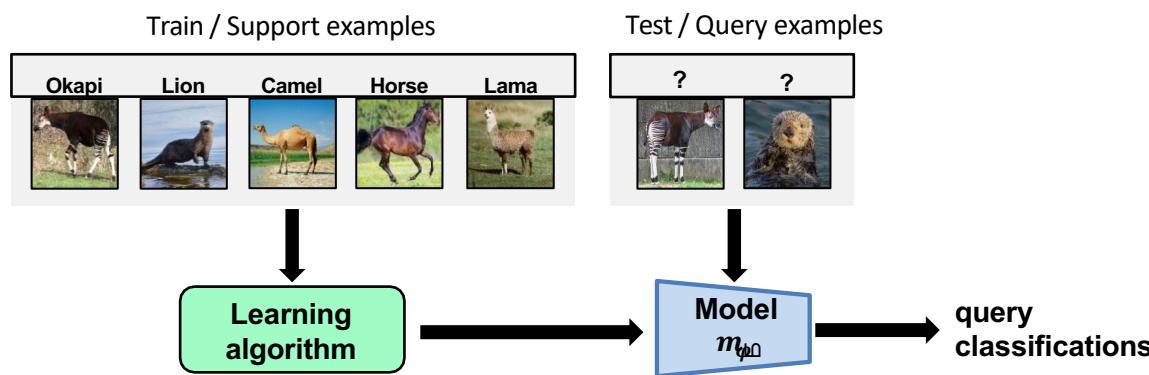
- **N** = number of classes
- **K** = training examples per class, **as small as 1 or 5!**



**Example:** 5-way 1-shot classification task

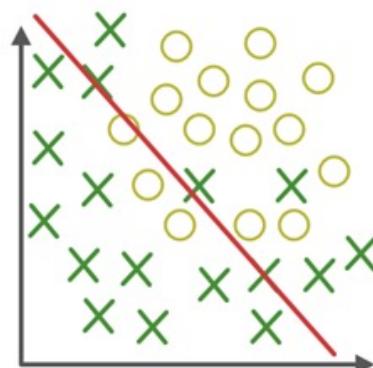
**Question:** is this possible with deep learning-based models?

## Train directly a deep learning model



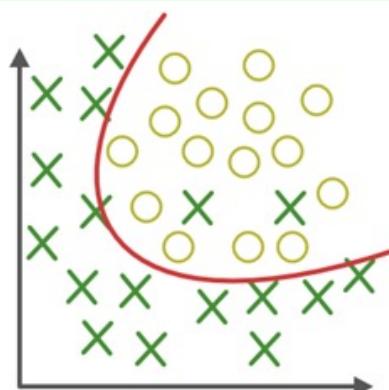
- Train from scratch a classification network
- Overfit to training data → poor accuracy on test data ☹

模型太簡單  
資料分布太複雜



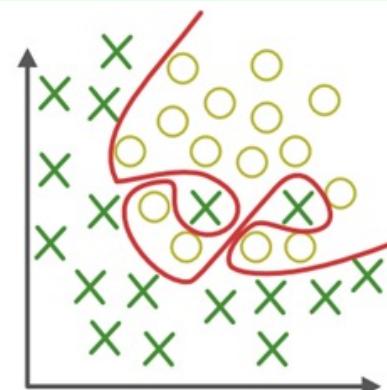
Under-fitting  
(too simple to  
explain the variance)

模型與資料  
擬合適當



Appropriate-fitting

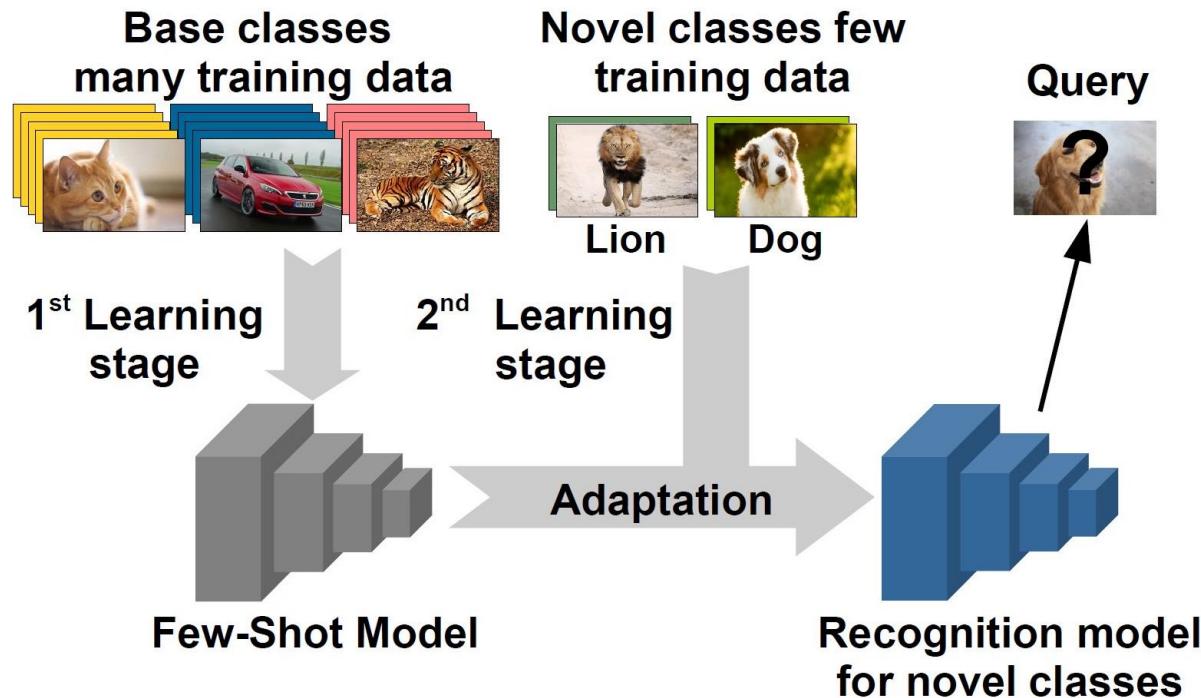
模型太複雜  
只能適用於訓練資料



Over-fitting  
(forcefitting--too  
good to be true) DG

— 模型  
○ 類別1  
× 類別2

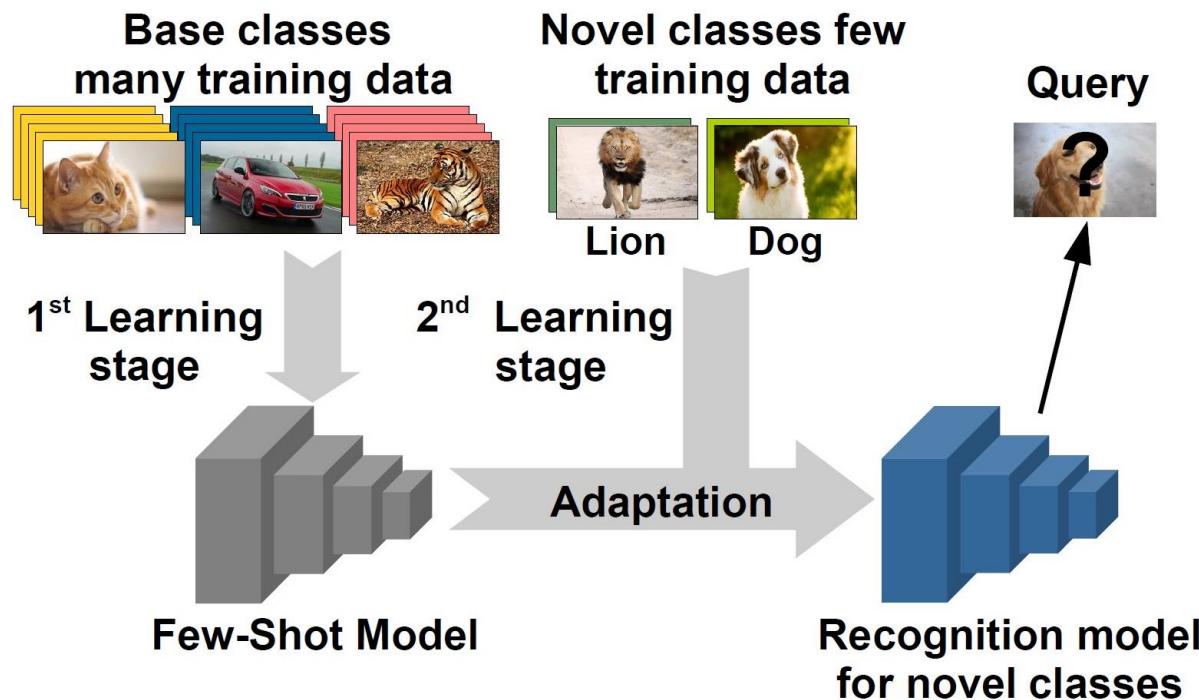
## Overcome data scarcity with transfer learning



### Overcome data scarcity with transfer learning

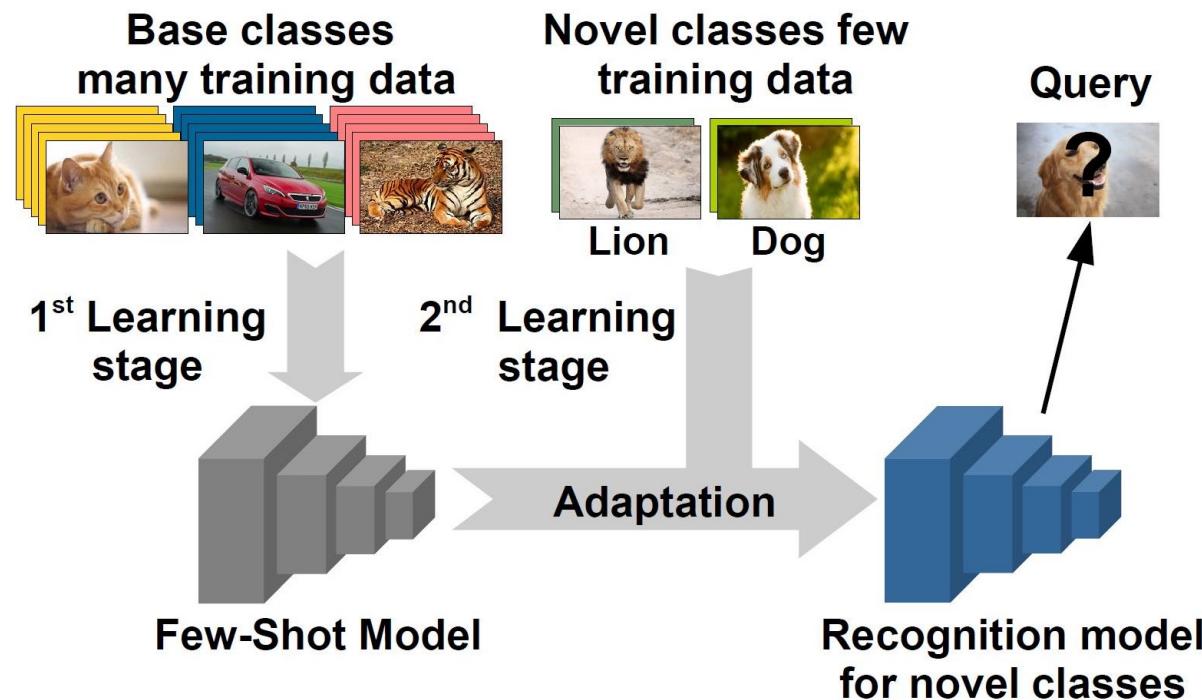
- Recipe followed by all few-shot learning methods

## Overcome data scarcity with transfer learning



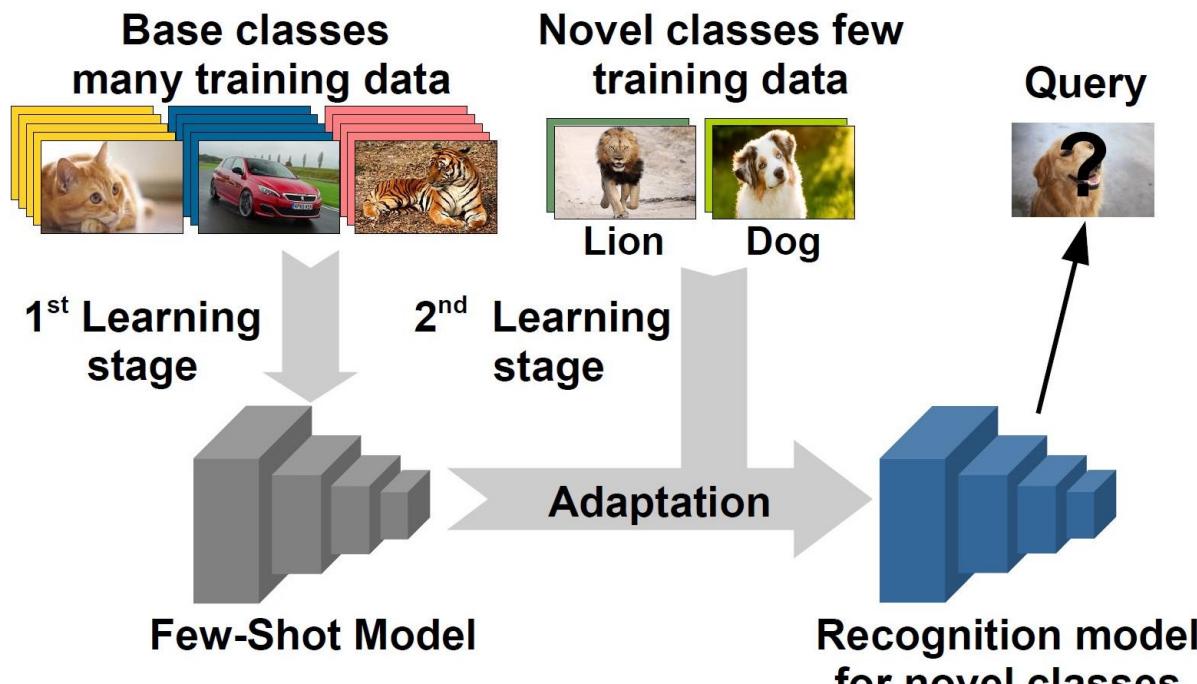
1. **Acquire knowledge:** train on other similar problems
2. **Transfer knowledge:** adapt to the problem of interest

## Overcome data scarcity with transfer learning



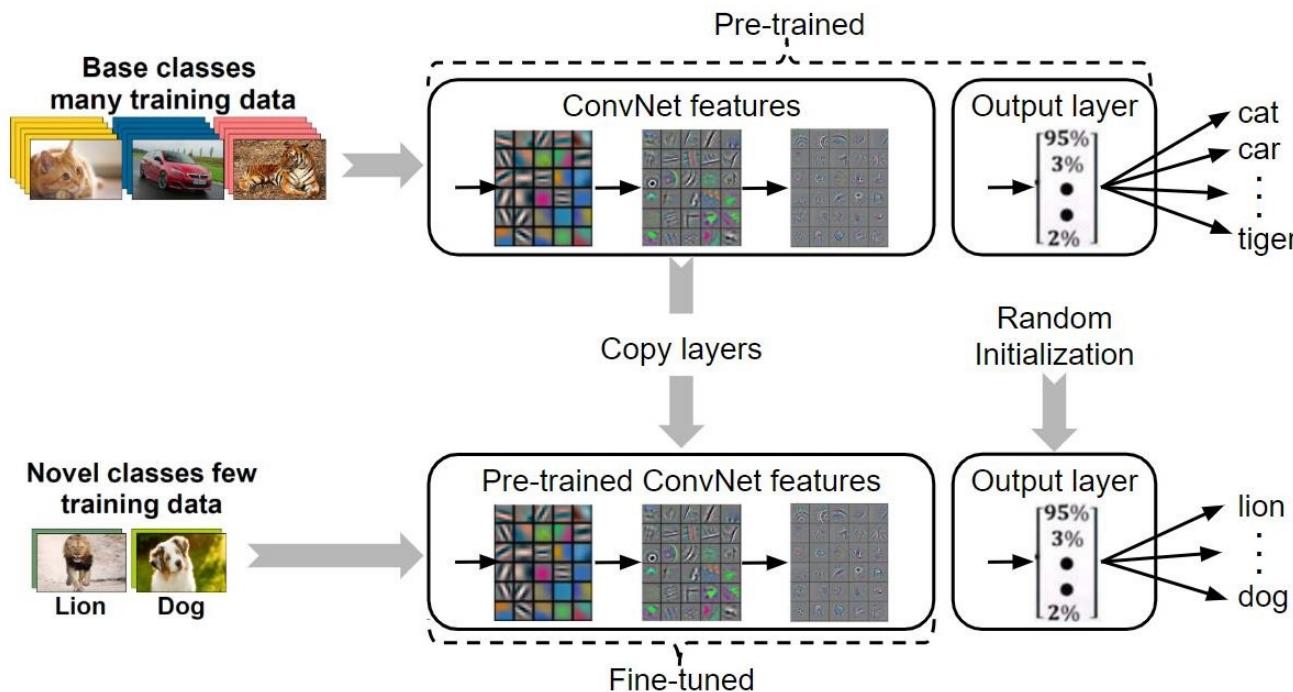
1. **Acquire knowledge:** use many training data from some **base classes**
2. **Transfer knowledge:** adapt to **novel classes** with few training data

## Overcome data scarcity with transfer learning



base classes == train classes  
novel classes == test classes } no overlap between them

## Common transfer learning example: Fine-tuning



1. **Acquire knowledge:** pre-train a network on the base class data
2. **Transfer knowledge:** fine-tune the network on novel class data

## Few-shot learning methods

**Fine-tuning:** **risk of overfitting** in case of extremely limited data (few-shot)

**Goal of few-shot learning:** devise transfer learning algorithms that would work well in the few-shot scenario, e.g., metric learning, meta-learning methods, ...

# Agenda

- Introduction
  - Few-shot learning problem
  - **Meta-learning paradigm**
  - How to evaluate
- Main types of few-shot learning algorithms
- Few-shot learning without forgetting

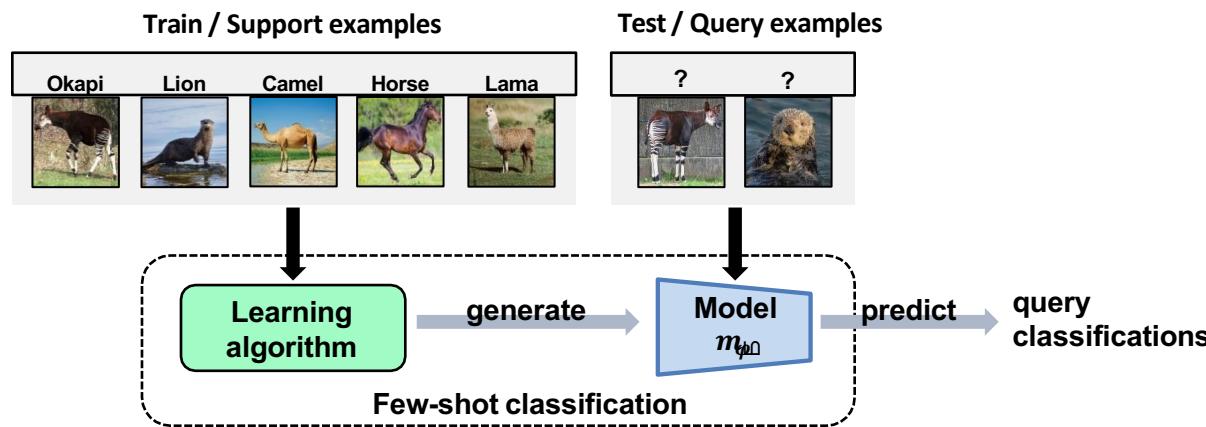
## Few-shot meta-learning

Most (but not all) few-shot methods use **meta-learning (learn-to-learn paradigm)**

- “Evolutionary principles in self-referential learning, or on learning how to learn”, Schmidhuber 1987
- “Meta-neural networks that learn by learning”, Naik et al. 1992
- “Lifelong learning algorithms”, Thrun 1998
- “Learning to learn by gradient descent by gradient descent”, Andrychowicz et al. 16
- ...

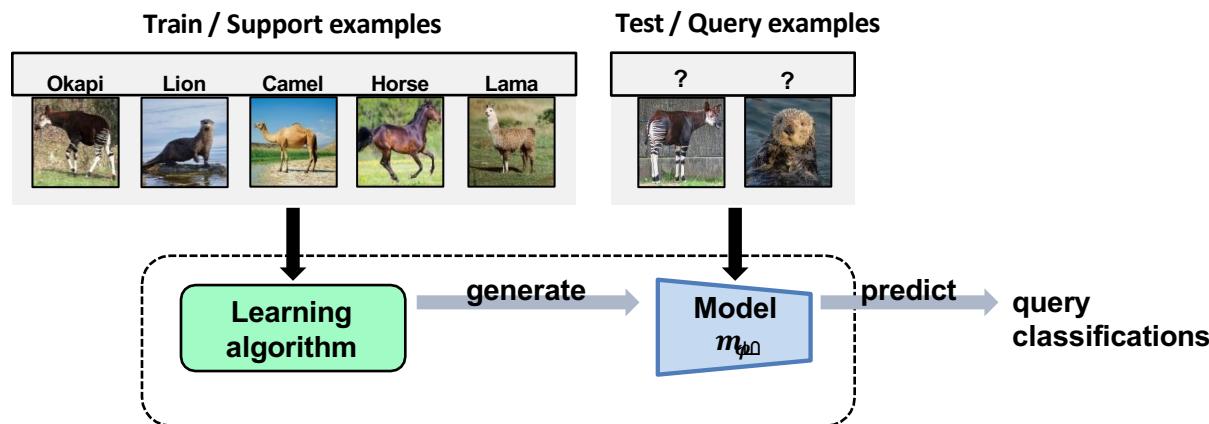
**What is few-shot meta-learning?**

# Few-shot classification



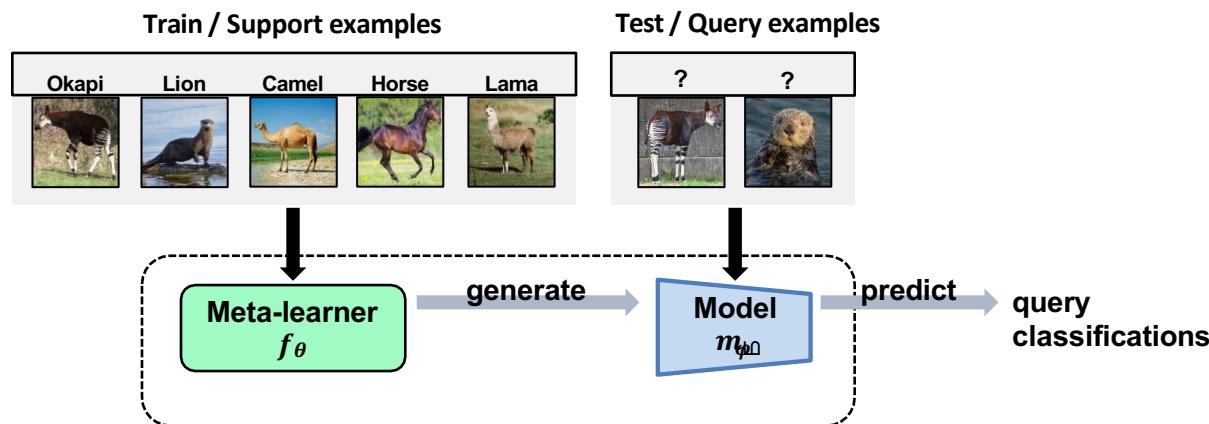
- **input:** labeled support data, unlabeled query data
- **intermediate output:** model for classifying the query images
- **output:** predicted query labels

# Few-shot classification with meta-learning



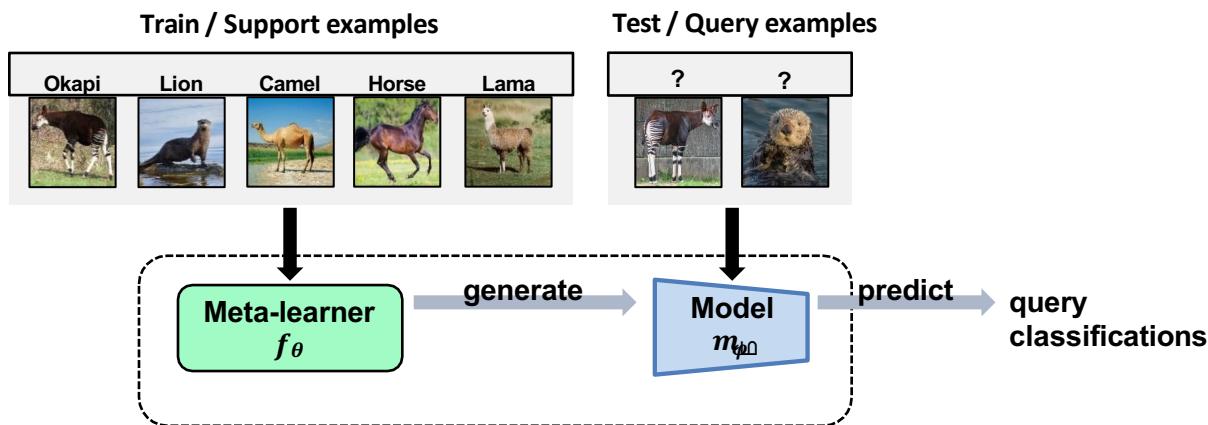
- **Train the learning algorithm** (instead of the classification model)
  - Implement it with a **meta-learner**  $f_{\theta}$
  - Optimize  $f_{\theta}$  on learning few-shot classification tasks (**learn-to-learn**)

## Few-shot classification with meta-learning



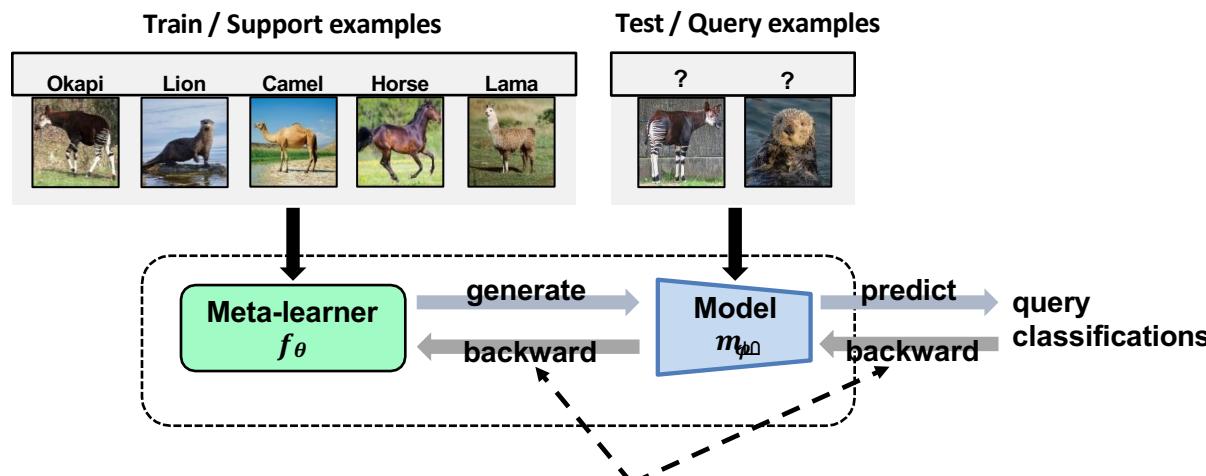
- **Train the learning algorithm** (instead of the classification model)
  - Implement it with a **meta-learner**  $f_\theta$  (**somewhat**)
  - Optimize  $f_\theta$  on solving few-shot classification tasks (**learn-to-learn**)

# Few-shot classification with meta-learning



- **Train the learning algorithm** (instead of the classification model)
  - Implement it with a **meta-learner  $f_\theta$**
  - Optimize  $f_\theta$  on solving few-shot classification tasks (**learn-to-learn**)

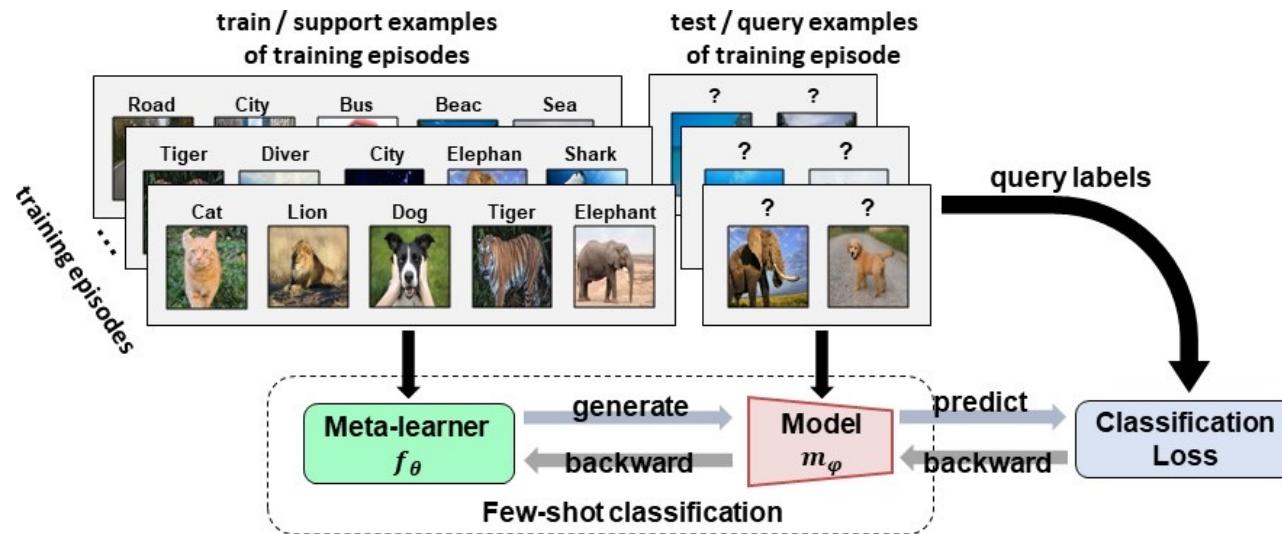
## Few-shot classification with meta-learning



must back-propagate the entire few-shot learning process

- **Train the learning algorithm** (instead of the classification model)
  - Implement it with a **meta-learner  $f_\theta$**
  - Optimize  $f_\theta$  on solving few-shot classification tasks (**learn-to-learn**)

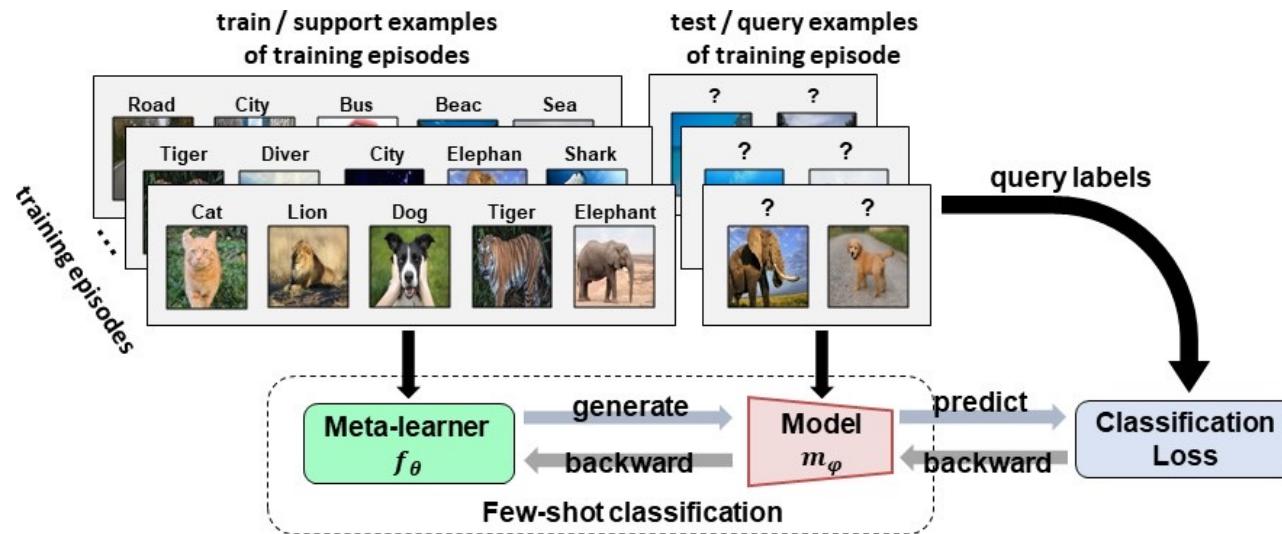
## Meta-learning: training time (1<sup>st</sup> learning stage)



### How to train the meta-learner?

- Train it on the same conditions it will be used in 2<sup>nd</sup> learning stage (meta-test)

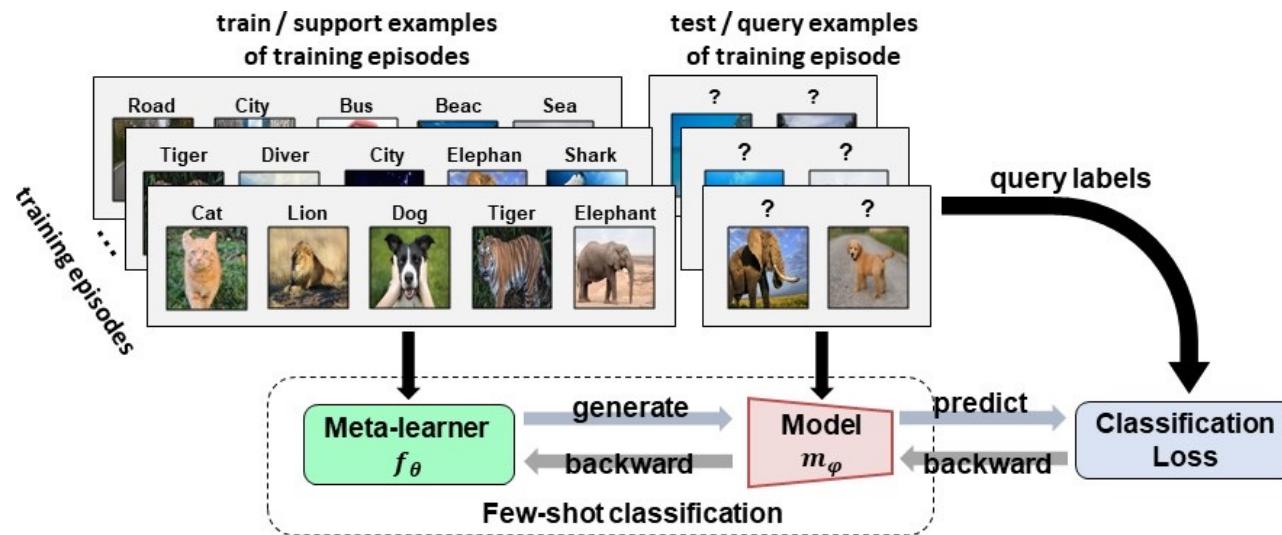
## Meta-learning: training time (1<sup>st</sup> learning stage)



### How to train the meta-learner?

- Train meta-learner  $f_\theta$  on solving a distribution of few-shot tasks (aka episodes)

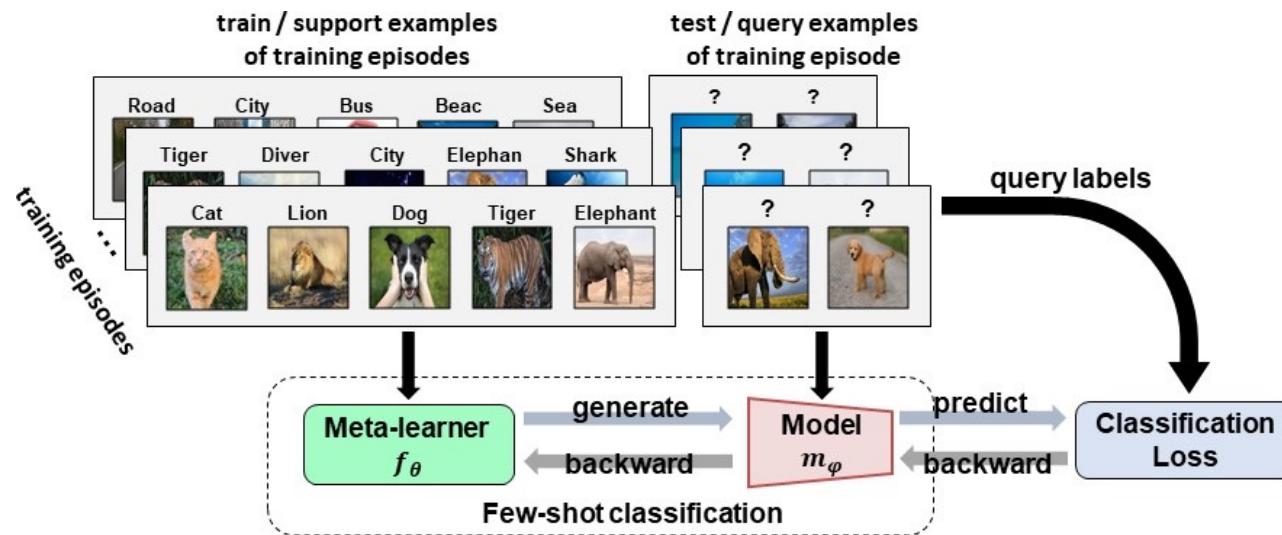
## Meta-learning: training time (1<sup>st</sup> learning stage)



### How to train the meta-learner?

- Train **meta-learner  $f_\theta$**  on solving a distribution of few-shot tasks (aka **episodes**)
- Construct such **training episodes** using the base class data

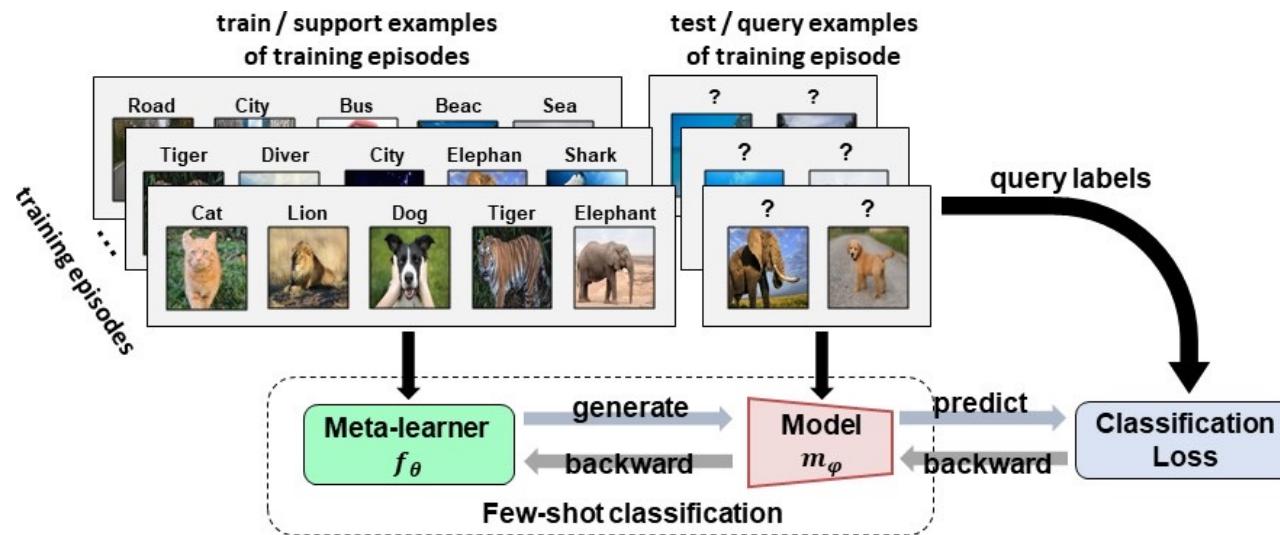
## Meta-learning: training time (1<sup>st</sup> learning stage)



### How to train the meta-learner?

- Train **meta-learner  $f_\theta$**  on solving a distribution of few-shot tasks (aka **episodes**)
- Construct such **training episodes** using the base class data
  - by sampling **N** classes x ( **K** support examples + **M** query examples)

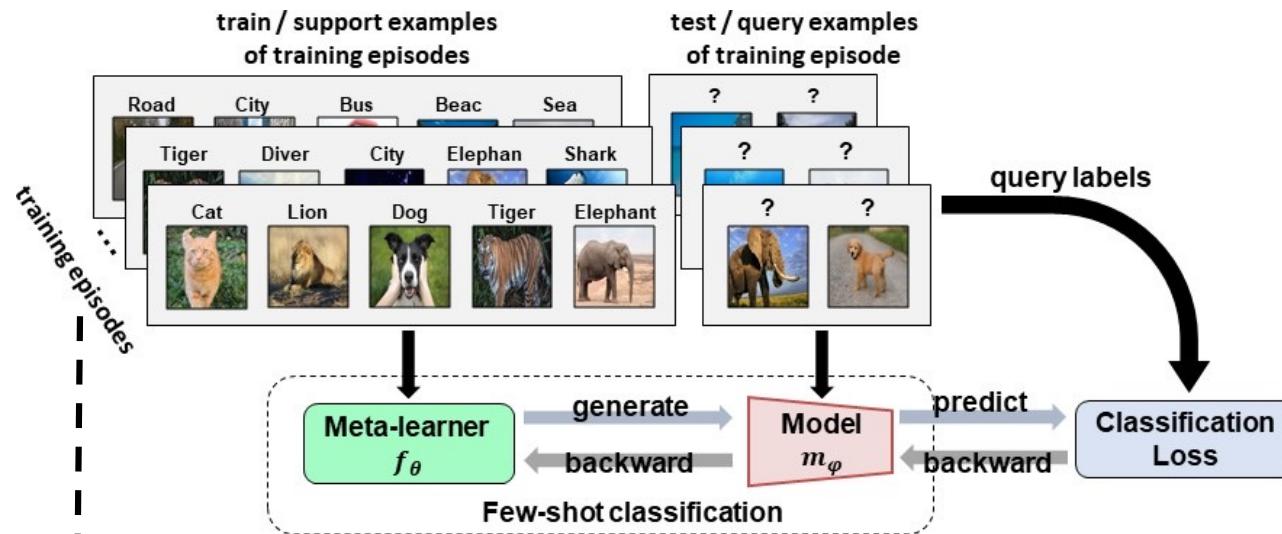
## Meta-learning: training time (1<sup>st</sup> learning stage)



**Objective:**

$$\min_{\theta} \mathbb{E}_{(S,Q)} L(f_\theta(S), Q)$$

## Meta-learning: training time (1<sup>st</sup> learning stage)

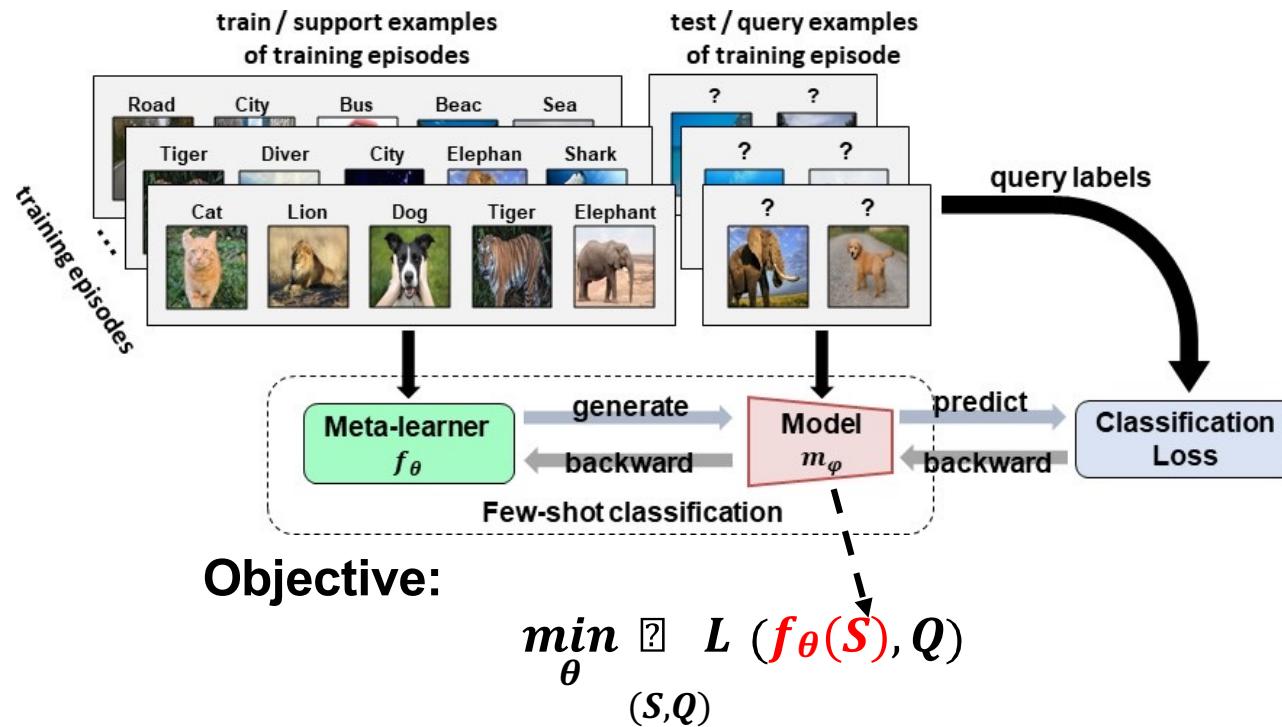


**Objective:**

$$\min_{\theta} \mathbb{E}_{(S,Q)} L(f_\theta(S), Q)$$

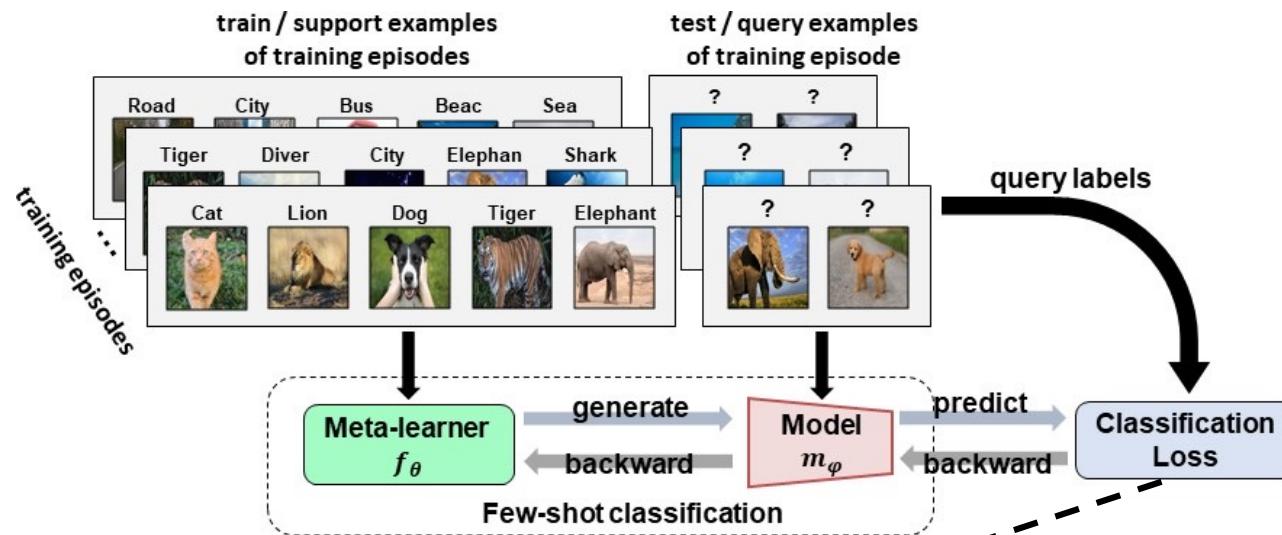
**Episode ( $S, Q$ ):** support set  $S = \{x_k^S, y_k^S\}_{k=1}^{N*K}$  and query set  $Q = \{x_m^Q, y_m^Q\}_{m=1}^{N*M}$

## Meta-learning: training time (1<sup>st</sup> learning stage)



**Inner part:** generate using the support set  $S$  the classification model  $m_\varphi = f_\theta(S)$

## Meta-learning: training time (1<sup>st</sup> learning stage)

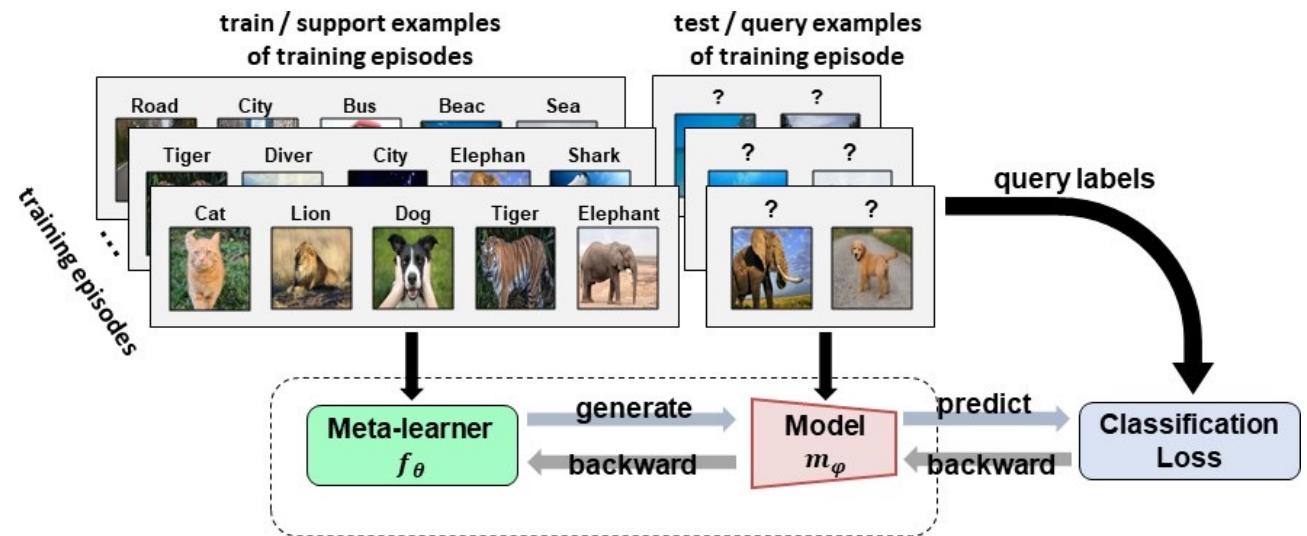


**Objective:**

$$\min_{\theta} \mathbb{E}_{(S, Q)} L(f_\theta(S), Q)$$

**Outer part:** optimize  $\theta$  w.r.t. the queries classification loss  $L(f_\theta(S), Q) = L(m_\varphi, Q)$

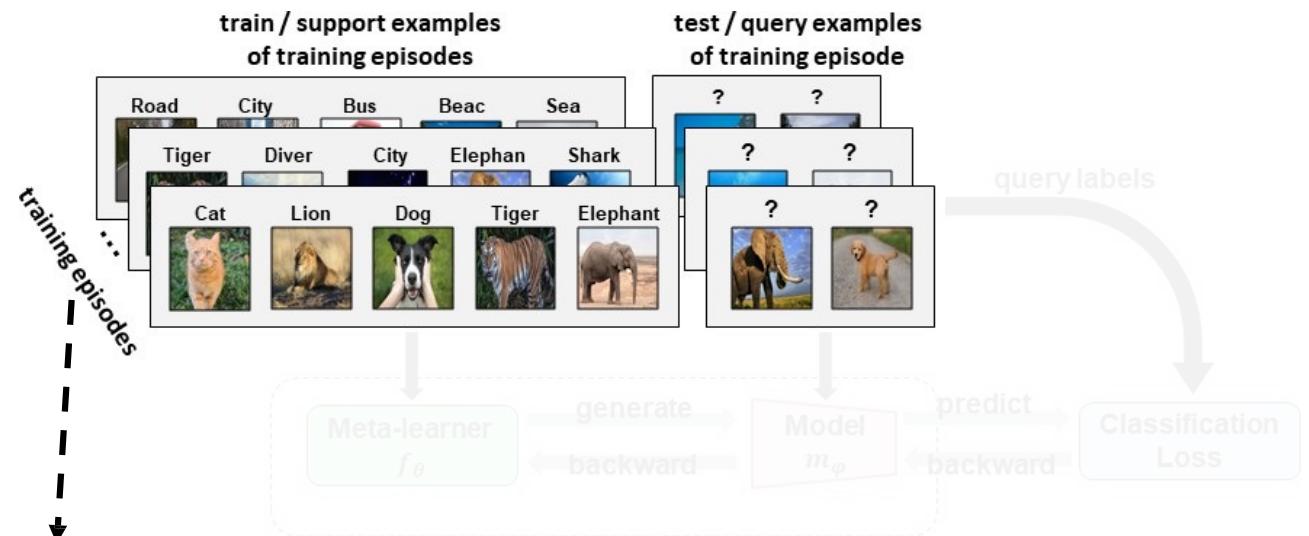
## Meta-learning: training time (1<sup>st</sup> learning stage)



### Meta-training routine:

1. Sample training episode  $(S, Q)$
2. Generate classification model  $m_\varphi = f_\theta(S)$
3. Predict classification scores  $p_m = m_\varphi(x_m^Q)$  for each  $x_m^Q$  in  $Q$
4. Optimize  $\theta$  w.r.t. the queries classification loss  $L(f_\theta(S), Q)$

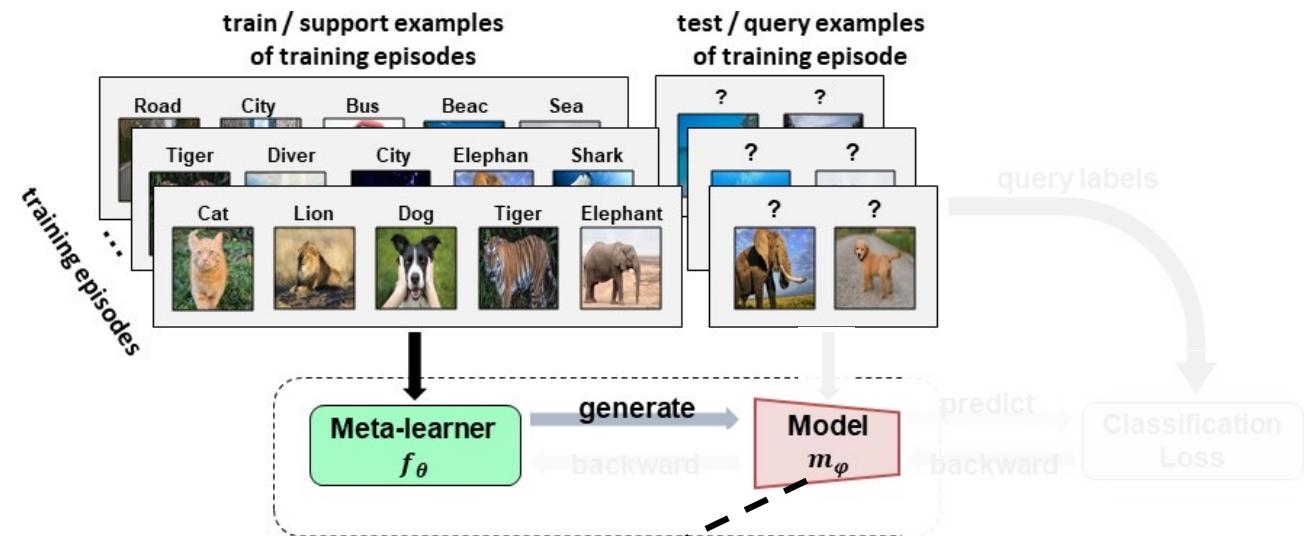
## Meta-learning: training time (1<sup>st</sup> learning stage)



**Meta-training routine:**

1. Sample training episode  $(S, Q)$
2. Generate classification model  $m_\varphi = f_\theta(S)$
3. Predict classification scores  $p_m = m_\varphi(x_m^Q)$  for each  $x_m^Q$  in  $Q$
4. Optimize  $\theta$  w.r.t. the queries classification loss  $L(f_\theta(S), Q)$

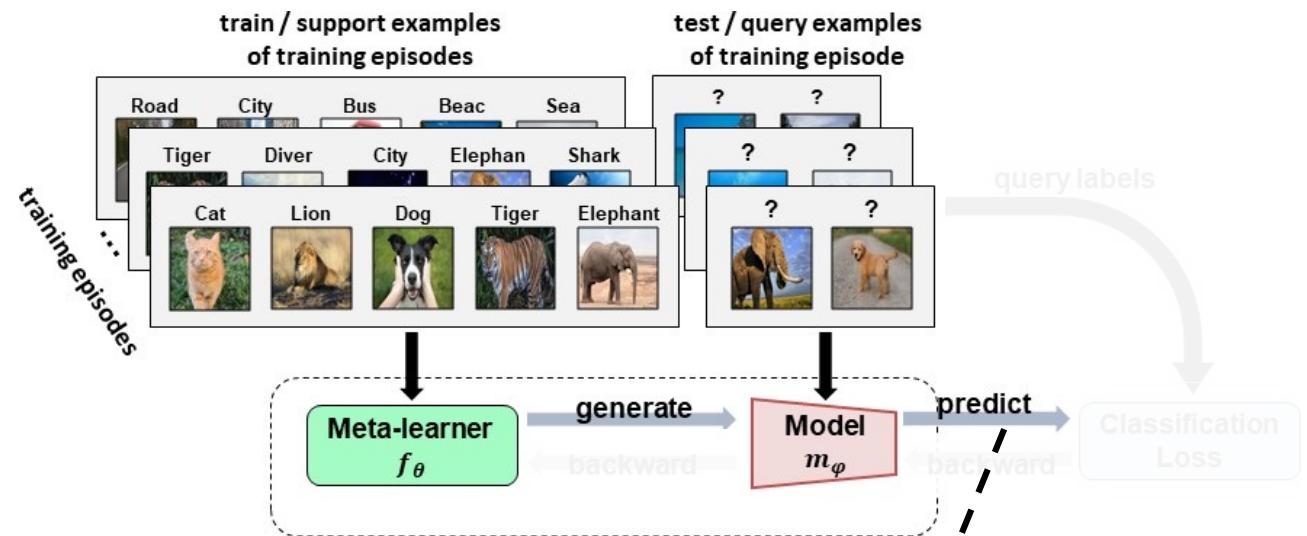
## Meta-learning: training time (1<sup>st</sup> learning stage)



### Meta-training routine:

1. Sample training episode  $(S, Q)$
2. **Generate classification model**  $m_\varphi = f_\theta(S)$
3. Predict classification scores  $p_m = m_\varphi(x_m^Q)$  for each  $x_m^Q$  in  $Q$
4. Optimize  $\theta$  w.r.t. the queries classification loss  $L(f_\theta(S), Q)$

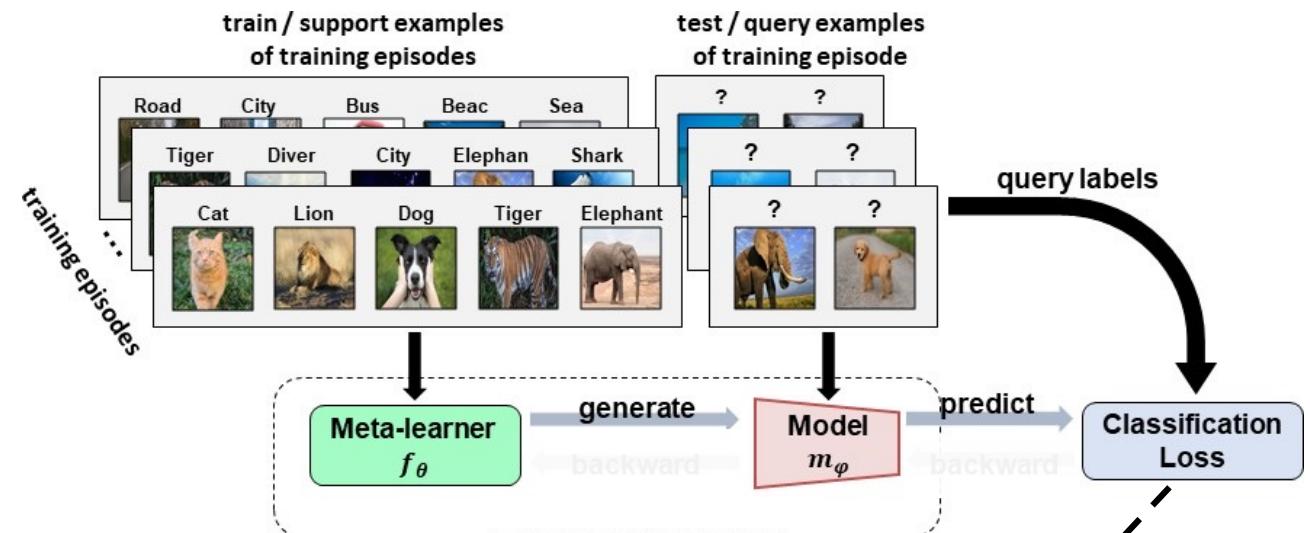
## Meta-learning: training time (1<sup>st</sup> learning stage)



### Meta-training routine:

1. Sample training episode  $(S, Q)$
2. Generate classification model  $m_\varphi = f_\theta(S)$
3. **Predict classification scores**  $p_m = m_\varphi(x_m^Q)$  for each  $x_m^Q$  in  $Q$
4. Optimize  $\theta$  w.r.t. the queries classification loss  $L(f_\theta(S), Q)$

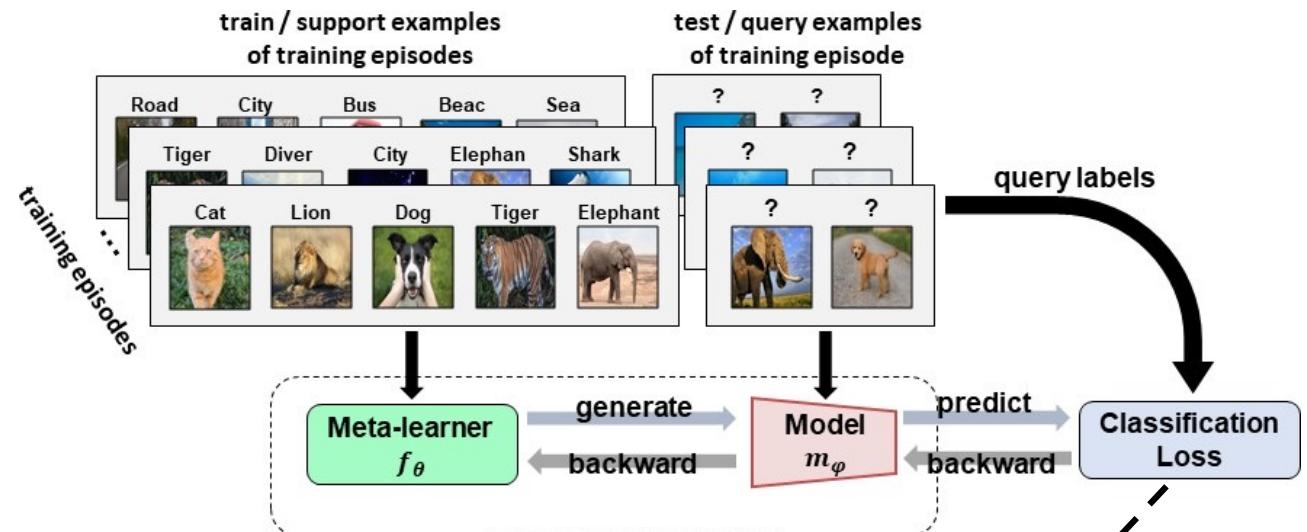
## Meta-learning: training time (1<sup>st</sup> learning stage)



### Meta-training routine:

1. Sample training episode  $(S, Q)$
2. Generate classification model  $m_\varphi = f_\theta(S)$
3. Predict classification scores  $p_m = m_\varphi(x_m^Q)$  for each  $x_m^Q$  in  $Q$
4. Optimize  $\theta$  w.r.t. the queries classification loss  $L(f_\theta(S), Q)$ 
  - e.g., cross entropy loss  $\sigma_m - \log(p_m[y_m^Q])$

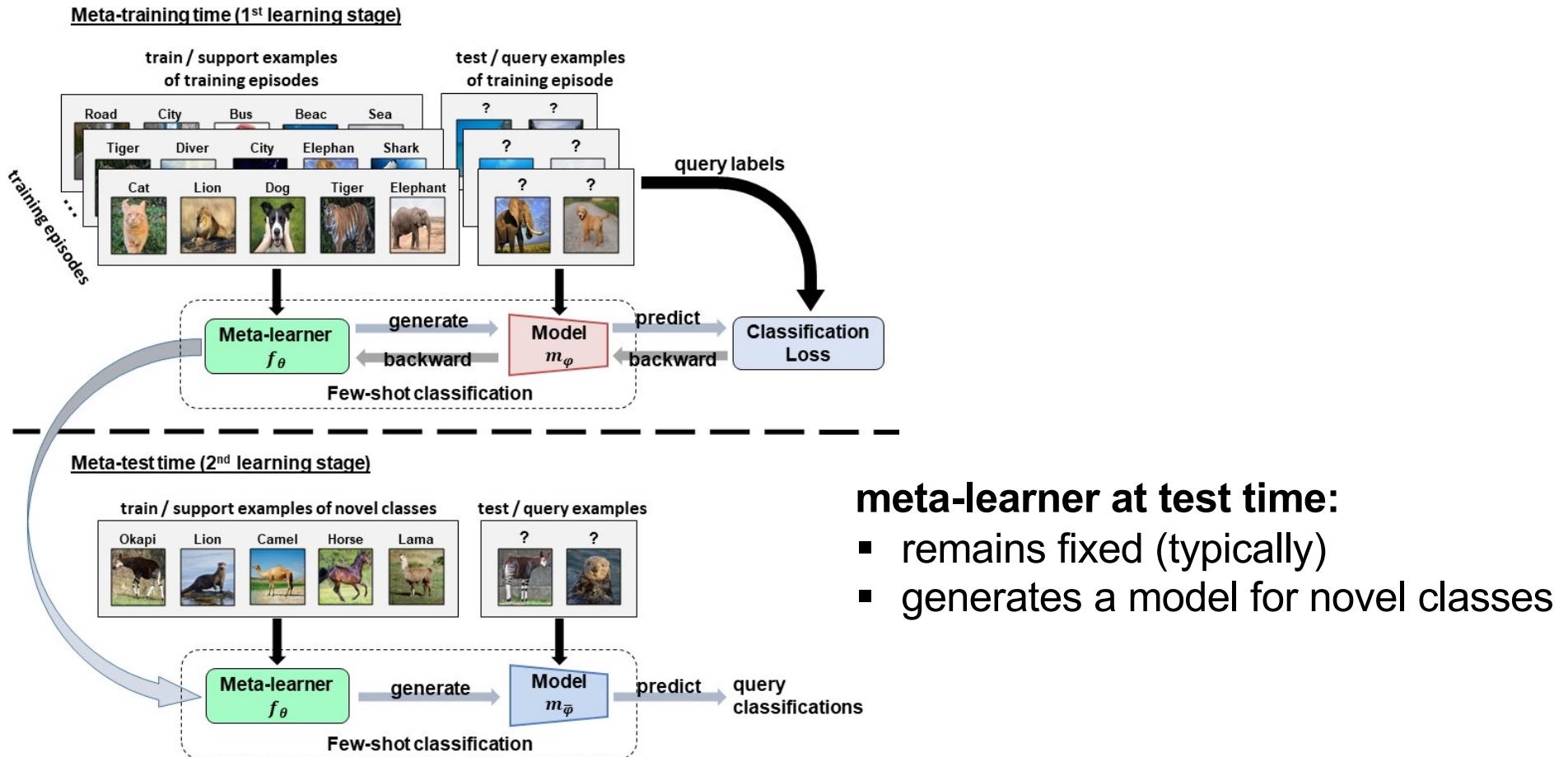
## Meta-learning: training time (1<sup>st</sup> learning stage)



### Meta-training routine:

1. Sample training episode  $(S, Q)$
2. Generate classification model  $m_\varphi = f_\theta(S)$
3. Predict classification scores  $p_m = m_\varphi(x_m^Q)$  for each  $x_m^Q$  in  $Q$
4. **Optimize  $\theta$  w.r.t. the queries classification loss  $L(f_\theta(S), Q)$** 
  - must back-propagate through the few-shot learning process

# Meta-learning: test time (2<sup>nd</sup> learning stage)



## From **Supervised Learning** to **Meta-Learning**

- **training** → meta-training
- **test time** → meta-test time
- **mini-batch of images** → mini-batch of few-shot episodes
- **training data** → meta-training data = all possible training episodes
- **test data** → meta-test data = test episodes

## Few-shot learning vs Meta-learning

### **Few-shot learning:**

- Any transfer learning method that targets on transferring well with limited data
- E.g.: pre-train + fine-tuning, or using metric learning, or using meta-learning

### **Meta-learning:**

- Learn the learning algorithm itself
  - “Learning to learn by gradient descent by gradient descent”, Andrychowicz et al. 16
- Ingredient of many few-shot algorithms,
- Also used in multi-task learning, RL, ...



Details of meta  
learning

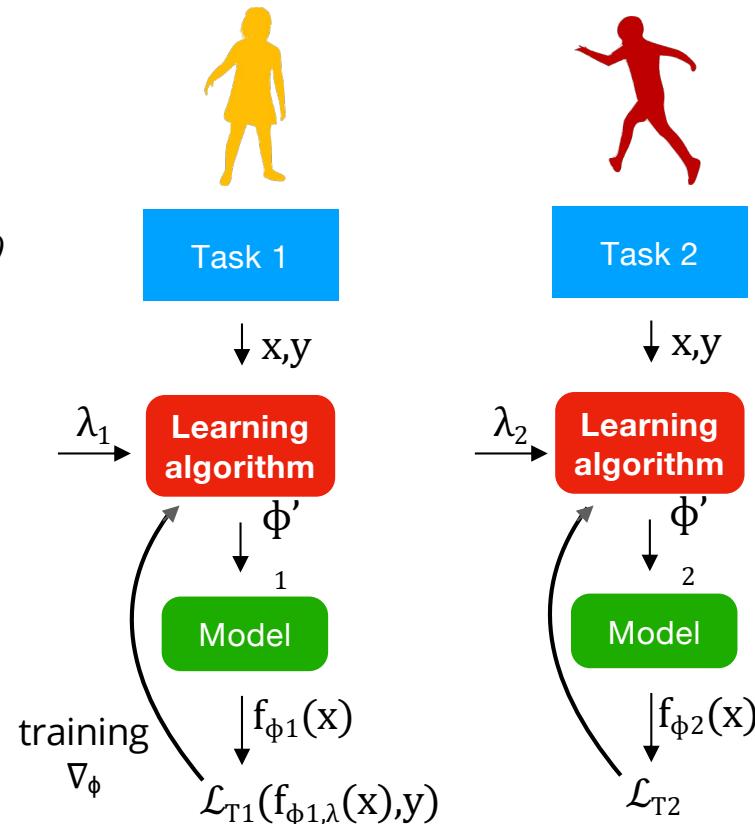
# Machine Learning

**Task:** distribution of samples  $q(x)$  outputs  $y$ , loss  $\mathcal{L}(x,y)$

**Learner:** model parameters  $\phi$ , hyper-parameters  $\lambda$  optimizer

$$\begin{aligned}\phi^* &= \operatorname{argmax}_\phi \log p(\phi | T) \\ &= \operatorname{argmin}_\phi \mathcal{L}(f_{\phi,\lambda}(x), y)\end{aligned}$$

**Model:**  $f_\phi(x) = y'$



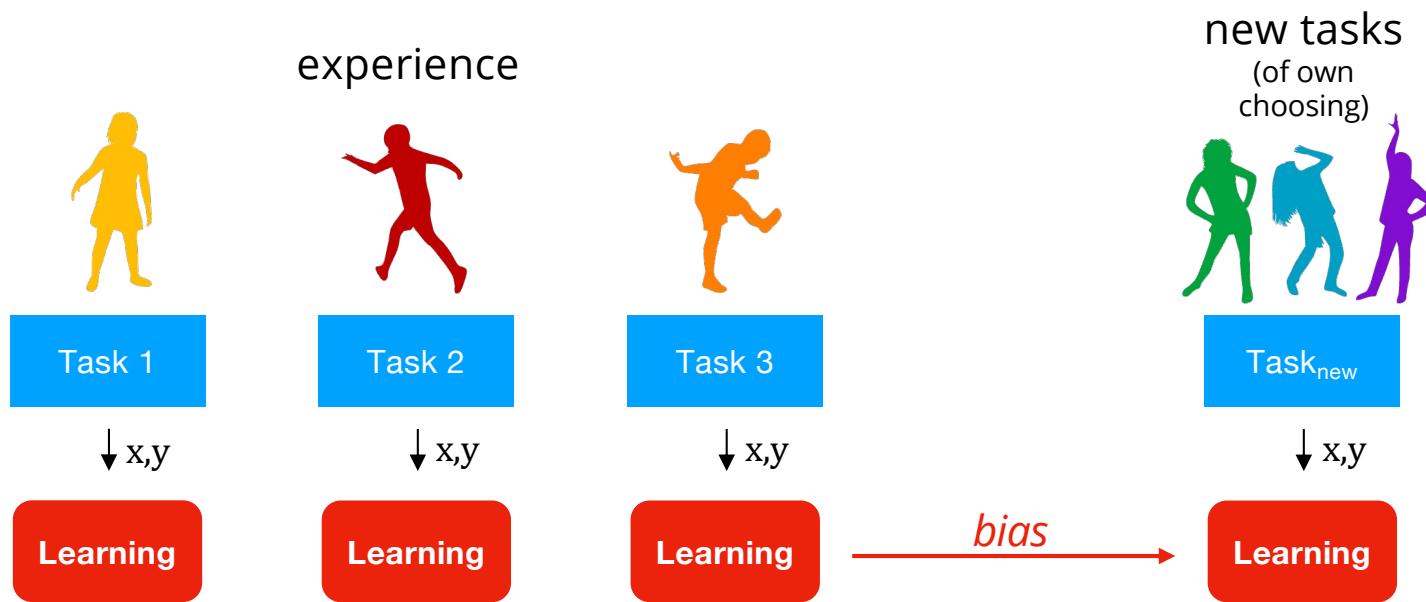
*for reinforcement learning:*  
 $q(x_t) + \text{transition } q(x_{t+1} | x_t, y_t)$   
 $\mathcal{L} = (\text{neg}) \text{ reward}$

Task 2 starts from scratch  
*(tabula rasa)*  
 ↓  
 unrealistic data and compute requirements

# Human-like Learning\*\*\*

Lake et al. 2017  
Kemp et al. 2007

learn *across* tasks: less trial-and-error, less data, less compute  
*tasks should be related to experience (doable, fun, interesting?)*

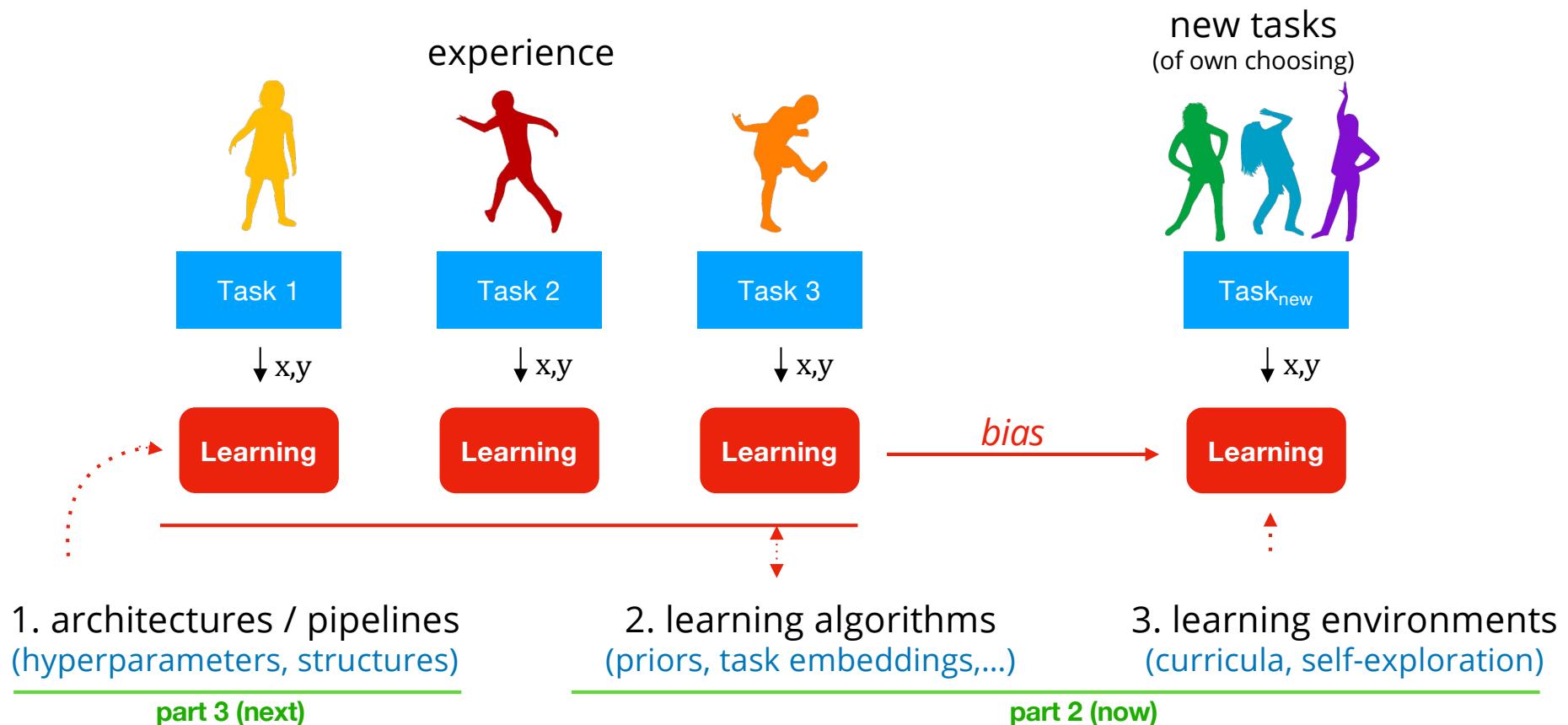


*key aspects of fast learning: compositionality, causality, learning to learn*

# What can we learn to learn?

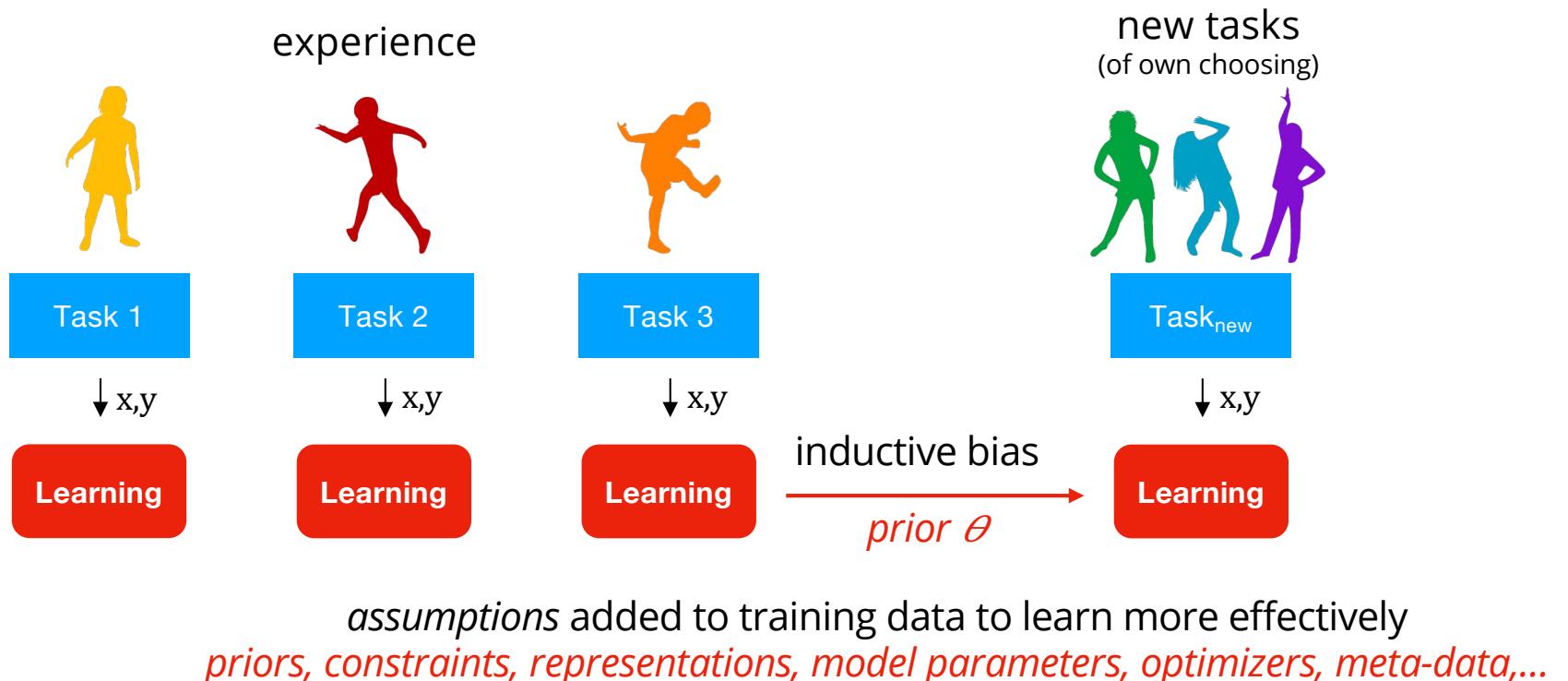
[Clune 2019](#)

*From hand-designed to learned learning algorithms ... to AI-generating algorithms?*



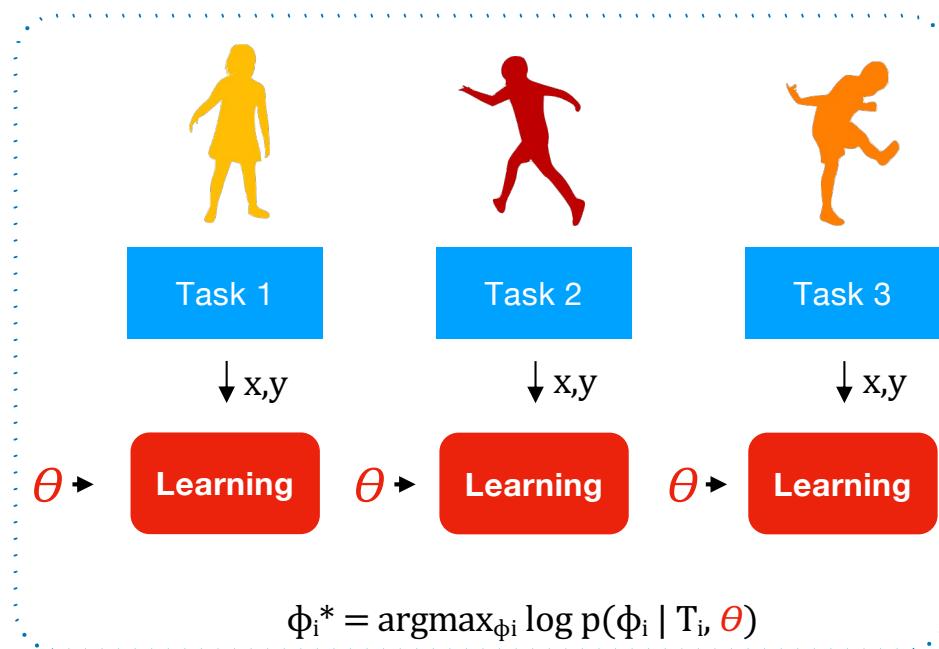
# Meta-learning (learning to learn)

learn *minimal* inductive biases instead of constructing *manual* ones  
should still generalize well (otherwise you meta-overfit)



# Strategy 1: bilevel optimization

*parameterize some aspect of the learner that we want to learn as meta-parameters  $\theta$  meta-learn  $\theta$  across tasks*



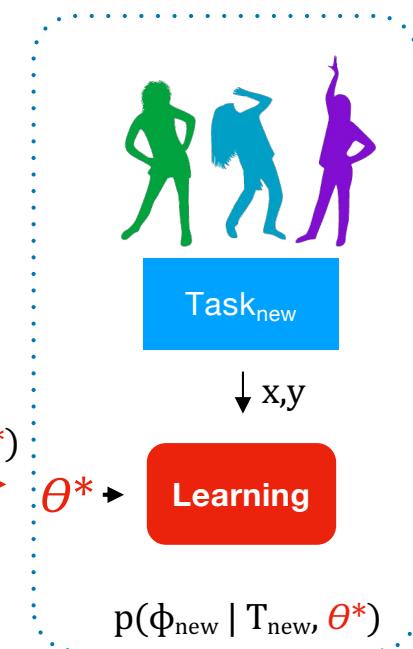
**Bilevel optimization:**

$$\theta^* = \text{argmax}_\theta \log p(\theta | T_i)$$

so that

$$\phi_i^* = \text{argmax}_\phi \log p(\phi | T_i, \theta^*)$$

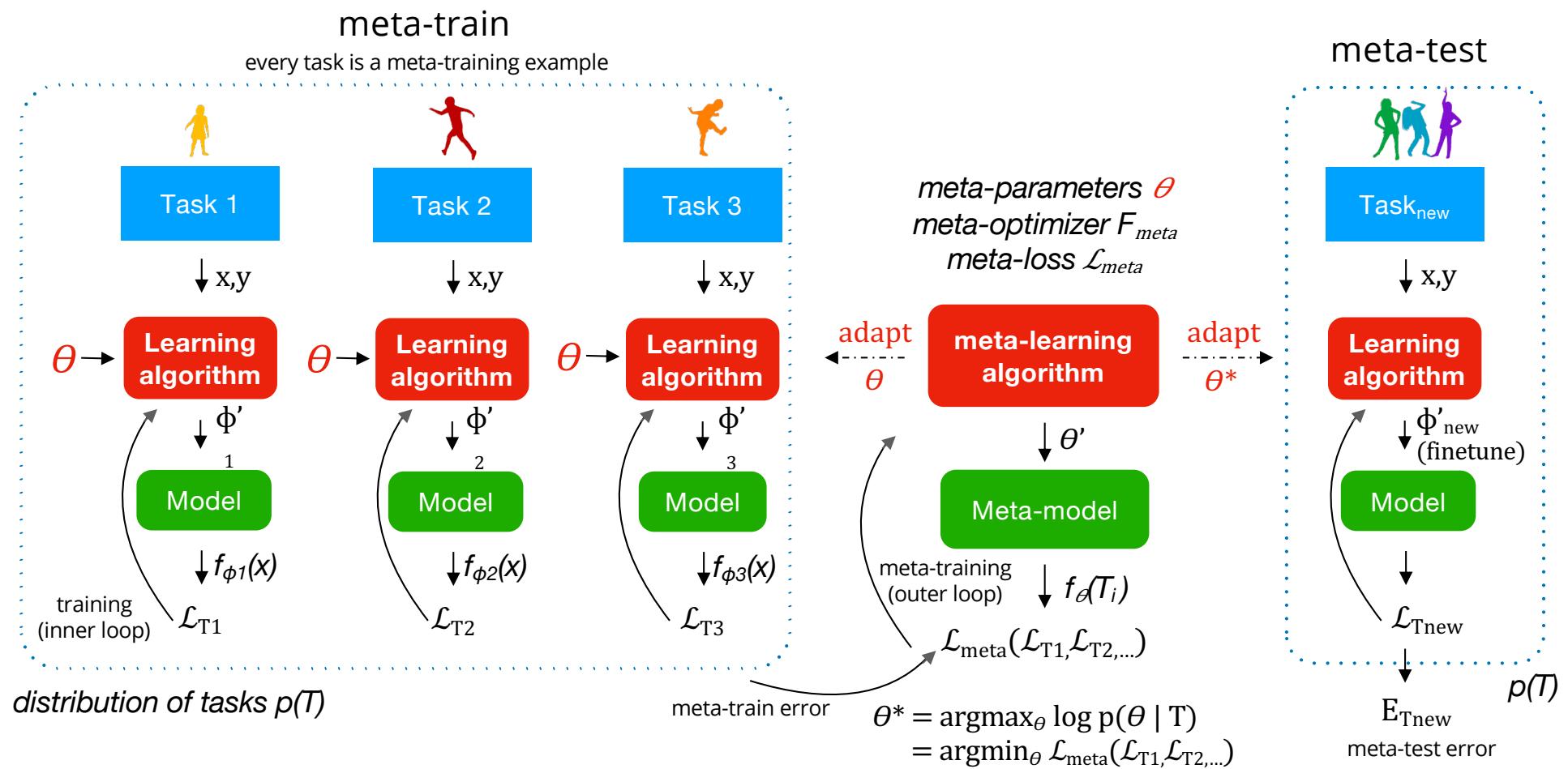
allows fine-tuning on  $T_{\text{new}}$



$\theta$  (prior), could encode an initialization  $\phi$ , the hyperparameters  $\lambda$ , the optimizer,...

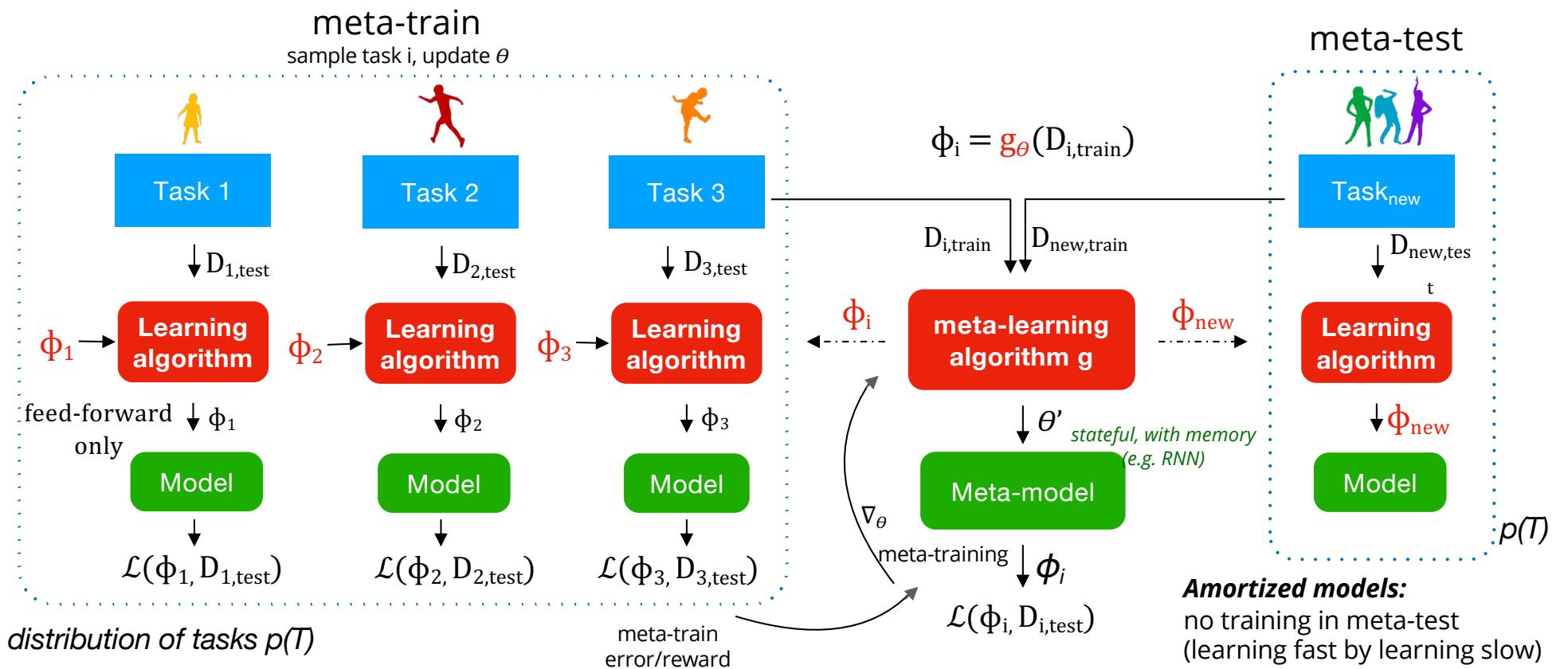
Learned  $\theta^*$  should *learn*  $T_{\text{new}}$  from small amount of data, yet *generalize* to a large number of tasks

# Meta-learning with bilevel optimization



# Strategy 2: black-box models

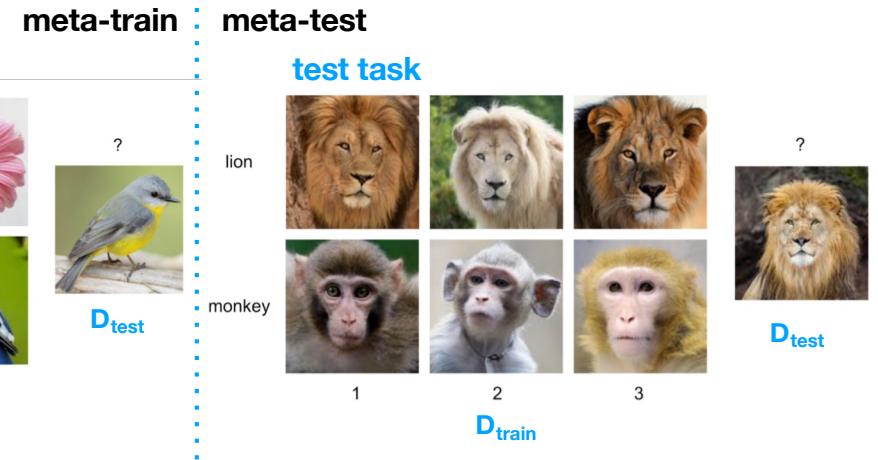
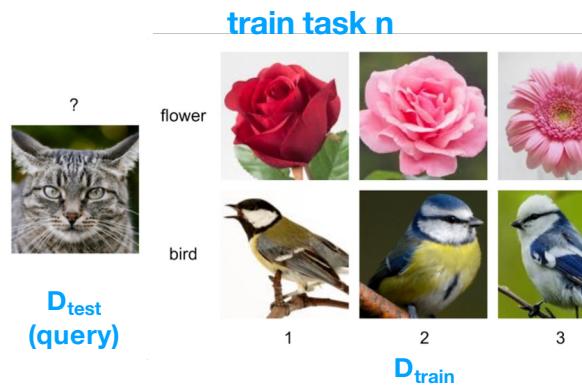
*black box meta-model  $g_\theta$ : predicts  $\phi$  given  $D_{train}$  ( $\theta$  is hidden)  
 hypernetwork where input embedding learned across tasks*



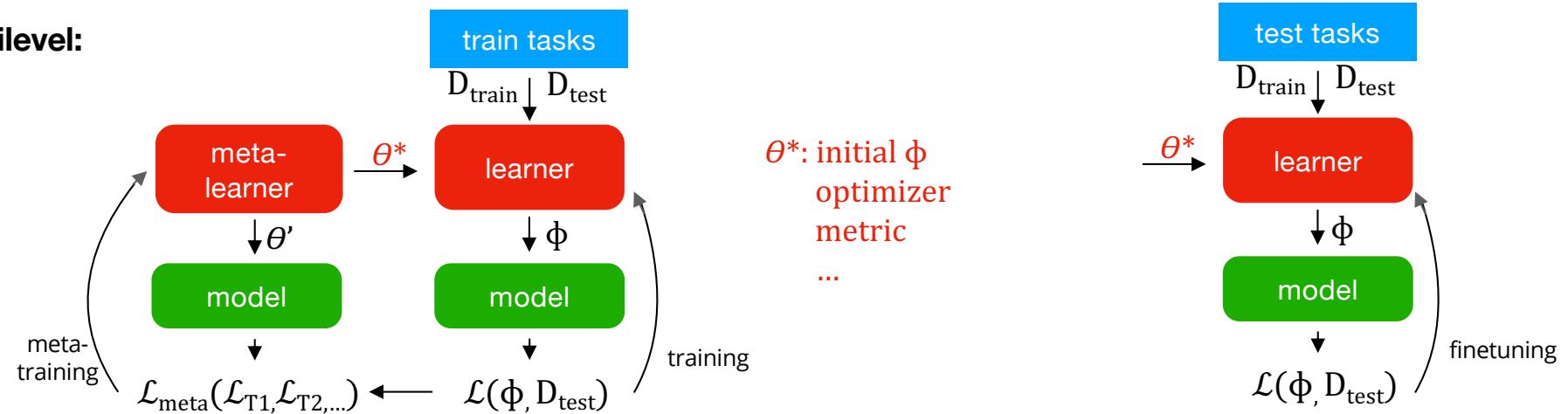
# Example: few-shot classification

Vinyals et al. 2016  
Triantafillou et al. 2019

2-way, 3-shot



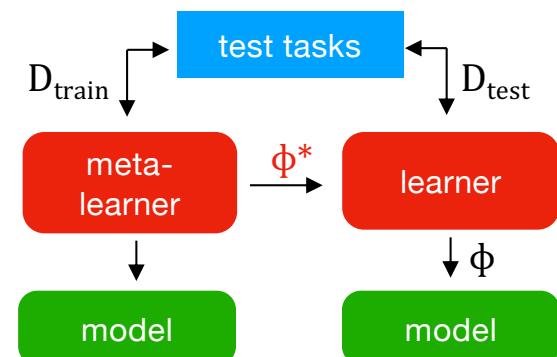
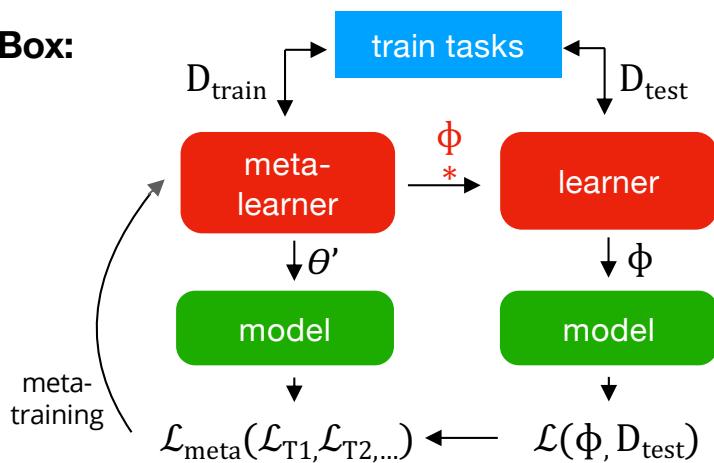
**Bilevel:**



# Example: few-shot classification



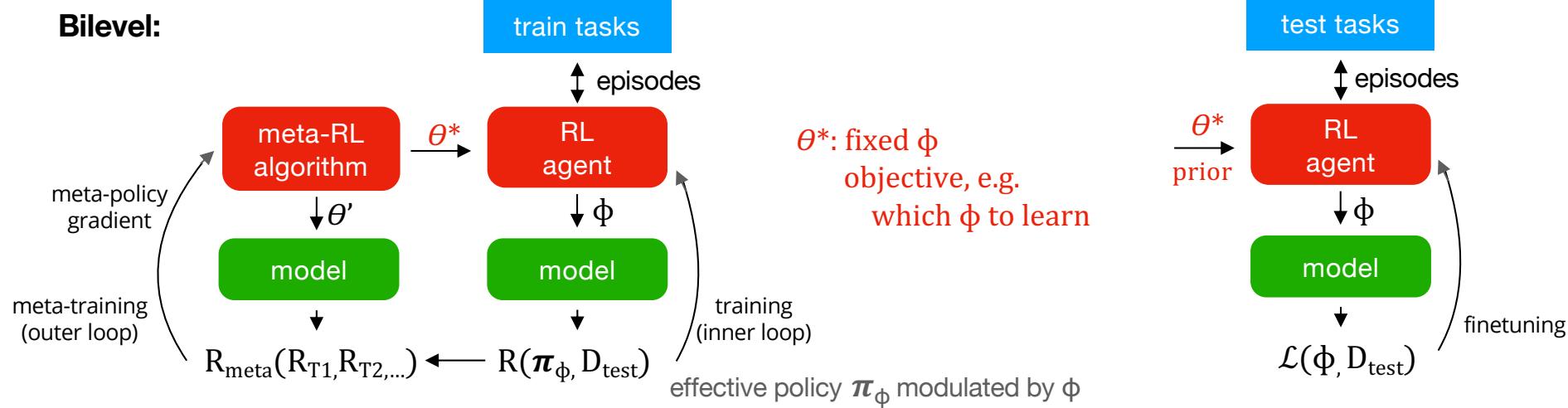
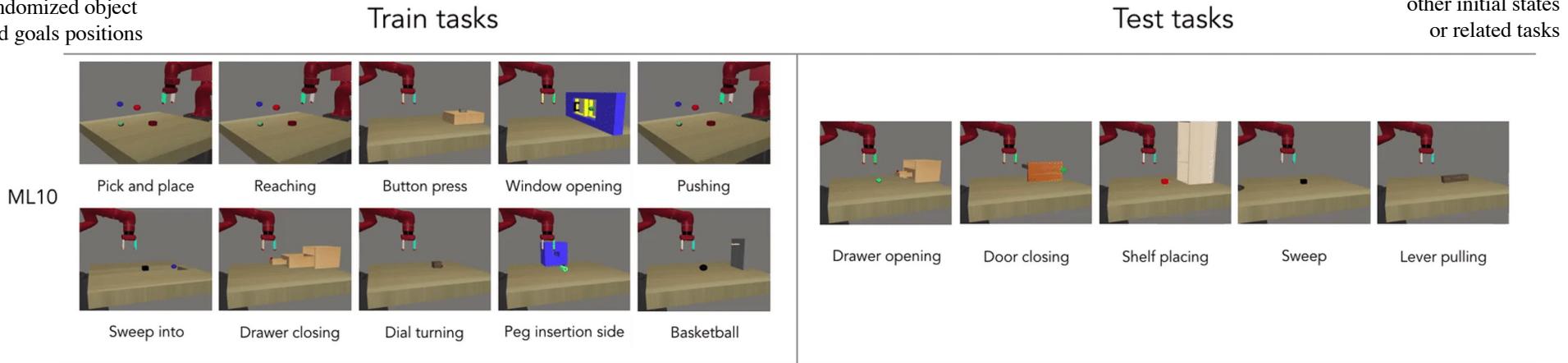
**Black Box:**



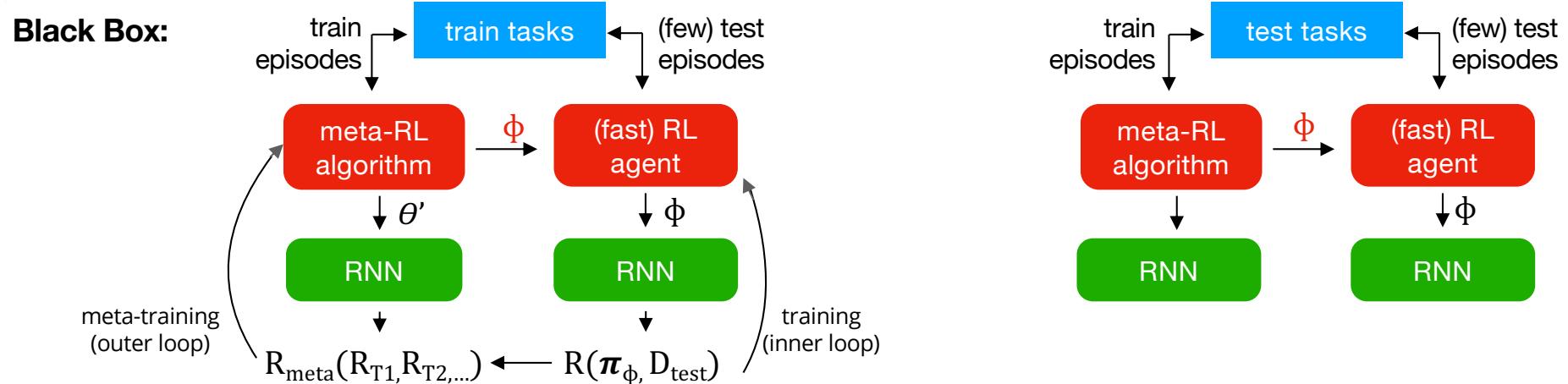
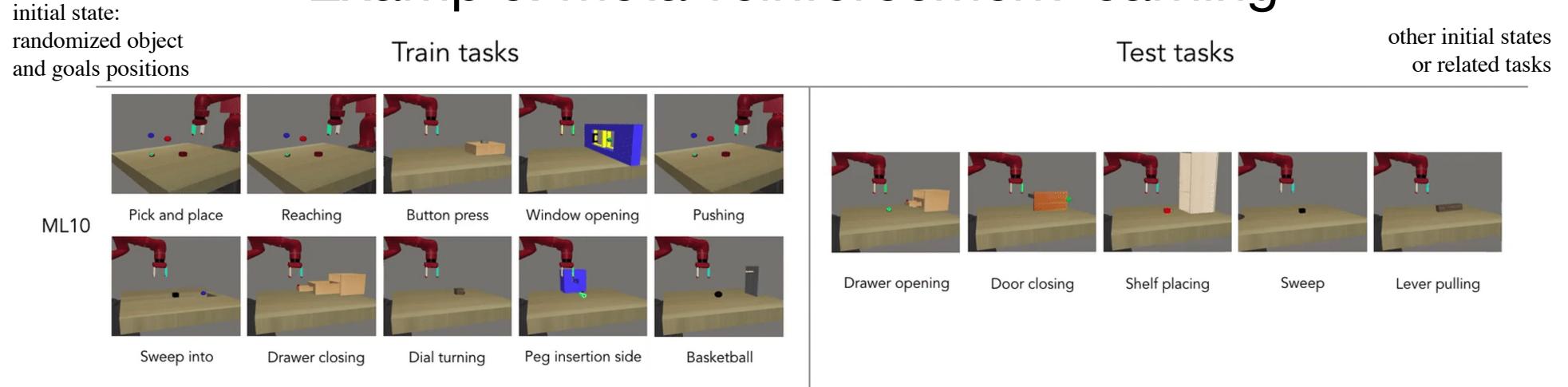
# Example: meta-reinforcement learning

Yu et al. 2019

initial state:  
randomized object  
and goals positions



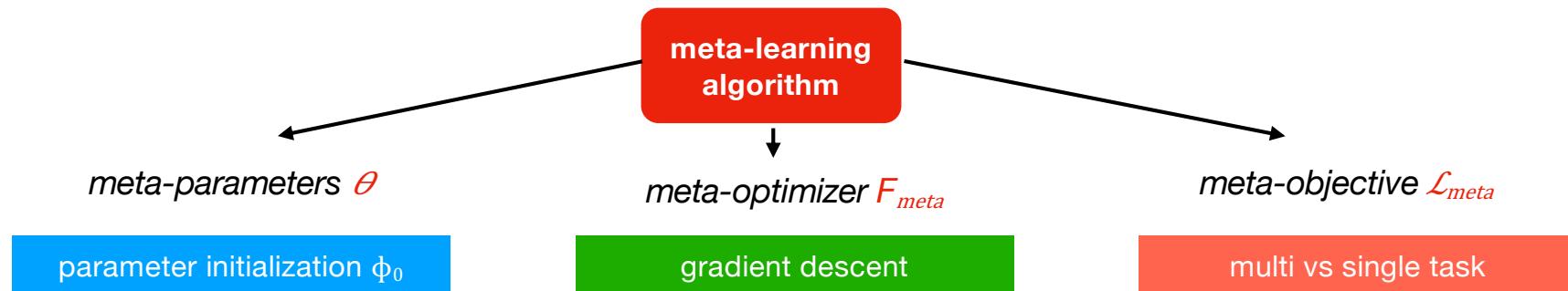
# Example: meta-reinforcement learning



# Taxonomy of meta-learning methods

[Hospedales et al. 2020](#)

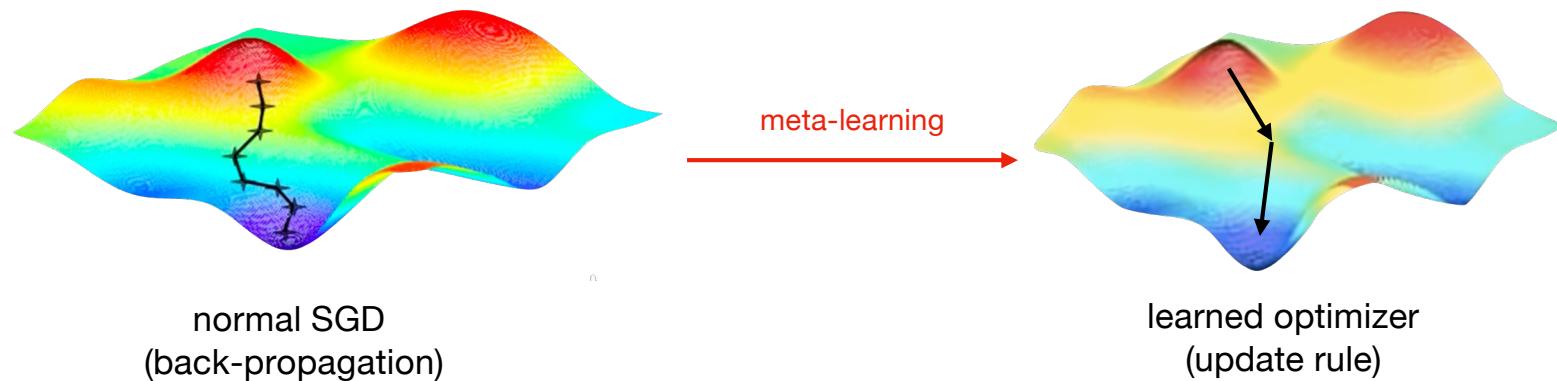
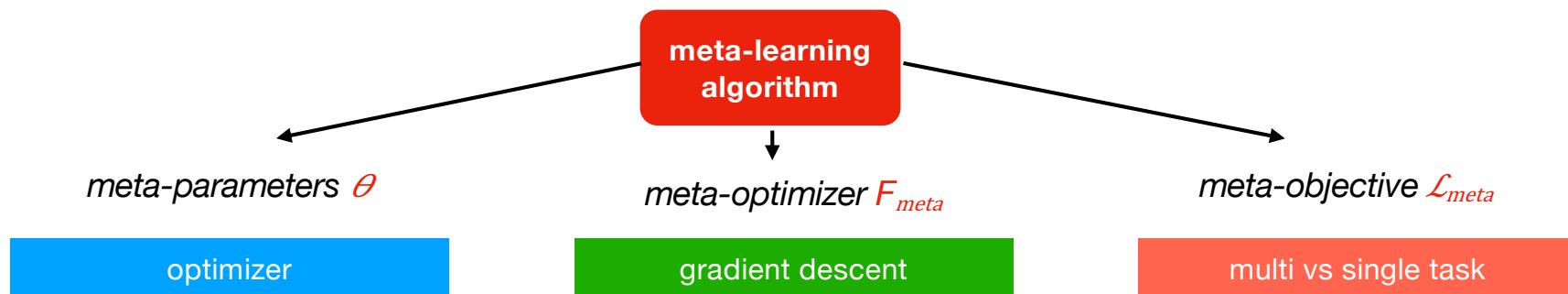
like base-learners, meta-learners consist of a representation, an objective, and an optimizer



# Taxonomy of meta-learning methods

[Hospedales et al. 2020](#)

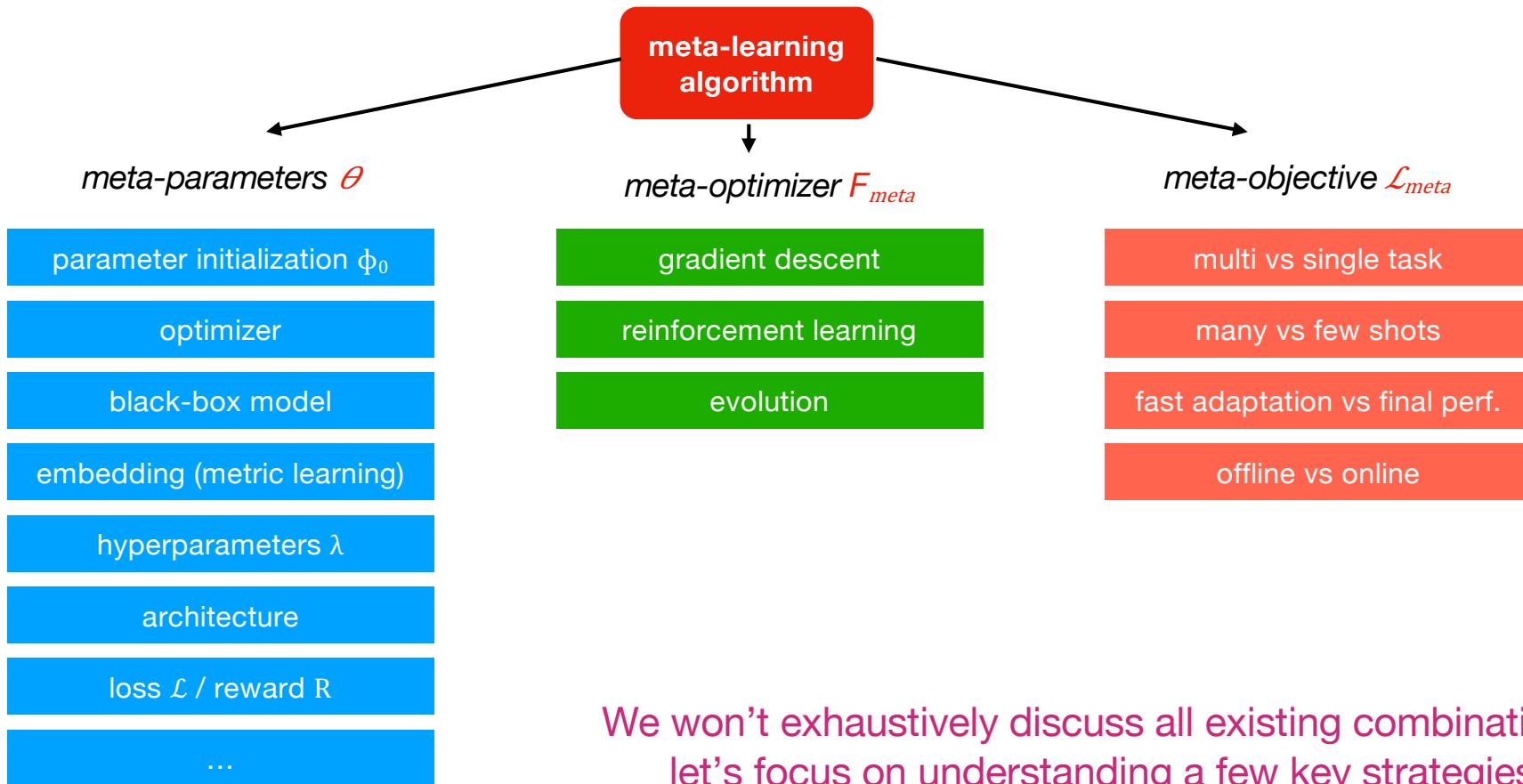
like base-learners, meta-learners consist of a representation, an objective, and an optimizer



# Taxonomy of meta-learning methods

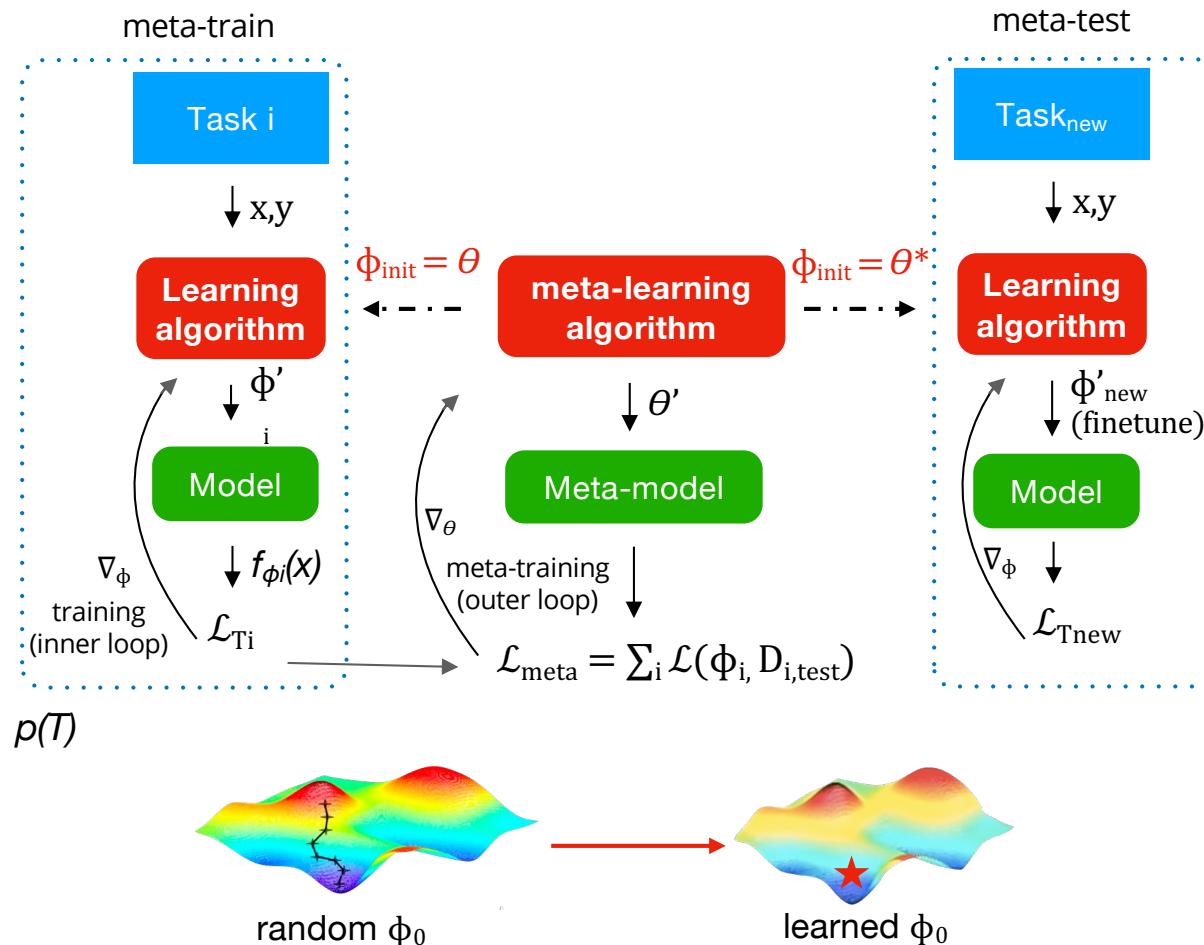
[Hospedales et al. 2020](#)  
[Huisman et al. 2020](#)

like base-learners, meta-learners consist of a representation, an objective, and an optimizer



We won't exhaustively discuss all existing combinations,  
let's focus on understanding a few key strategies

# Gradient-based methods: learning $\Phi_{\text{init}}$



- $\theta$  (prior): model initialization  $\Phi_{\text{init}}$ 
  - learn representation suitable for many tasks (e.g. pretrained CNN)
  - maximize rapid learning
- Each task  $i$  yields task-adapted  $\Phi_i$ 
  - Update algorithm  $u$
- Finetune  $\Phi^* = u(\theta^*, D_{\text{new,train}})$  (few steps)

$$\Phi'_{\text{new}} = u(\theta^*, D_{\text{new,train}})$$

Can be seen as bilevel optimization:

$$\theta^* = \operatorname{argmin}_{\theta} \sum_i \mathcal{L}_i(\Phi_i, D_{i,\text{test}})$$

$$\Phi_i = u(\theta, D_{i,\text{train}})$$

# Model agnostic meta-learning (MAML)

[Finn et al. 2017](#)

## Meta-training

- Current initialization  $\theta$ , model  $f_\theta$
- On  $i$  tasks, perform  $k$  SGD steps to find  $\phi_i^*$ , then evaluate  $\nabla_\theta \mathcal{L}_i(f_{\phi_i^*})$
- Update task-specific parameters:  $\phi_i = \theta - \alpha \nabla_\theta \mathcal{L}_i(f_{\phi_i^*})$
- Update  $\theta$  to minimize sum of per-task losses, repeat

$$\theta \leftarrow \theta - \beta \nabla_\theta \sum_i \mathcal{L}_i(f_{\phi_i})$$

derivative of test-set loss



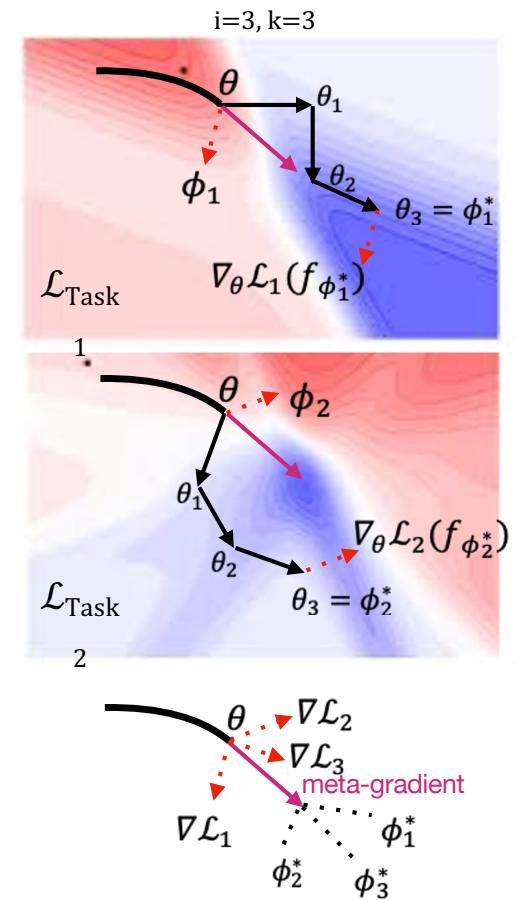
$$\theta \leftarrow \theta - \beta \nabla_\theta \sum_i \mathcal{L}_i(f(\theta - \alpha \nabla_\theta \mathcal{L}_i(f_{\phi_i^*})))$$

$\alpha, \beta$ : learning rates

meta-gradient: second-order gradient + backpropagate  
compute how changes in  $\theta$  affect the gradient at new  $\theta$

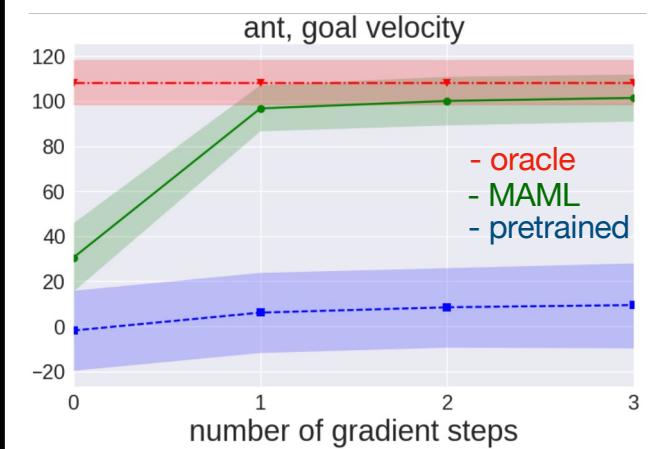
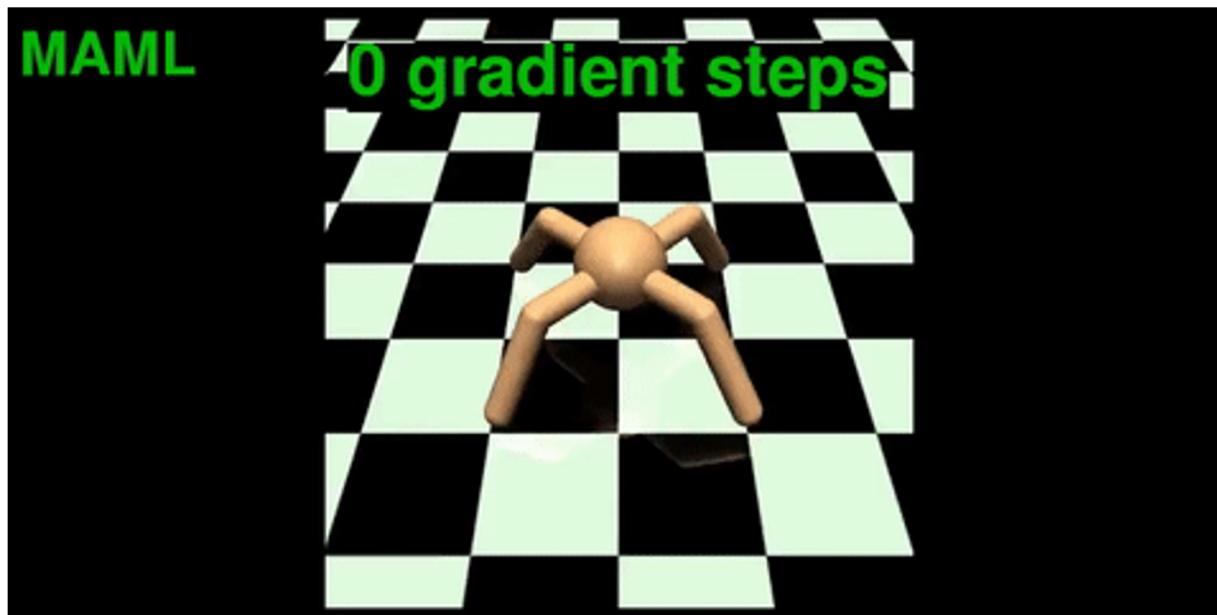
## Meta-testing

- Training data of new task  $D_{\text{train}}$
- $\theta^*$ : pre-trained parameters
- Finetune:  $\phi = \theta^* - \alpha \nabla_\theta \mathcal{L}(f_\theta)$



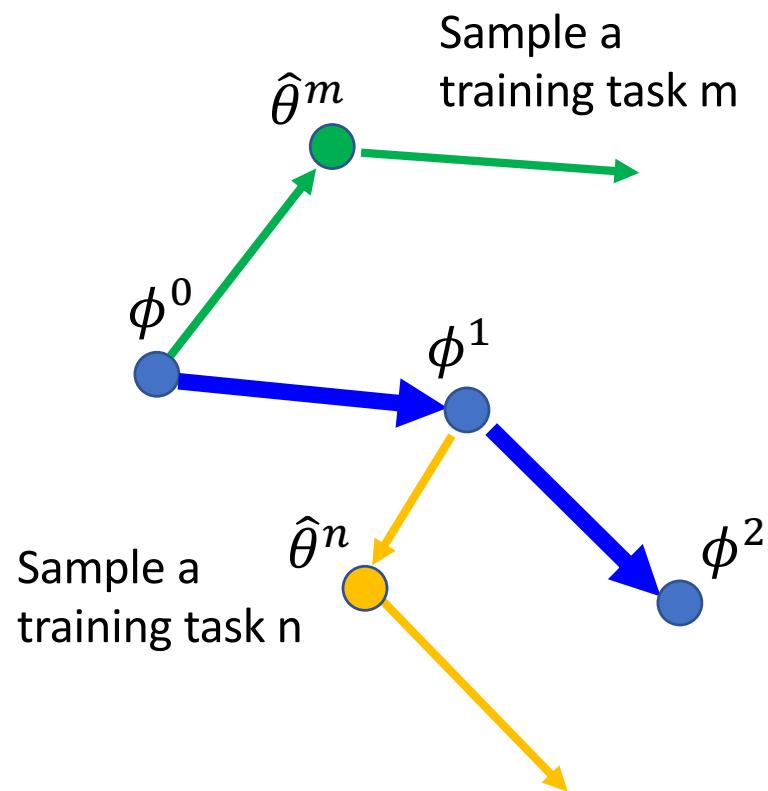
# Model agnostic meta-learning (MAML)

- Example for reinforcement learning:
  - Goal: reach certain velocity in certain direction

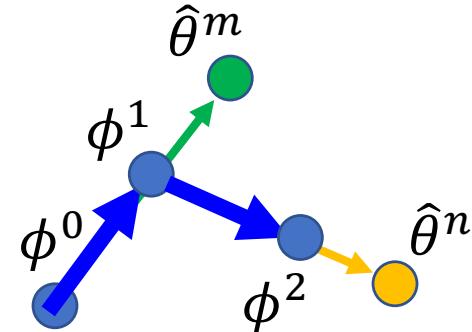


Source: Finn et al. 2017

# MAML – Real Implementation



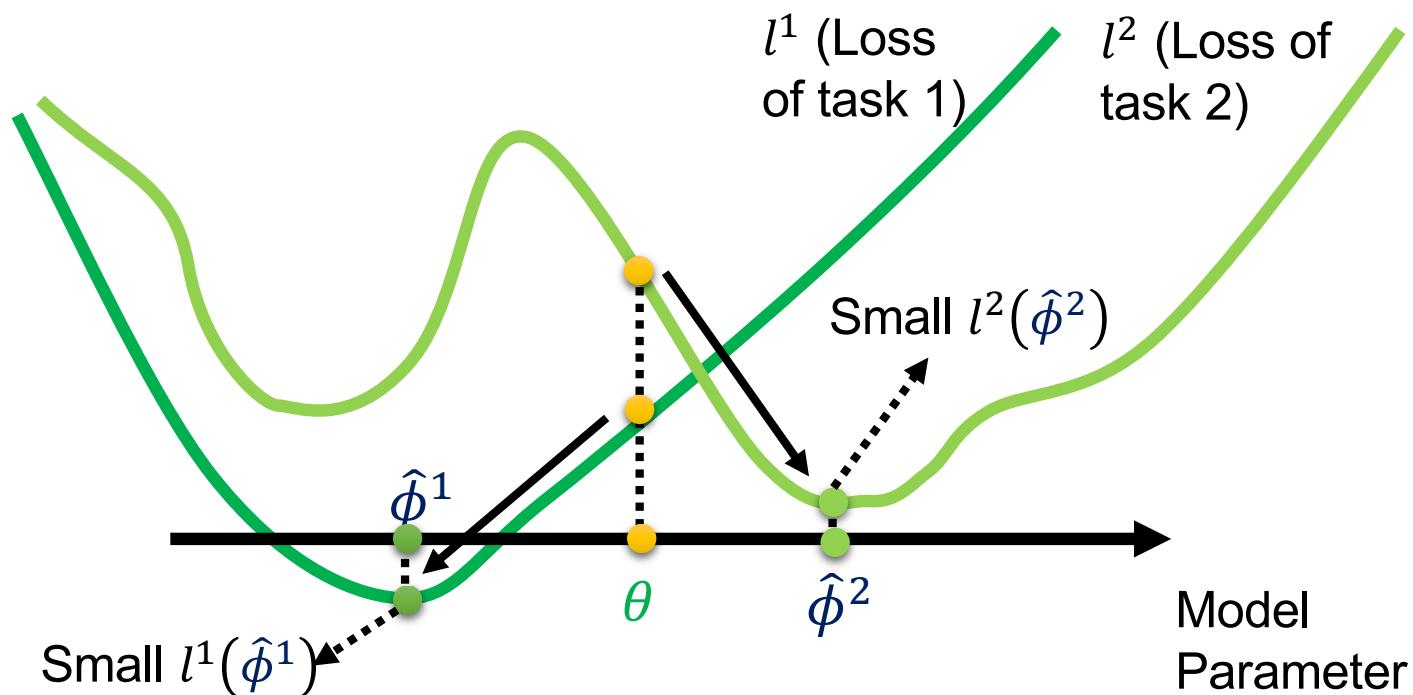
**Model Pre-training**



## MAML

Do not care about how  $\theta$  performs on training task

We care about how the performance of  $\hat{\phi}^n$  based on  $\theta$

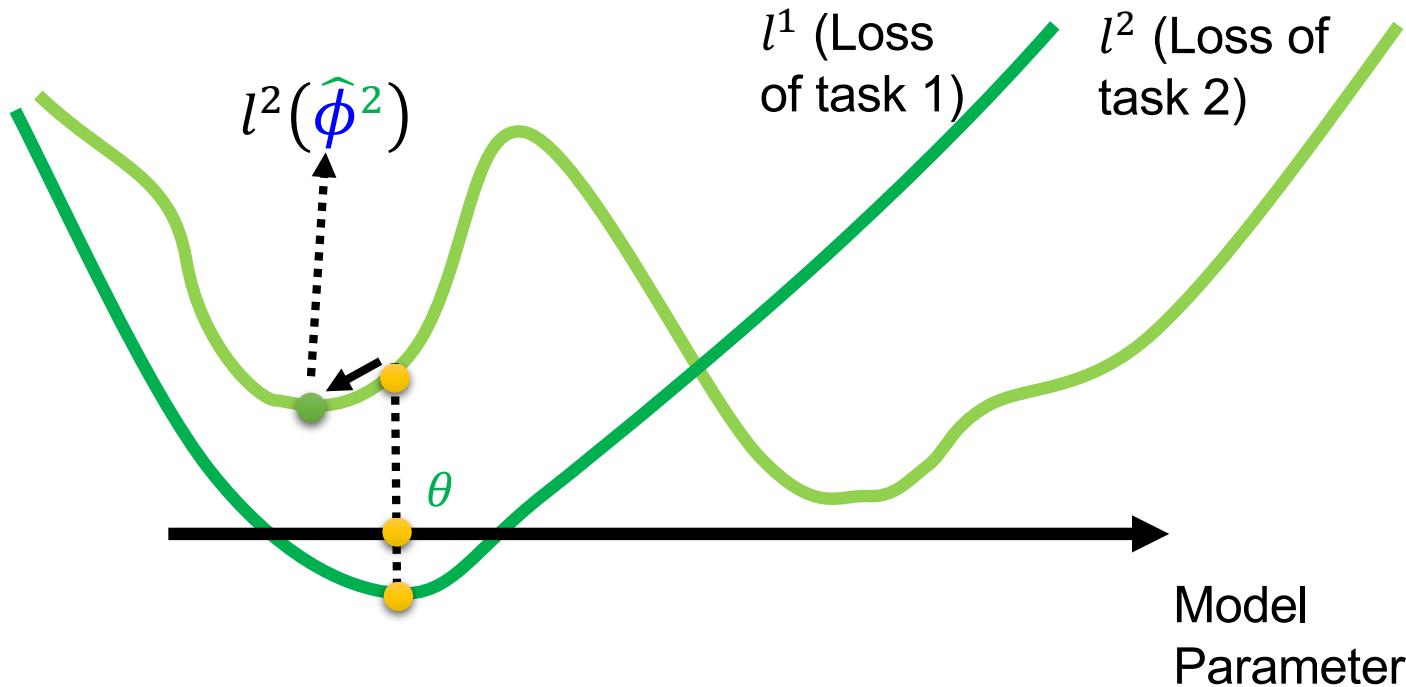


Source: Prof. Hung-Yi Lee

## Model Pre-training

Find the best  $\theta$  for all tasks

Using  $\theta$  as initialization does not guarantee the good performance of  $\hat{\phi}^n$



Source: Prof. Hung-Yi Lee

# Other gradient-based techniques

- Changing update rule yield different variants:

- MAML <sup>1,6</sup>

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_i \mathcal{L}_i(f_{\phi_i})$$

$$\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$$

- MetaSGD <sup>2</sup>

$$\phi_i = \theta - \alpha \text{diag}(w) \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$$

*w: weight per parameter*

- Tnet <sup>3</sup>

$$\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*}, w)$$

- Meta curvature <sup>4</sup>

$$\phi_i = \theta - \alpha B(\theta, w) \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$$

- WarpGrad <sup>5</sup>

$$\phi_i = \theta - \alpha P(\theta, \phi_i) \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$$

All use second-order gradients,  
Meta-learn a transformation of the gradient  
for better adaptation

- Online MAML (Follow the Meta Leader) <sup>7</sup>
  - Minimizes regret
  - Robust, but computation costs grow over time

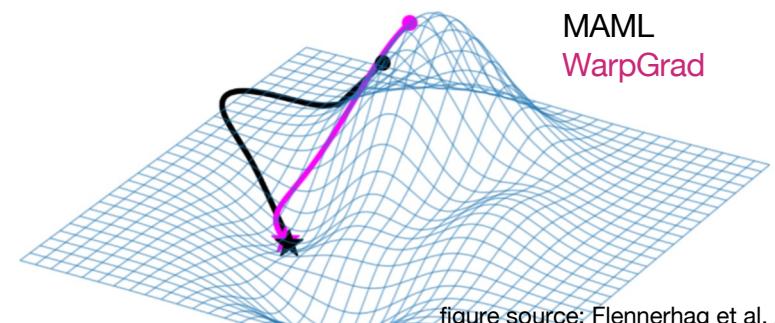


figure source: Flennerhag et al. 2019

<sup>1</sup> Finn et al. 2017

<sup>2</sup> Li et al. 2017

<sup>3</sup> Lee et al. 2018

<sup>4</sup> Park et al. 2019

<sup>5</sup> Flennerhag et al. 2019

<sup>6</sup> Antoniou et al. 2019

<sup>7</sup> Finn et al. 2019

<sup>1</sup> Finn et al.

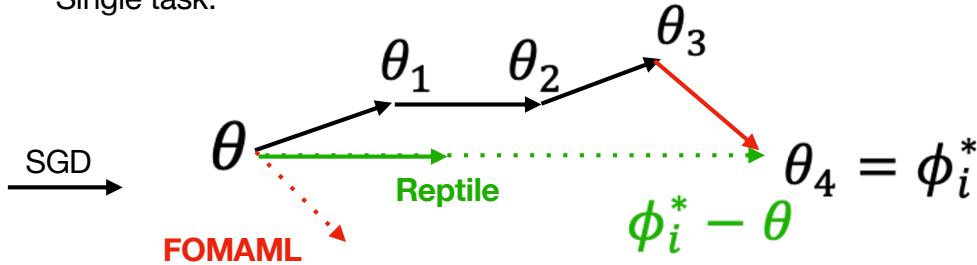
<sup>2</sup> Nichol et al. 2018

<sup>3</sup> Rajeswaran et al. 2019

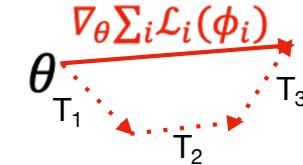
# Scalability

- Backpropagating derivates of  $\phi_i$  wrt  $\theta$  is compute + memory intensive (for large models)
- First order approximations of MAML:
  - First order MAML<sup>1</sup> uses only the last inner gradient update:  $\theta \leftarrow \theta - \beta \sum_i \mathcal{L}_i(f_{\phi_i^*})$
  - Reptile<sup>2</sup>: iterate over tasks, update  $\theta$  in direction of  $\phi_i^*$ :  $\theta \leftarrow \theta - \beta(\phi_i^* - \theta)$

Single task:



FOMAML meta-gradient:



# Scalability (2)

<sup>1</sup> [Rajeswaran et al. 2019](#)

- Implicit MAML<sup>1</sup>: uses an approximate gradient
  - Compute derivative of  $\phi_i^*$  wrt  $\theta$
  - $\phi_i^*$  could be anywhere. Hence, add penalty:  $\| \phi_i^* - \theta \|^2$
  - Accurate if we stay close to  $\theta$
  - $\mathcal{L}_i(\phi_i) + \lambda \| \phi_i^* - \theta \|^2$  has closed form solution

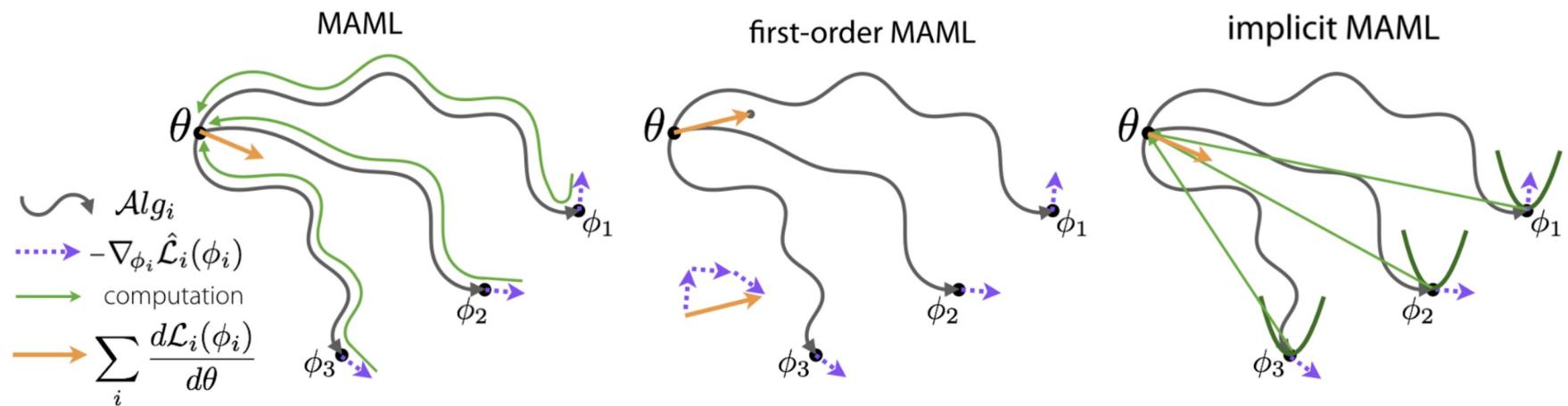
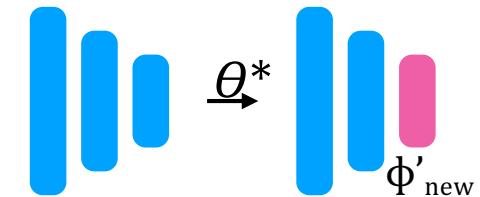


figure source: Rajeswaran et al. 2019

# Generalizability

- <sup>1</sup> [Finn et al. 2018](#)
- <sup>2</sup> [Raghu et al. 2020](#)
- <sup>3</sup> [Tian et al. 2020](#)
- <sup>4</sup> [Stadie et al. 2019](#)

- MAML is more resilient to overfitting than many other meta-learning techniques <sup>1</sup>
- Effectiveness seems mainly due to feature reuse (finding  $\theta$ ) <sup>2</sup>
  - Fine-tuning *only the last layer* equally good
- On few-shot learning benchmarks, a good embedding outperforms most meta-learning <sup>3</sup>
  - Learn representation on entire meta-dataset (merged into single task)
  - Train a linear classifier on embedded few-shot  $D_{\text{train}}$ , predict  $D_{\text{test}}$



- For meta-RL, also learn how to *explore* (how to sample new environments)
  - E-MAML: add exploration to meta-objective (allows longer term goals) <sup>4</sup>

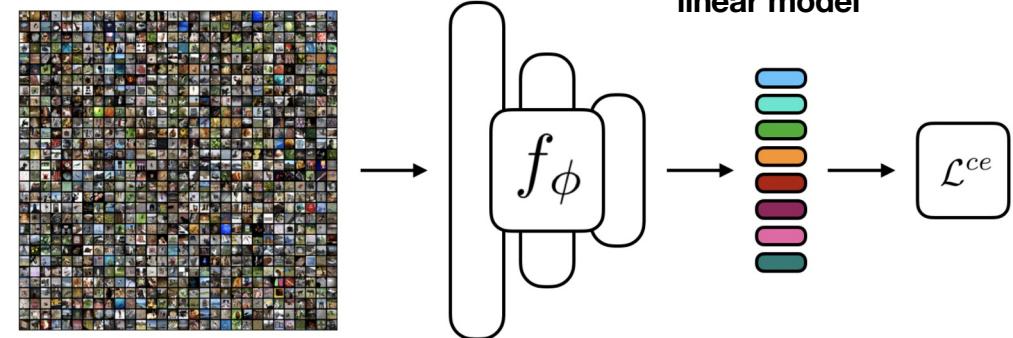
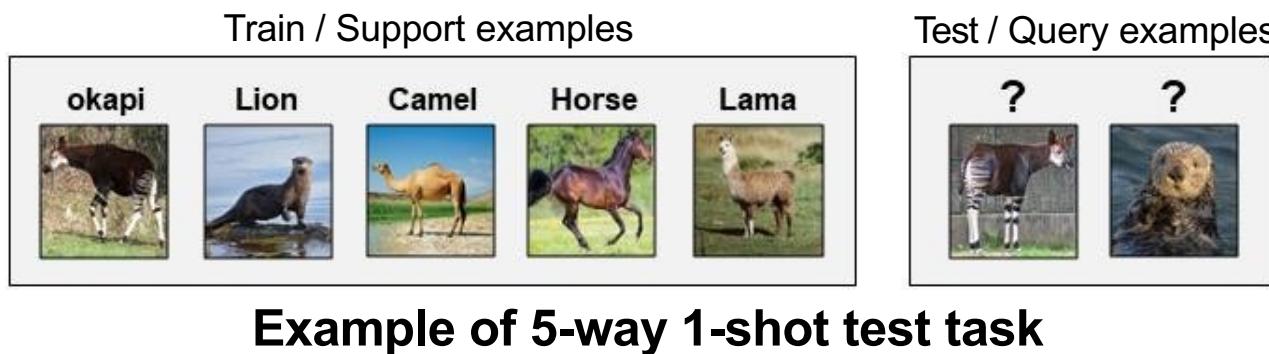


figure source: [Tian et al. 2020](#)

# Agenda

- Introduction
  - Few-shot learning problem
  - Meta-learning paradigm
  - **How to evaluate**
- Main types of few-shot learning algorithms
- Few-shot learning without forgetting

## How to evaluate few-shot algorithms



**2<sup>nd</sup> learning stage** (meta-test time for meta-learning):

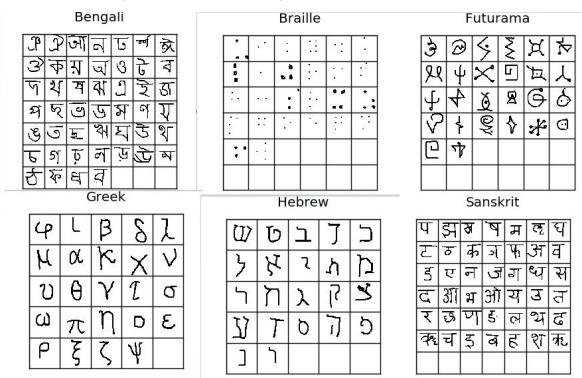
- Use a **held out set of classes**
- Sample a **large number of N-way K-shot few-shot tasks**
- **Report average accuracy** on the  $N \times M$  query examples of all tasks

# How to evaluate few-shot algorithms

## Datasets / benchmarks

### Omniglot: Lake et al. 11

- 1623 characters from 50 alphabets
- 20 instances per character / class
- 5-way and 20-way 1-shot or 5-shot tasks



### MinilmageNet: Ravi et al. 17

- 84x84 sized images
- 100 classes: 64 train, 16 val, 20 test
- 1-shot 5-way & 5-shot 5-way tasks



### ImageNet-FS: Hariharan et al. 17

- normal ImageNet images
- classes: 389 train, 300 val, 311 test
- 311-way 1, 2, 5, 10, or 20 shot tasks
- more realistic & challenging setting



Also: *tiered-MinilmageNet* (Ren et. al. 18), *CIFAR-FS* (Bertinetto et al 19), *CUB*, *Tracking in the wild* (Valmadre et al. 18), ...

# Agenda

- Introduction
- Main types of few-shot learning algorithms
  - Metric learning
  - Meta-learning with memory modules
  - Learn to predict model parameters
- Few-shot learning without forgetting

# Agenda

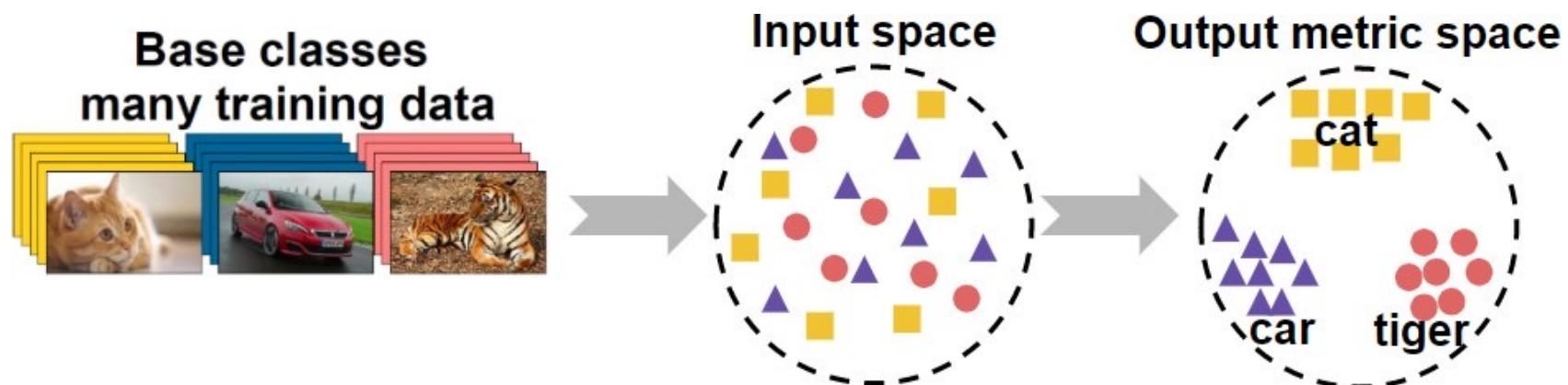
- Introduction
- Main types of few-shot learning algorithms
  - Metric learning
  - Meta-learning with memory modules
  - Learn to predict model parameters
- Few-shot learning without forgetting

**Disclaimer:** loose categorization, many combine elements of several types, not exhaustive enumeration

# Agenda

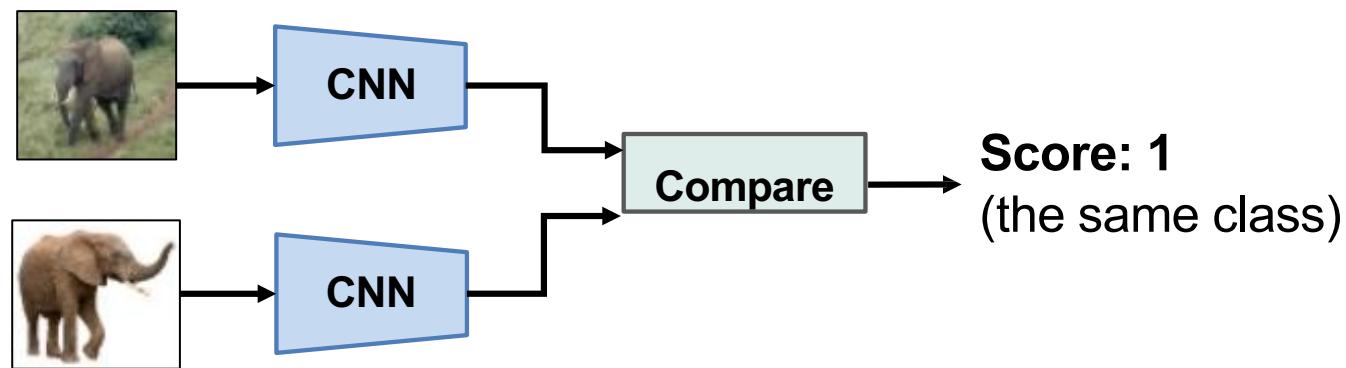
- Introduction
- Main types of few-shot learning algorithms
  - **Metric learning**
  - Meta-learning with memory modules
  - Learn to predict model parameters
- Few-shot learning without forgetting

## Metric learning for few-shot classification



- **1<sup>st</sup> learning stage:** train a deep metric function on the base class data
- **2<sup>nd</sup> learning stage:** use it as a nearest neighbor classifier to novel classes
  - **Non-parametric** at this stage
  - **Simple and works well with limited data**

## Siamese neural networks



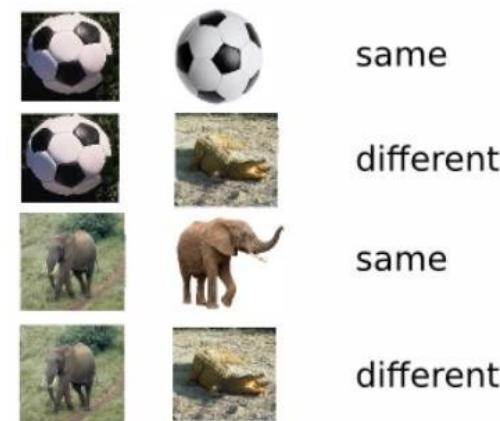
### Siamese network:

- Given two images: outputs a similarity / distance score.
- Similarity score: 1 if the two images belong to the same class, 0 otherwise

“Siamese neural networks for one-shot image recognition”, O. Koch et. al. 2015

## Siamese neural networks

**1<sup>st</sup> learning stage – verification task:**  
Learn with a siamese convnet if 2 images belong to same / different classes.



**2<sup>nd</sup> stage (convnet is fixed):**  
Classify query to most similar support image



“Siamese neural networks for one-shot image recognition”, O. Koch et. al. 2015

# Metric learning

## **Extensive work on (deep) metric learning:**

- “Neighborhood Component Analysis”, Goldberger et. al. 05
- “Dimensionality Reduction by Learning an Invariant Mapping”, Hadsell et. al. 06
- “Distance Metric Learning for Large Margin Nearest Neighbor Classification”, Weinberger et. al. 09
- “Deep Metric Learning Using Triplet Network”, Hoffer et. al. 15
- “Web-Scale Training for Face Identification”, Taigman et. al. 15
- “FaceNet: A Unified Embedding for Face Recognition and Clustering”, Schroff et al 15
- ...

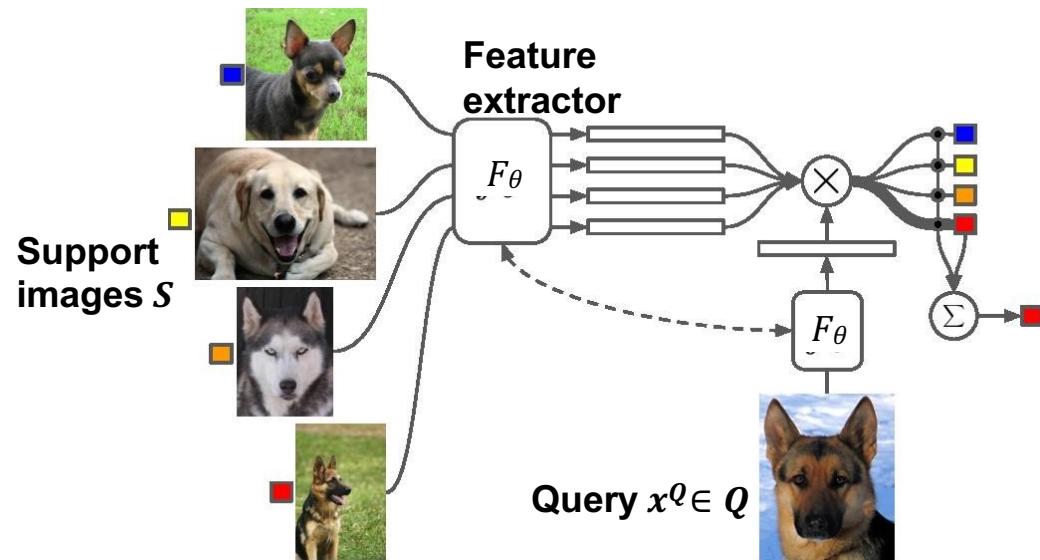
## Meta-training based metric learning

Train the metric model on the same way it would be used at 2<sup>nd</sup> learning stage

- “Matching Networks for one-shot learning”, O. Vinyals et al. 16

# Matching Networks

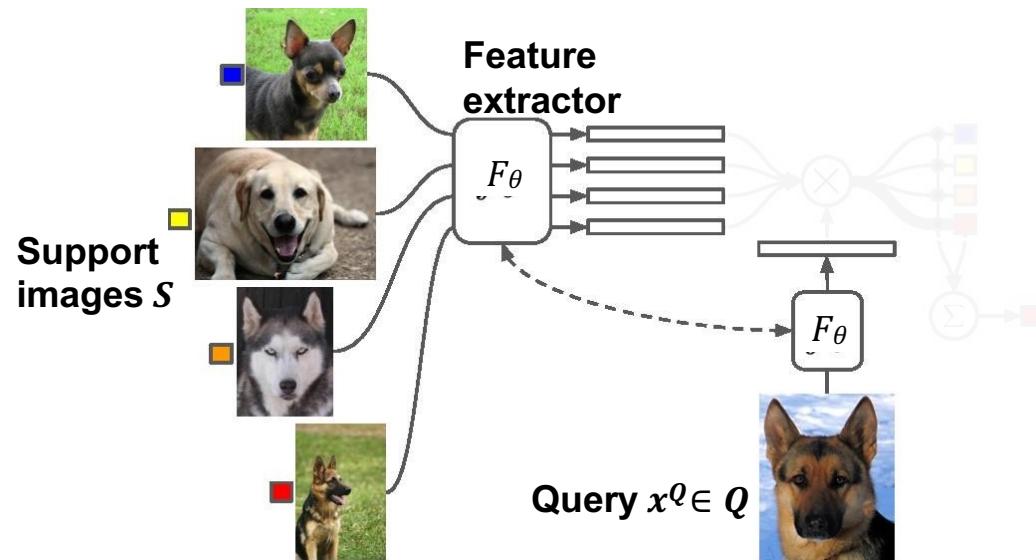
- Learn to match



“Matching networks for one shot learning”, Vinyals et al. 2016

## Matching Networks

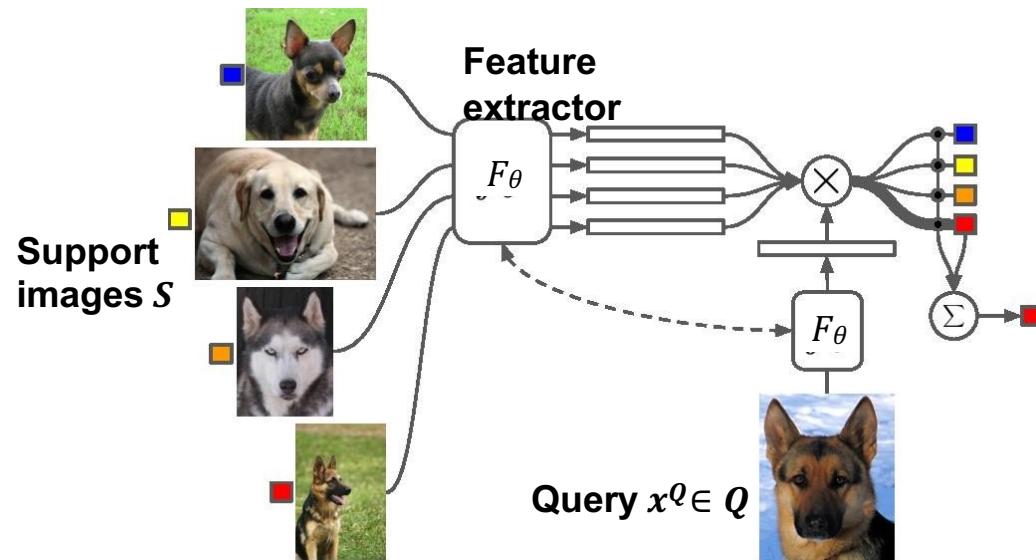
- Learn to match
  - Extract features from the query and support images



“Matching networks for one shot learning”, Vinyals et al. 2016

# Matching Networks

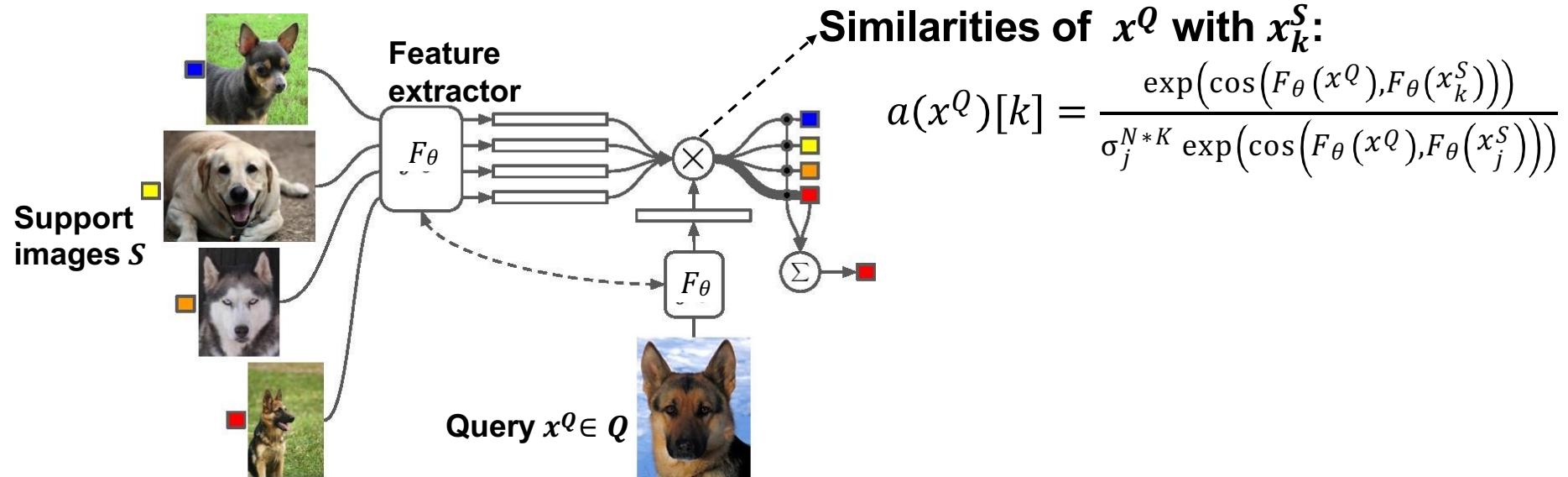
- **Learn to match**
  - Extract features from the query and support images
  - Classify with **differentiable (soft) nearest neighbor classifier**



“Matching networks for one shot learning”, Vinyals et al. 2016

## Matching Networks

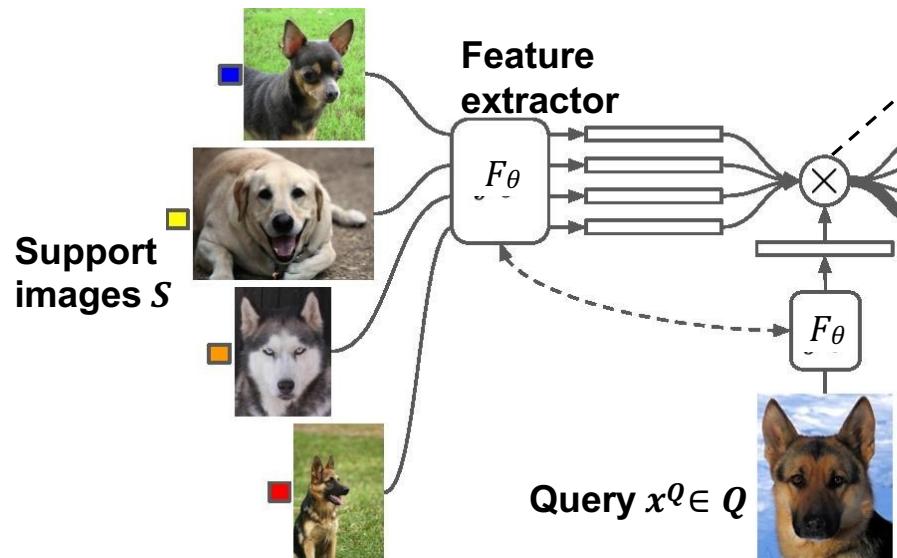
- **Learn to match**
  - Extract features from the query and support images
  - Classify with **differentiable (soft) nearest neighbor classifier**



“Matching networks for one shot learning”, Vinyals et al. 2016

## Matching Networks

- **Learn to match**
  - Extract features from the query and support images
  - Classify with **differentiable (soft) nearest neighbor classifier**



**Similarities of  $x^Q$  with  $x_k^S$ :**

$$a(x^Q)[k] = \frac{\exp(\cos(F_\theta(x^Q), F_\theta(x_k^S)))}{\sum_j^{N*K} \exp(\cos(F_\theta(x^Q), F_\theta(x_j^S)))}$$

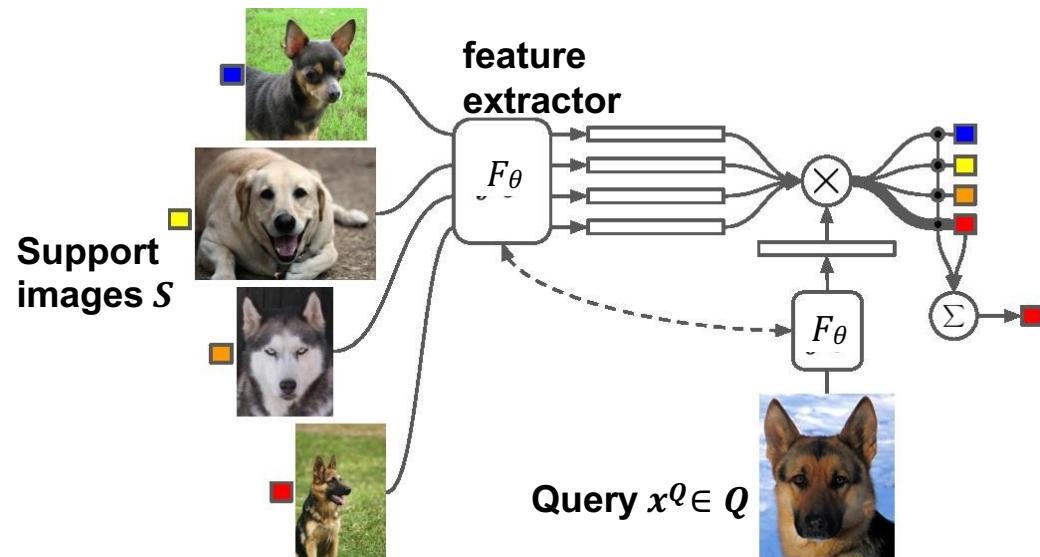
**Classification probabilities for  $x^Q$ :**

$$p = m_\varphi(x^Q) = \sum_k a(x^Q)[k] \cdot \text{one\_hot}(y_k^S)$$

“Matching networks for one shot learning”, Vinyals et al. 2016

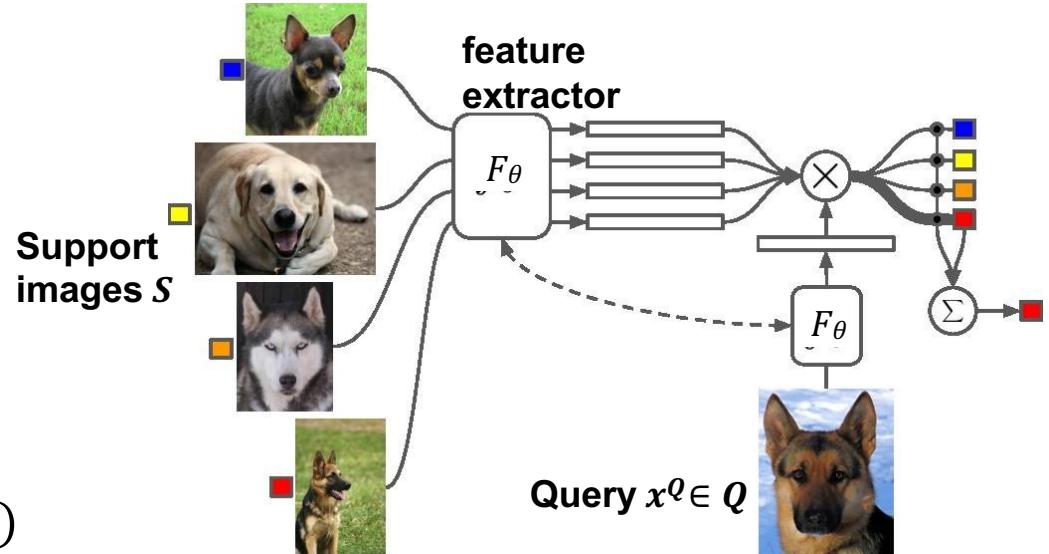
## Meta-training in Matching Networks

- **Meta-learner  $f_\theta$ :** feature extractor  $F_\theta(\cdot)$
- **Generated model  $m_\varphi$ :** extractor  $F_\theta(\cdot)$  with support features  $\{F_\theta(x_k^S), y_k^S\}_{k=1}^{N*K}$



“Matching networks for one shot learning”, Vinyals et al. 2016

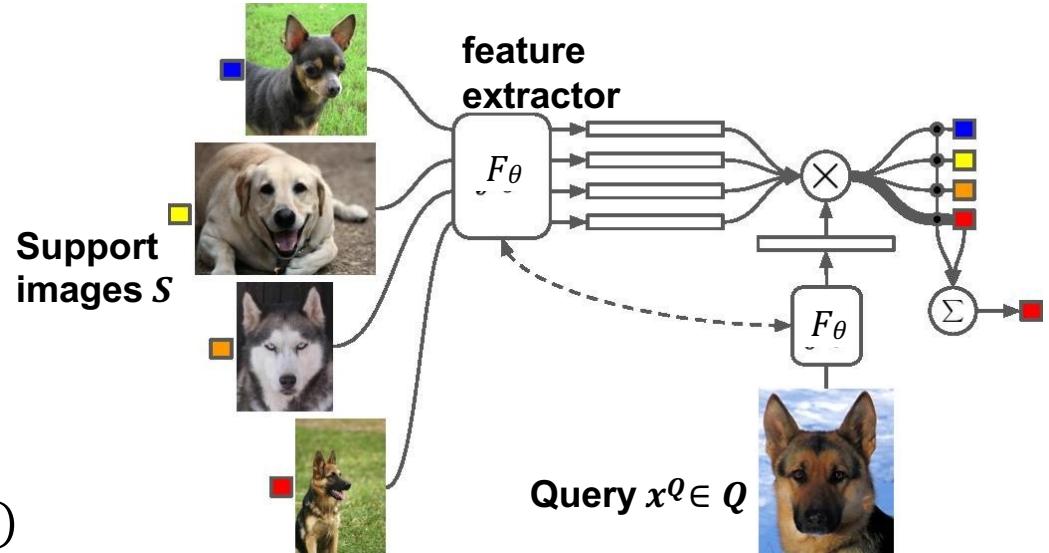
## Meta-training in Matching Networks



### Meta-training routine:

1. Sample training episode  $(S, Q)$
2. Generate classification model  $m_\varphi = f(S) = \left\{F_\theta(\cdot), \{F_\theta(x_k^S), y_k^S\}_{k=1}^{N*K}\right\}$
3. Predict classification scores  $p_m \stackrel{\theta}{=} m_\varphi(x_m^Q) = \sigma_k a(x_m^Q)[k] \cdot \text{one\_hot}(y_k^S)$
4. Optimize  $\theta$  w.r.t. the query classification loss  $L(f_\theta(S), Q) = \sigma_m - \log(p_m[y_m^Q])$

## Meta-training in Matching Networks



### Meta-training routine:

1. Sample training episode  $(S, Q)$
2. **Generate classification model**  $m_\varphi = f_\theta(S) = \left\{F_\theta(\cdot), \{F_\theta(x_k^S), y_k^S\}_{k=1}^{N*K}\right\}$
3. **Predict classification scores**  $p_m = m_\varphi(x_m^Q) = \sigma_k a(x_m^Q)[k] \cdot \text{one\_hot}(y_k^S)$
4. Optimize  $\theta$  w.r.t. the query classification loss  $L(f_\theta(S), Q) = \sigma_m - \log(p_m[y_m^Q])$

# Matching Networks

Model	Fine Tune	5-way Acc		20-way Acc	
		1-shot	5-shot	1-shot	5-shot
<b>BASELINE CLASSIFIER</b>	Y	86.0%	97.6%	72.9%	92.3%
<b>MANN (NO CONV) [21]</b>	N	82.8%	94.9%	—	—
<b>CONVOLUTIONAL SIAMESE NET [11]</b>	N	96.7%	98.4%	88.0%	96.5%
<b>CONVOLUTIONAL SIAMESE NET [11]</b>	Y	97.3%	98.4%	88.1%	97.0%
<b>MATCHING NETS (OURS)</b>	N	<b>98.1%</b>	<b>98.9%</b>	<b>93.8%</b>	98.5%
<b>MATCHING NETS (OURS)</b>	Y	97.9%	98.7%	93.5%	<b>98.7%</b>

Table 1: Results on the Omniglot dataset.

Model	Fine Tune	5-way Acc	
		1-shot	5-shot
<b>BASELINE CLASSIFIER</b>	Y	38.4%	51.2%
<b>MATCHING NETS (OURS)</b>	N	44.2%	57.0%
<b>MATCHING NETS (OURS)</b>	Y	<b>46.6%</b>	<b>60.0%</b>

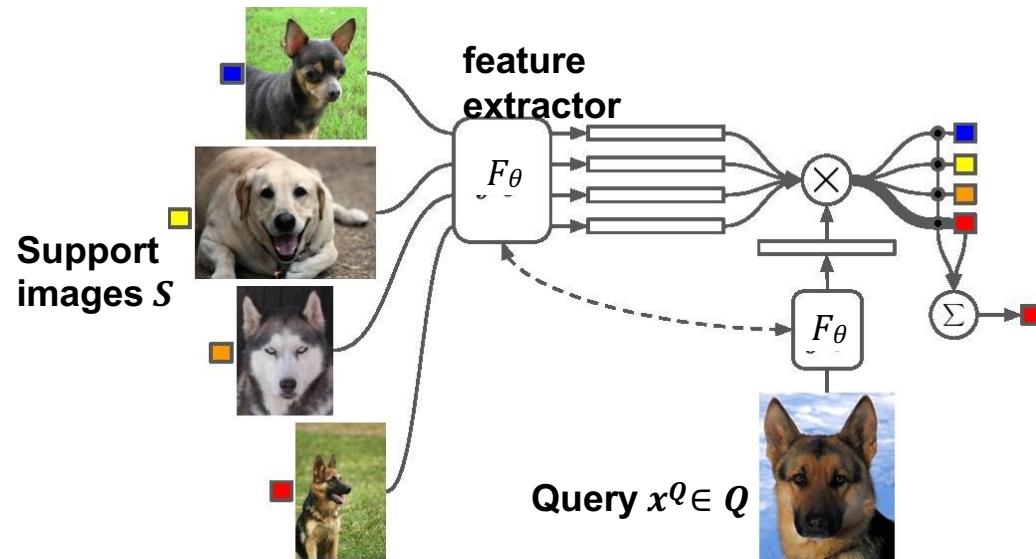
Table 2: Results on *miniImageNet*.

- **Metric learning:**  
better results than pre-training & fine-tuning
- **Meta-training:**  
improves over siamese networks

## Matching Networks

$K > 1$  support example per class:

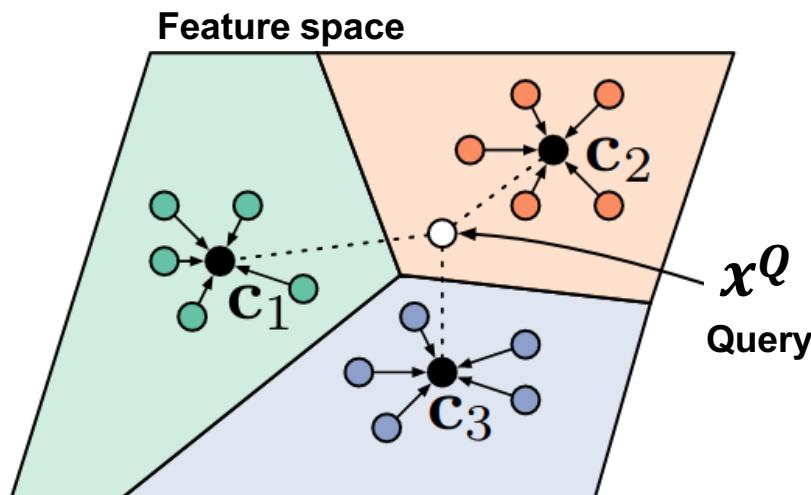
- **Independently matches a query with each support example**
- Can we do something smarter?



“Matching networks for one shot learning”, Vinyals et al. 2016

## Prototypical Networks

- Learn to extract class prototypes for comparisons:
  - prototype: aggregates information of all support images in a class

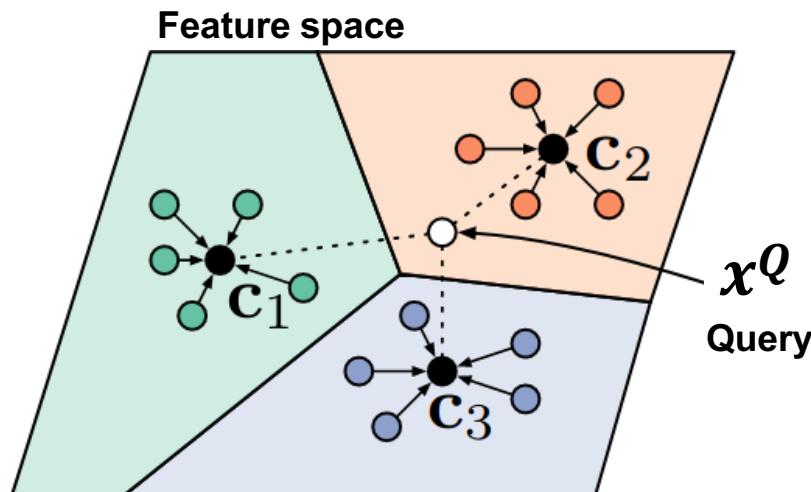


“Prototypical Networks for Few-Shot Learning”, Snell et al. 2017

## Prototypical Networks

- **prototype  $i$ -th class = mean training feature vector of its support set  $S_i$** 
  - K=1: the same as matching networks

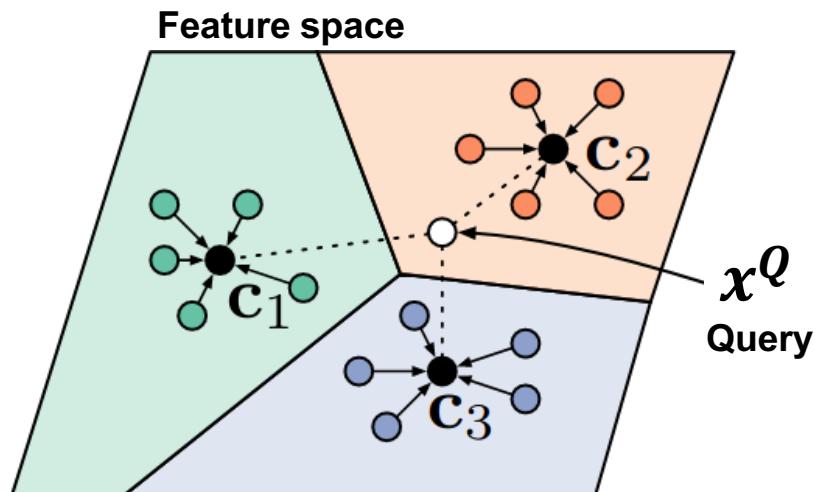
$$c_i = \frac{1}{|S_i|} \sum_{(x_k^S, y_k^S) \in S_i} \theta(x_k^S)$$



## Prototypical Networks

- prototype  $i$ -th class = mean training feature vector of its support set  $S_i$

$$c_i = \frac{1}{|S_i|} \underset{(x_k^S, y_k^S) \in S_i}{\underset{F}{\mathop{\Sigma}}} \theta(x_k^S)$$



- Classify to closest prototype with prob.

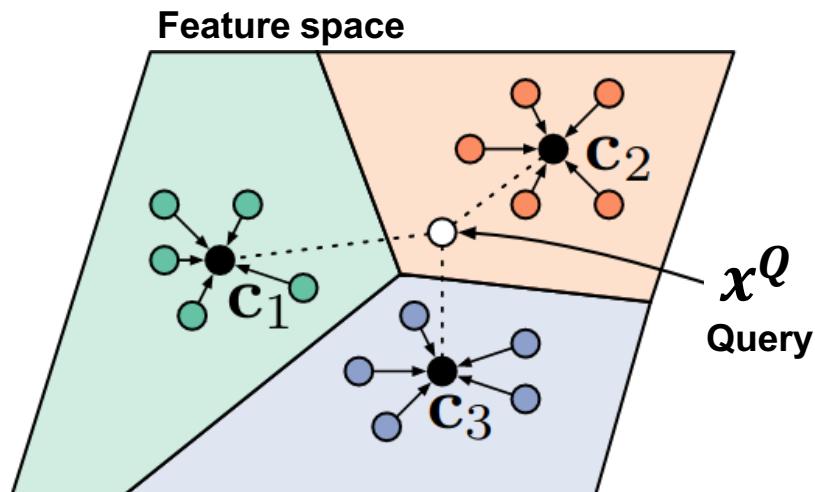
$$p[i] = m_\varphi(x^Q)[i] = \frac{\exp(-\text{dist}(F_\theta(x^Q), c_i))}{\sum_j \exp(-\text{dist}(F_\theta(x^Q), c_j))}$$

Distance  $\text{dist}(\cdot, \cdot)$ : Euclidean or cosine

## Prototypical Networks

- prototype  $i$ -th class = mean training feature vector of its support set  $S_i$

$$c_i = \frac{1}{|S_i|} \underset{(x_k^S, y_k^S) \in S_i}{\underset{F}{\mathop{\Sigma}}} \theta(x_k^S)$$



- Classify to closest prototype with prob.

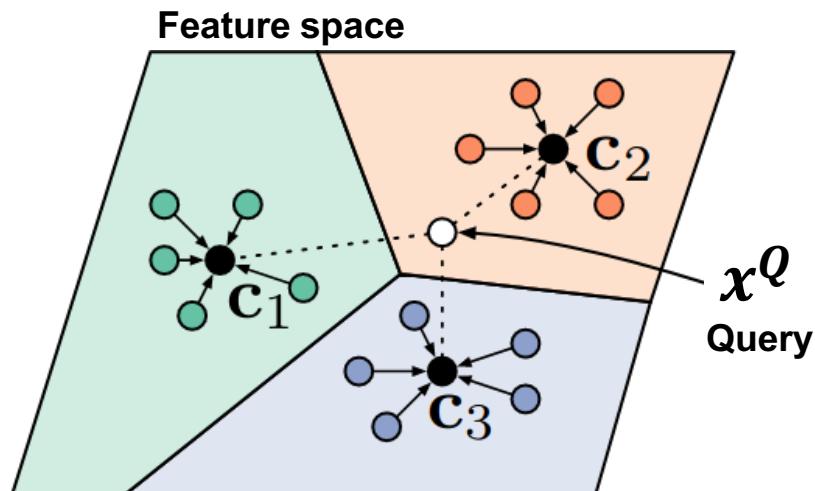
$$p[i] = m_\varphi(x^Q)[i] = \frac{\exp(-\text{dist}(F_\theta(x^Q), c_i))}{\sum_j \exp(-\text{dist}(F_\theta(x^Q), c_j))}$$

Prototypes: similar to output weights of a classification network with bias = 0

## Prototypical Networks

- prototype  $i$ -th class = mean training feature vector of its support set  $S_i$

$$c_i = \frac{1}{|S_i|} \underset{(x_k^S, y_k^S) \in S_i}{\underset{F}{\mathop{\Sigma}}} \theta(x_k^S)$$

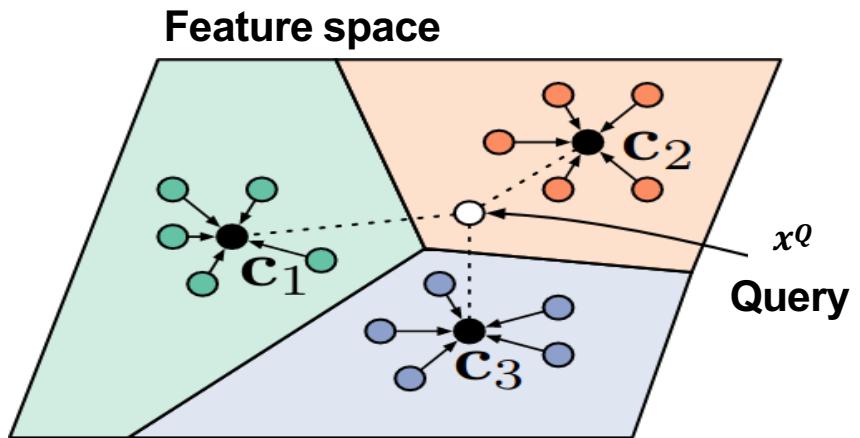


- Classify to closest prototype with prob.

$$p[i] = m_\varphi(x^Q)[i] = \frac{\exp(-\text{dist}(F_\theta(x^Q), c_i))}{\sum_j \exp(-\text{dist}(F_\theta(x^Q), c_j))}$$

During meta-training (optimizing  $F_\theta$ ):  
back-propagate through the prototypes too

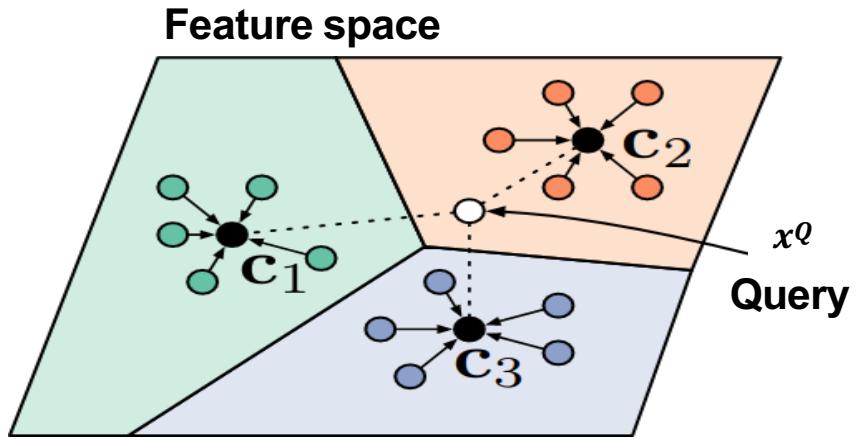
## Meta-training in Prototypical Networks



### Meta-training routine:

1. Sample training episode  $(S, Q)$
2. Generate classification model  $m_\varphi = f_\theta(S) = \{F_\theta(\cdot), \{c_i\}_{i=1}^N\}$
3. Predict classification scores  $p_m = m_\varphi(x_m^Q) = \frac{\exp(-\text{dist}(F_\theta(x_m^Q), c_i))}{\sum_j^N \exp(-\text{dist}(F_\theta(x_m^Q), c_j))}$
4. Optimize  $\theta$  w.r.t. the query classification loss  $L(f_\theta(S), Q) = \sigma_m - \log(p_m[y_m^Q])$

## Meta-training in Prototypical Networks



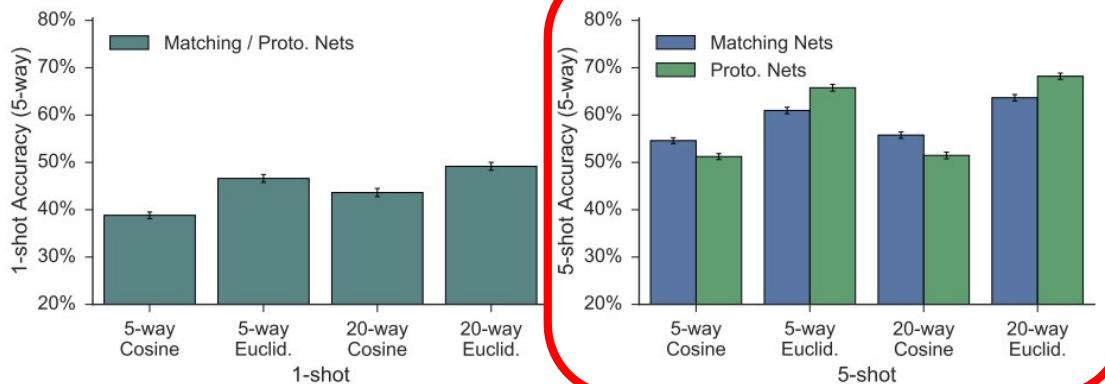
### Meta-training routine:

1. Sample training episode  $(S, Q)$
2. **Generate classification model**  $m_\varphi = f_\theta(S) = \left\{ F_\theta(\cdot), \{c_i\}_{i=1}^N \right\}$
3. **Predict classification scores**  $p_m = m_\varphi(x_m^Q) = \frac{\exp(-\text{dist}(F_\theta(x^Q), c_i))}{\sum_j^N \exp(-\text{dist}(F_\theta(x^Q), c_j))}$
4. Optimize  $\theta$  w.r.t. the query classification loss  $L(f_\theta(S), Q) = \sigma_m^j - \log(p_m[y_m^Q])$

# Prototypical Networks

Table 2: Few-shot classification accuracies on *miniImageNet*. All accuracy results are averaged over 600 test episodes and are reported with 95% confidence intervals.

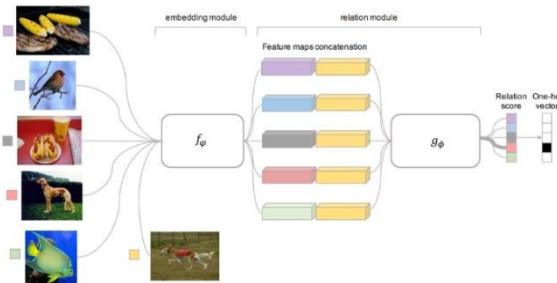
Model	Dist.	Fine Tune	5-way Acc.	
			1-shot	5-shot
<b>BASELINE NEAREST NEIGHBORS*</b>	Cosine	N	$28.86 \pm 0.54\%$	$49.79 \pm 0.79\%$
<b>MATCHING NETWORKS [29]*</b>	Cosine	N	$43.40 \pm 0.78\%$	$51.09 \pm 0.71\%$
<b>MATCHING NETWORKS FCE [29]*</b>	Cosine	N	$43.56 \pm 0.84\%$	$55.31 \pm 0.73\%$
<b>META-LEARNER LSTM [22]*</b>	-	N	$43.44 \pm 0.77\%$	$60.60 \pm 0.71\%$
<b>PROTOTYPICAL NETWORKS (OURS)</b>	Euclid.	N	<b><math>49.42 \pm 0.78\%</math></b>	<b><math>68.20 \pm 0.66\%</math></b>



**For K>1 shots per class: prototype vectors with Euclidean distance have better accuracy than individual comparison with each support example (Matching Nets)**

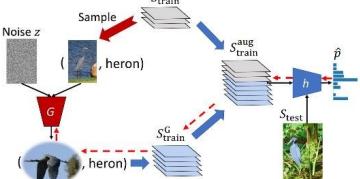
# Meta-training based metric learning

**Implement distance function in prototypical nets with a relation network**  
 "Learning to Compare: Relation Network for Few-Shot Learning", Sung et. al. 18

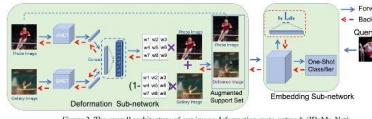


**Learn to synthesize additional support examples for the metric function**

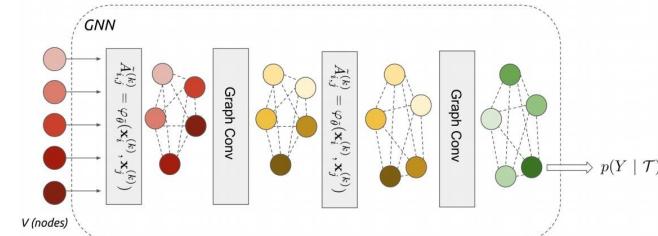
"Low-shot learning from imaginary data", Wang et.al. 18



"Image deformation meta-networks for one-shot learning", Chen et.al. 19

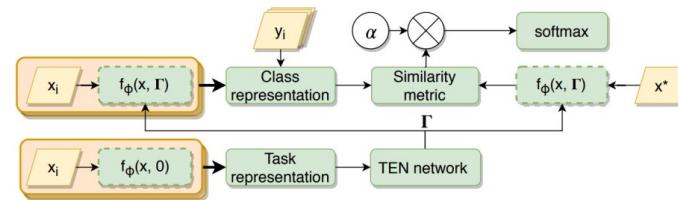


**Propagate with a GNN information from the labeled support set to the query**  
 "Few-shot Learning with Graph Neural Networks", Garcia et al. 18



**Task-adaptive metric function based on task-context representations**

"TADAM: Task dependent adaptive metric for improved few-shot learning", Oreshkin et. al. 18



## Meta-training based metric learning

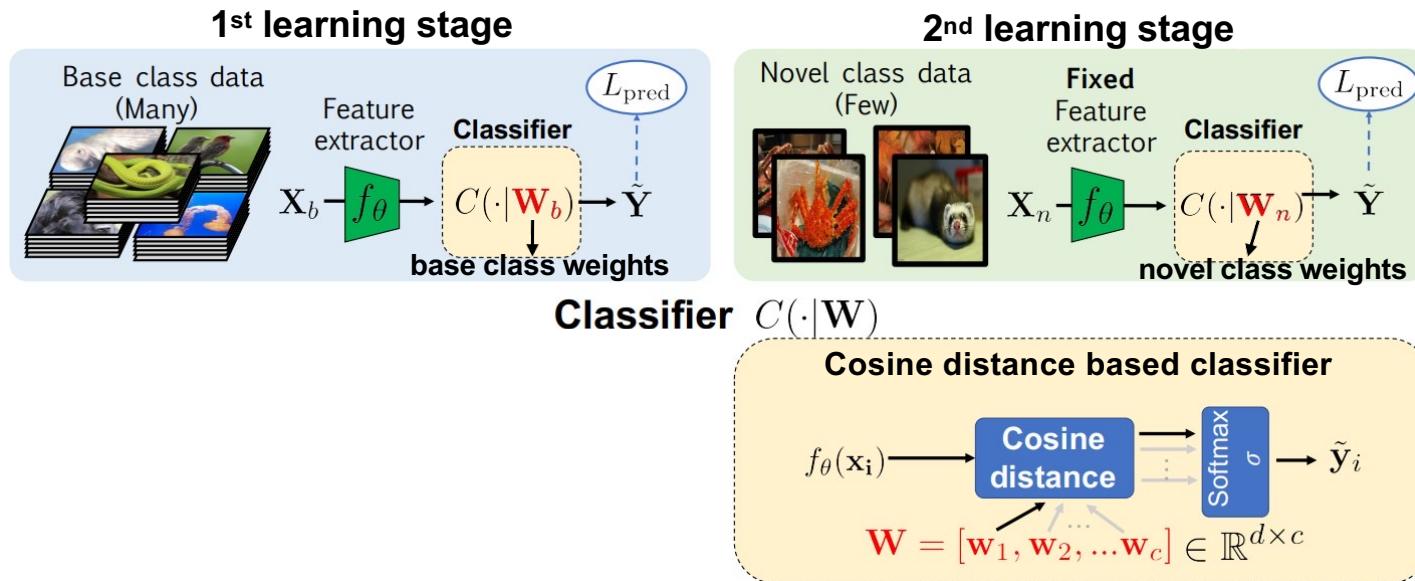
- In general: simple and effective methods
- But, meta-training can be bothersome:
  - Train a different metric function for each possible K or N
  - For small N → training with easy examples
  - Not all methods follow this rule, but then, how to tune N, K and M?

## Meta-training based metric learning

- In general: simple and effective methods
- But, meta-training can be bothersome
- Is meta-training really necessary for learning good features?

# Cosine distance based classification network

- Train typical classification network: **feature extractor + classification head**
- **Classification head:** replace dot-product (i.e., linear layer) with cosine distance



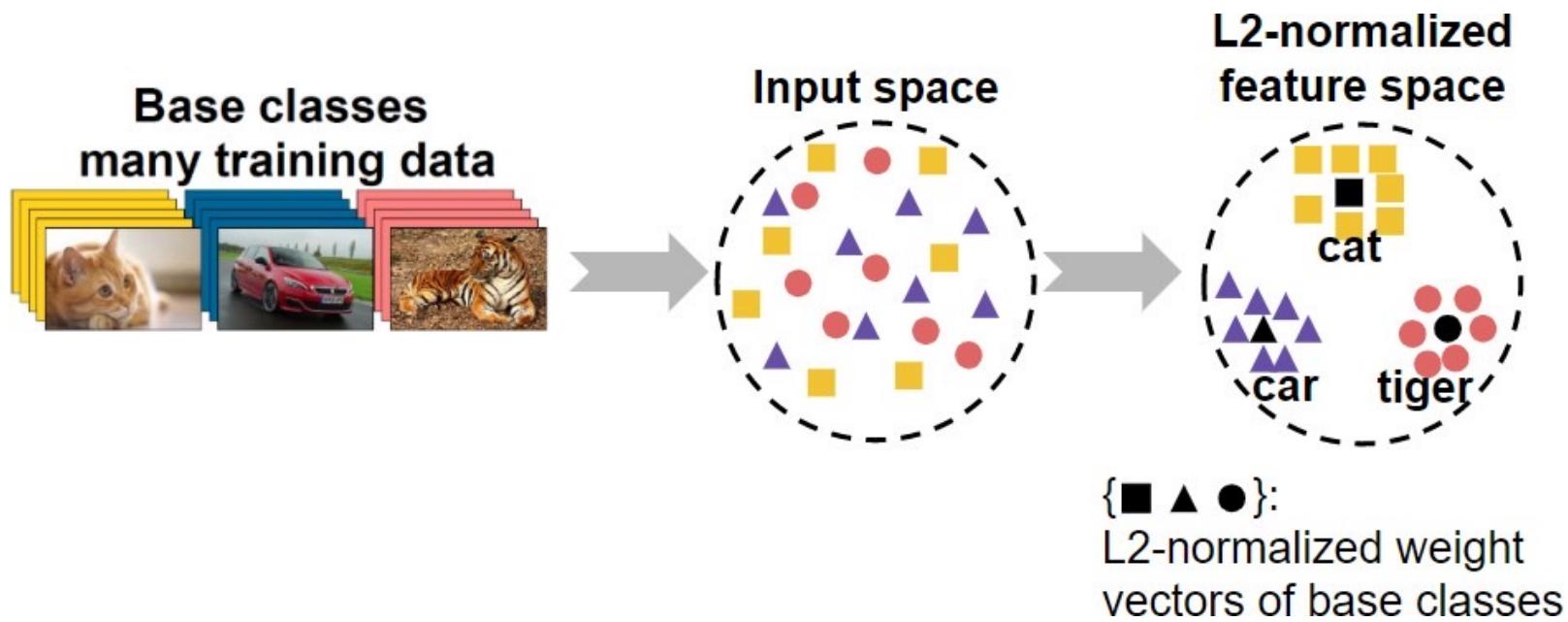
“Dynamic Few-Shot Visual Learning without Forgetting”, Gidaris et al. 2018

“Low-Shot Learning with Imprinted Weights”, Qi et al. 2018

## Why distance based classification head?

Enforces **similar behavior as metric learning models**:

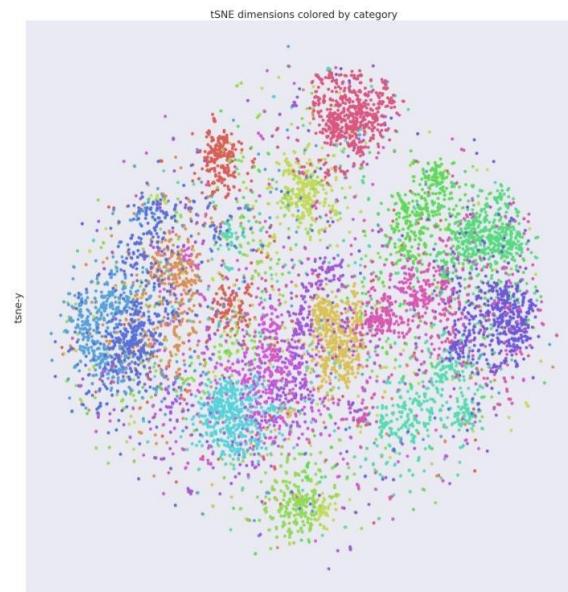
- Given an image, the learned feature must maximize (minimize) cosine similarity with weight vector of the correct class (incorrect classes)



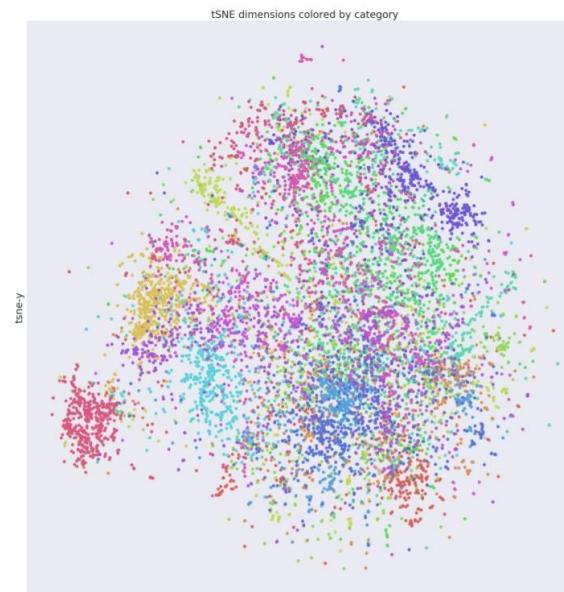
## Why distance-based classification head?

Enforces **similar behavior** as metric learning models

- Learn features with **reduced intra-class variance** →
- **Better generalization to novel classes**



(a) Cosine-similarity based features of novel categories



(b) Dot-product based features of novel categories

Source:  
Gidaris et al. 2018

## Cosine distance based classification network

- **1<sup>st</sup> learning stage:** standard training using the base class data
  - Trains the extractor  $f_\theta$  and classification weights  $\mathbf{W}_b$  of base classes
  - Much simpler than meta-training based metric methods

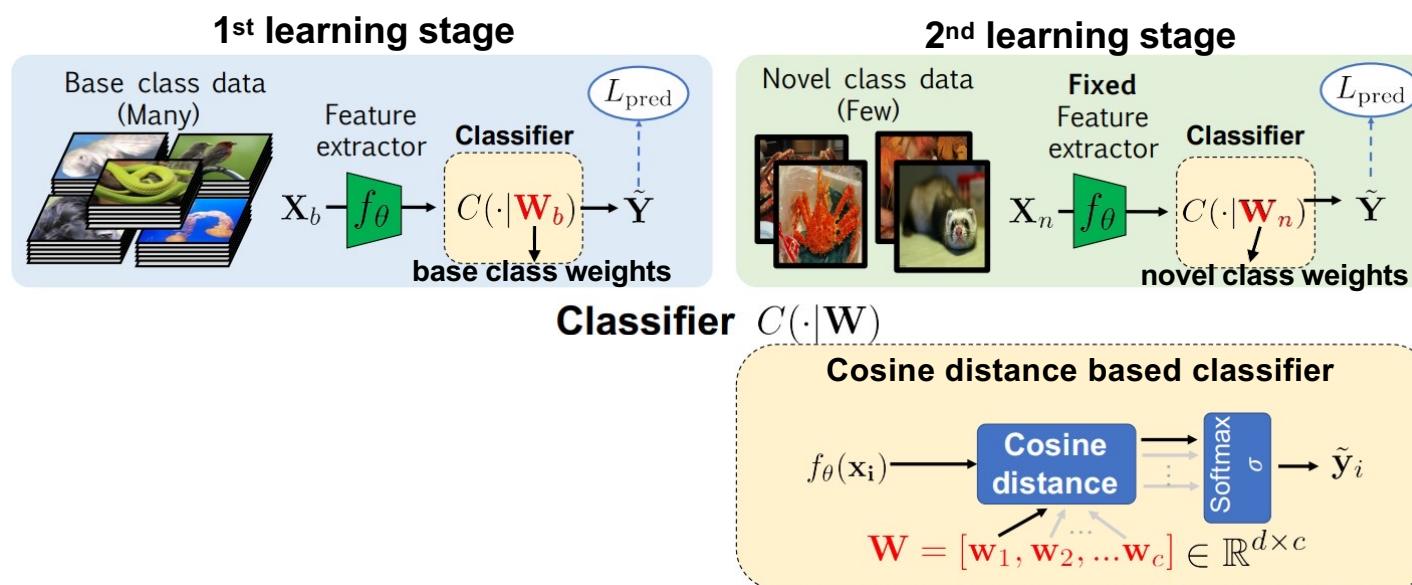


Image source (modified):  
Chen et al. 2019

## Cosine distance based classification network

- **2<sup>nd</sup> stage:** fix extractor  $f_\theta$  + “learn” only the classification weights  $W_n$ :
  - **compute  $W_n$  with prototypical feature averaging**

$$w_i = \frac{1}{|S_i|} \sum_{(x_k^S, y_k^S) \in S_i} f_\theta(x_k^S), \forall w_i \in W_n$$

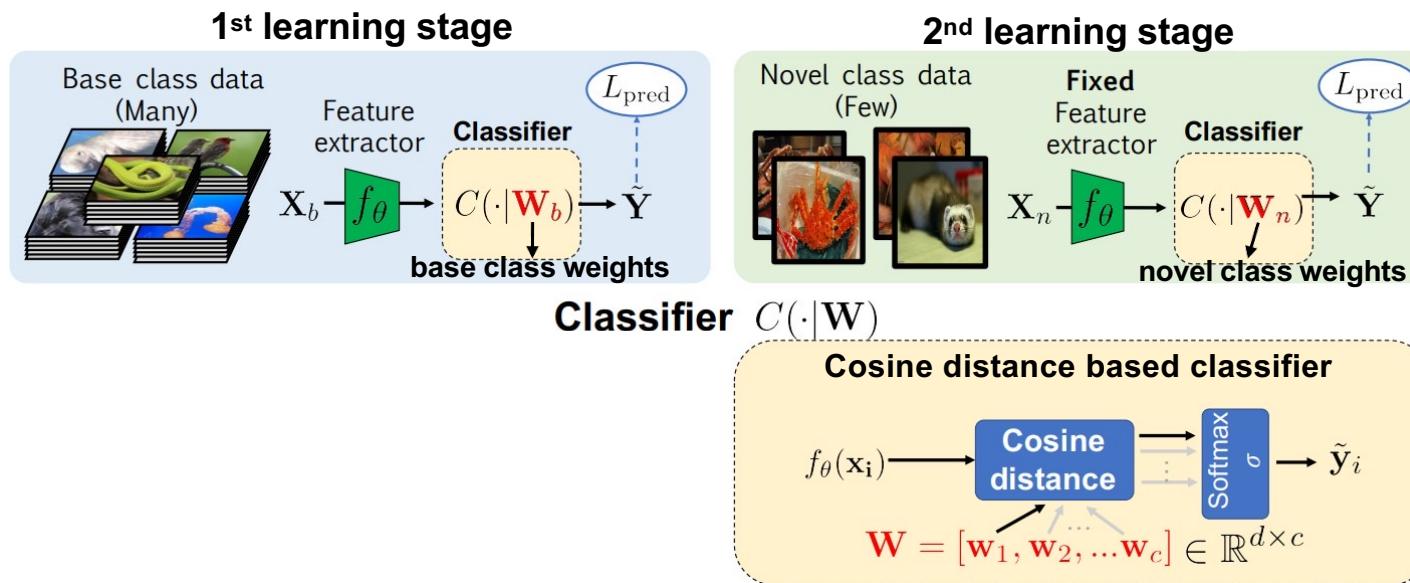


Image source (modified):  
Chen et al. 2019

## Cosine classifier

Models	1-Shot	5-Shot
Matching-Nets [26]	$55.53 \pm 0.48\%$	$68.87 \pm 0.38\%$
Prototypical-Nets [23]	$54.44 \pm 0.48\%$	$72.67 \pm 0.37\%$
Cosine Classifier	$54.55 \pm 0.44\%$	$72.83 \pm 0.35\%$

**Table 1:** 5-way accuracies on MinilmageNet.

**Simpler training with better results than Matching and Prototypical Nets**

Approach	$K=1$	2	5	10	20
<i>Prior work</i>					
Prototypical-Nets	39.3	54.4	66.3	71.2	73.9
Matching Networks	43.6	54.0	66.0	72.5	76.9
Cosine Classifier	45.23	56.90	68.68	74.36	77.69

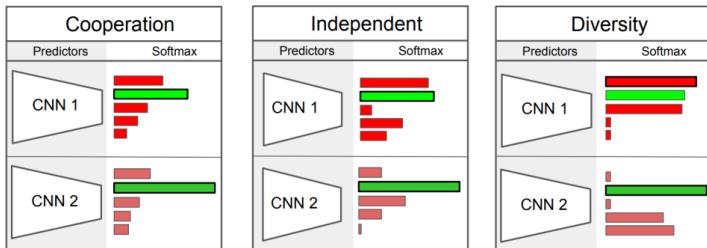
**Table 2:** 311-way accuracies on ImageNet-FS for  $K=1, 2, 5, 10$ , or 20 examples per novel class.

Source: "Dynamic Few-Shot Visual Learning without Forgetting", Gidaris et al. 18

# Cosine classifier

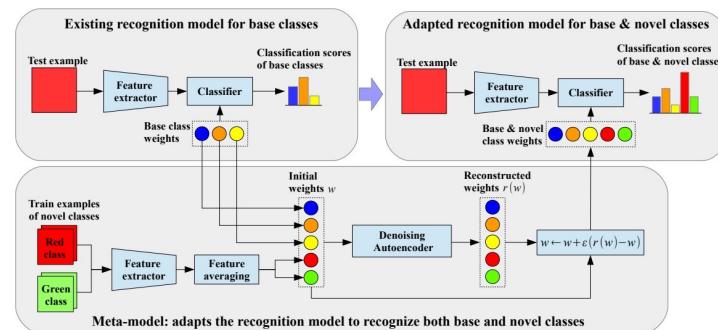
## Learn an ensemble of cosine classifiers

“Diversity with cooperation: ensemble methods for few-shot classification”, Dvornik et al. 18



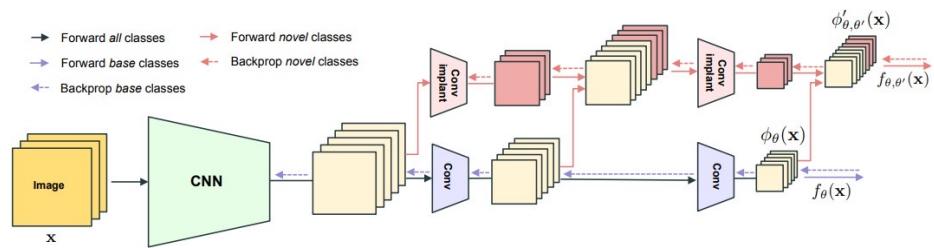
## Learn to predict class prototypes for pre-trained cosine classifiers

Gidaris et al. 18, Gidaris et al. 19

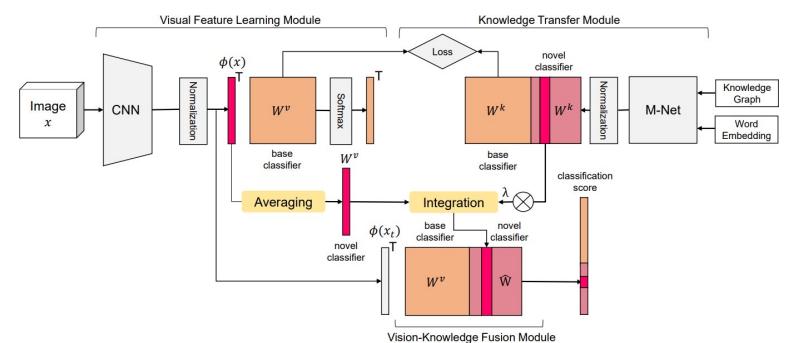


## Dense (cosine-based) classification & implanting new task-specific layers

“Dense classification and implanting for few-shot learning”, Lifchitz et al. 19



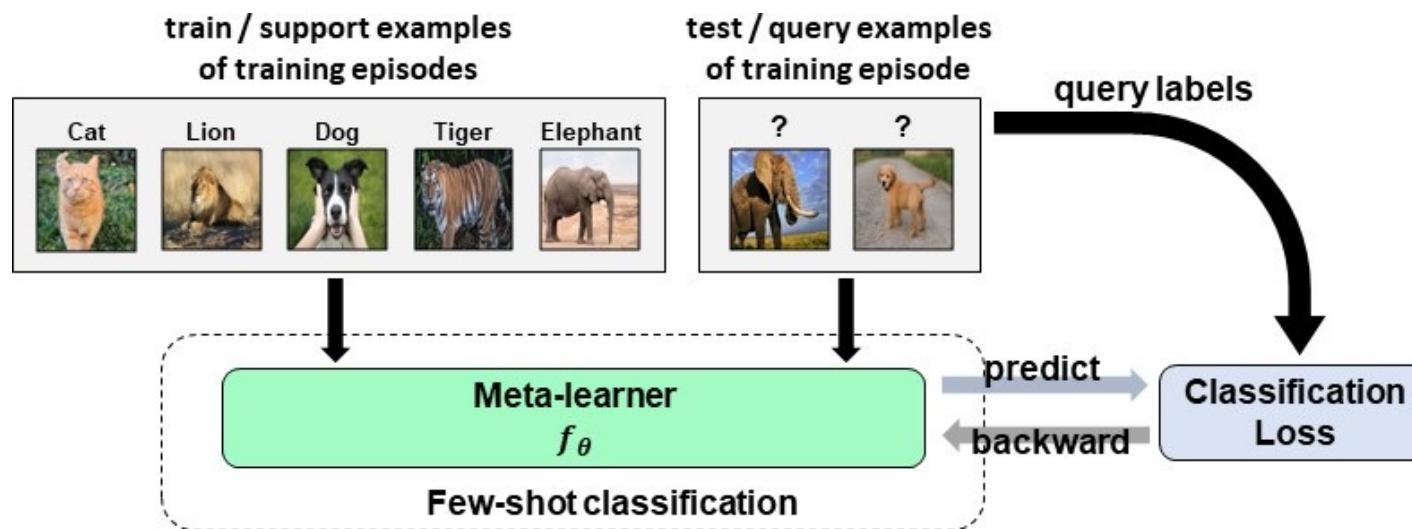
## Learn to predict class prototypes for cosine classifiers leveraging word embedding based knowledge graphs Peng et al. 19



# Agenda

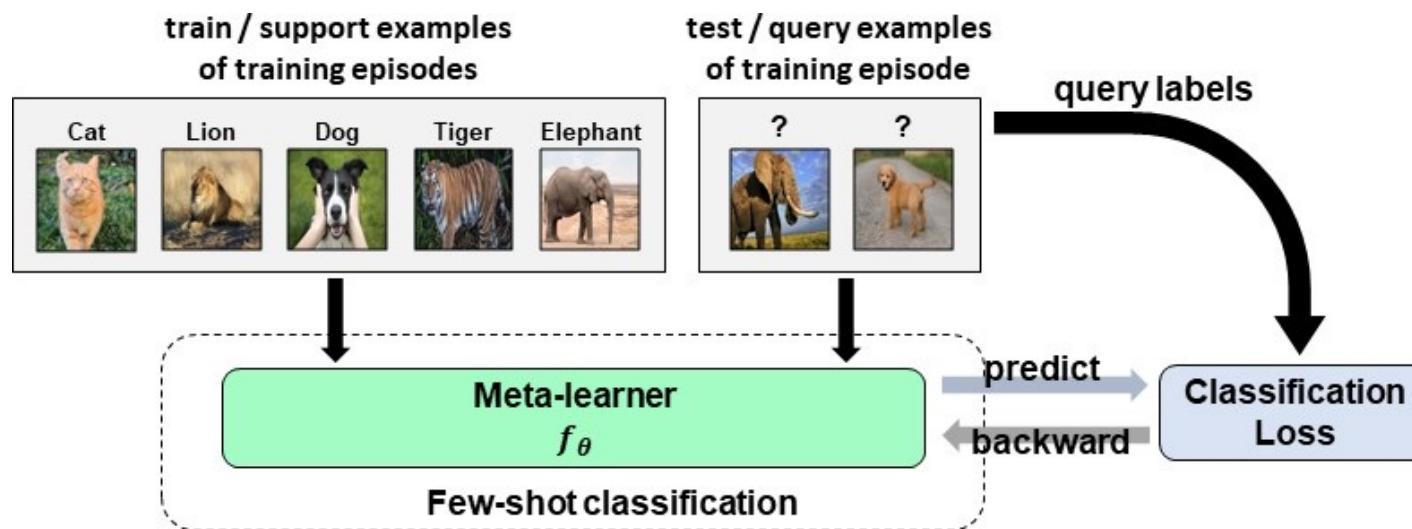
- Introduction
- Main types of few-shot learning algorithms
  - Metric learning
  - **Meta-learning with memory modules**
  - Learn to predict model parameters
- Few-shot learning without forgetting

## Meta-learning with memory modules



- **Few-shot classification:**
  - **input:** labeled support data, unlabeled query data
  - **intermediate step:** generate model
  - **output:** predicted query labels

## Meta-learning with memory modules

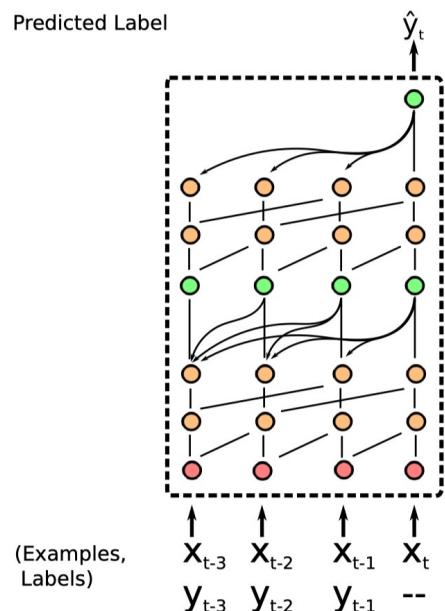


- **Few-shot classification:**
  - **input:** labeled support data, unlabeled query data
  - **intermediate step:** generate model → store support data to memory
  - **output:** predicted query labels by accessing the memory

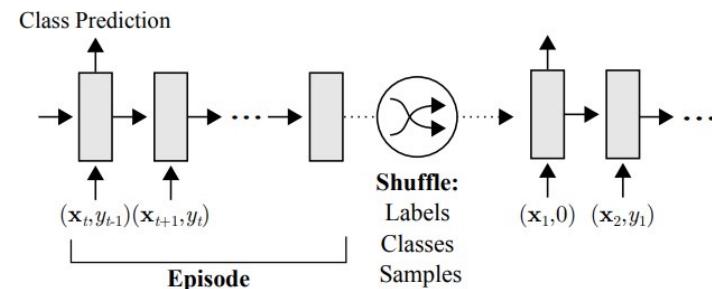
Treats **few-shot classification** as a “black box” prediction problem

# Meta-learning with memory modules

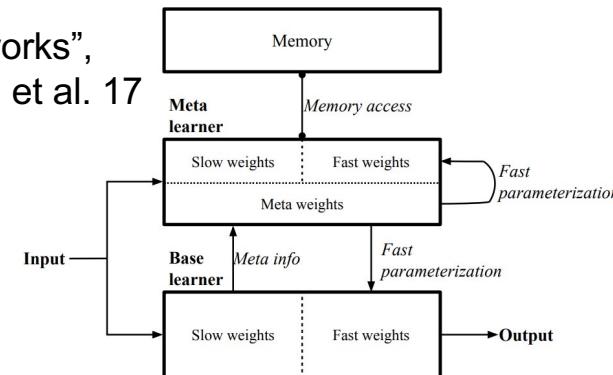
“A Simple Neural Attentive Meta-Learner”,  
Mishra et al. 18



“Meta-Learning with Memory-Augmented Neural Networks”, Santoro et al. 16

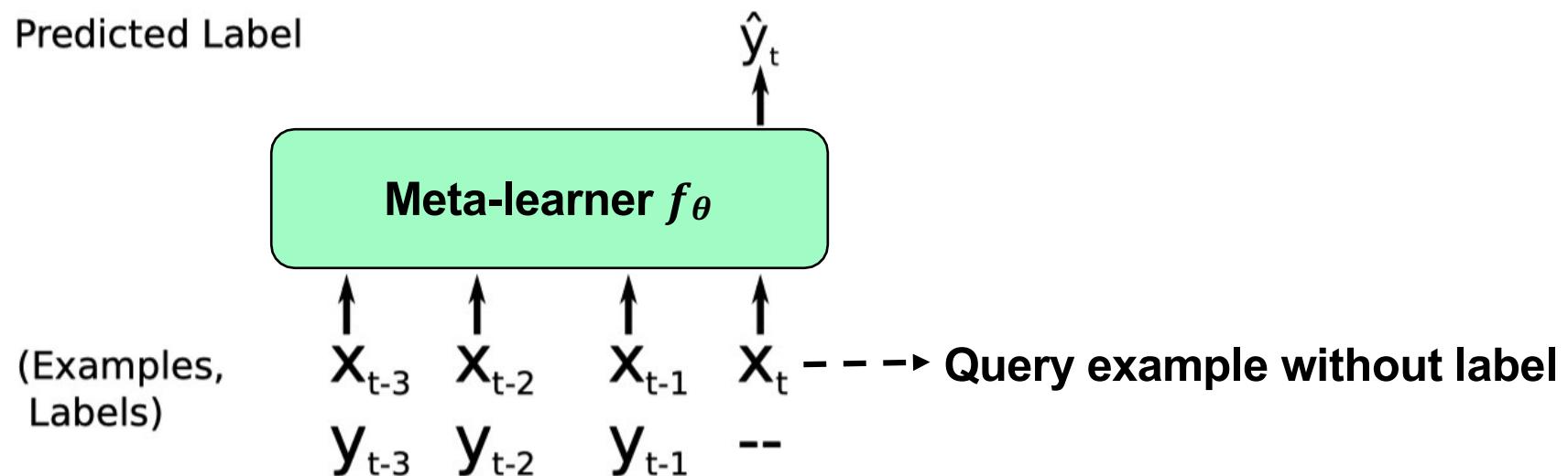


“Meta Networks”,  
Munkhdalai et al. 17



## Example: Simple Neural Attentive Meta-Learner

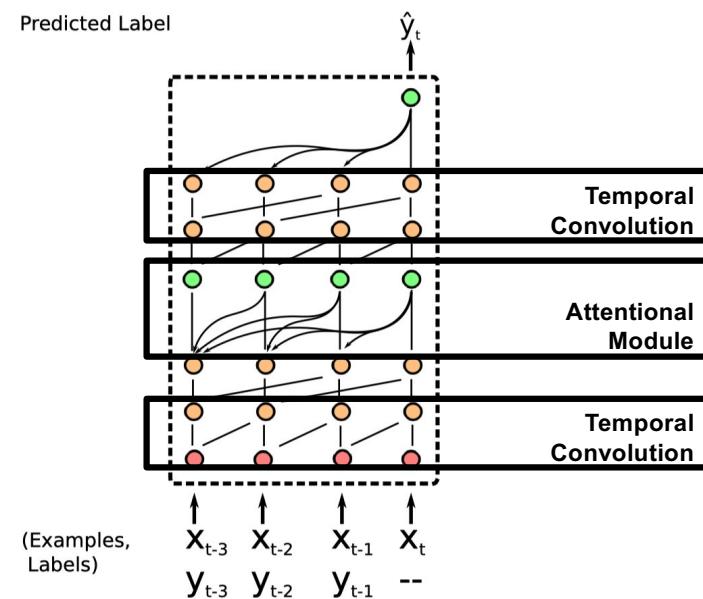
- Few-shot as a sequence labeling task:
  - Given past labeled images, what is the label of the current query image



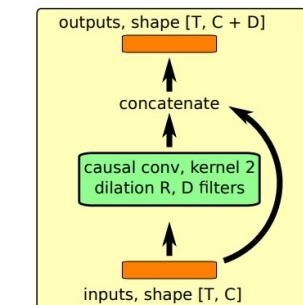
“A Simple Neural Attentive Meta-Learner”, Mishra et al. 18

## Example: Simple Neural Attentive Meta-Learner

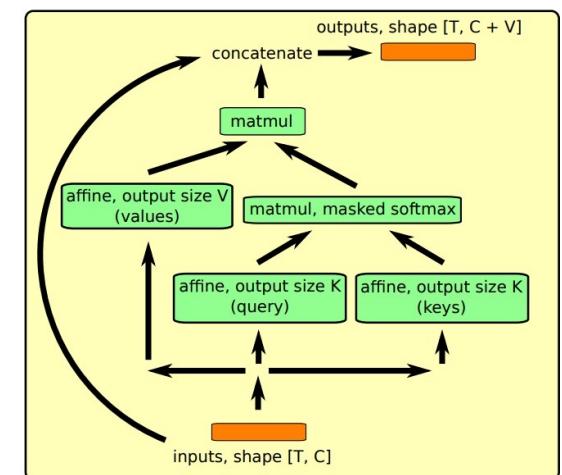
- **Meta-learner implementation:**
  - **Temporal convolutions:** aggregates past information
  - **Attentional Module:** pinpoints to query-specific past information
    - “Attention is all you need”, Vaswani et al. 17



**Temporal  
Conv. Module**

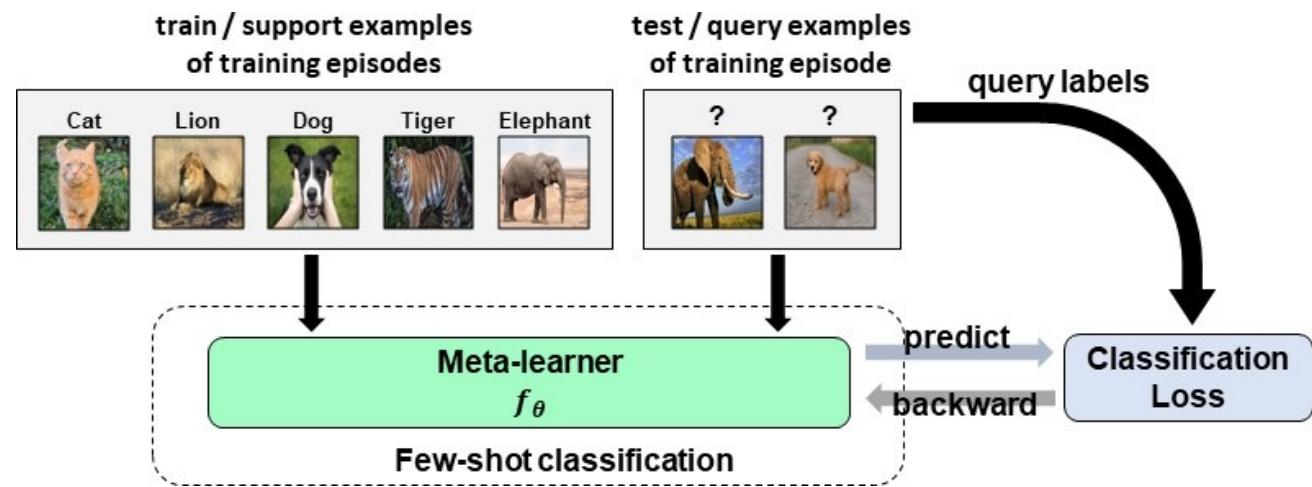


**Attention Module**



“A Simple Neural Attentive Meta-Learner”, Mishra et al. 18

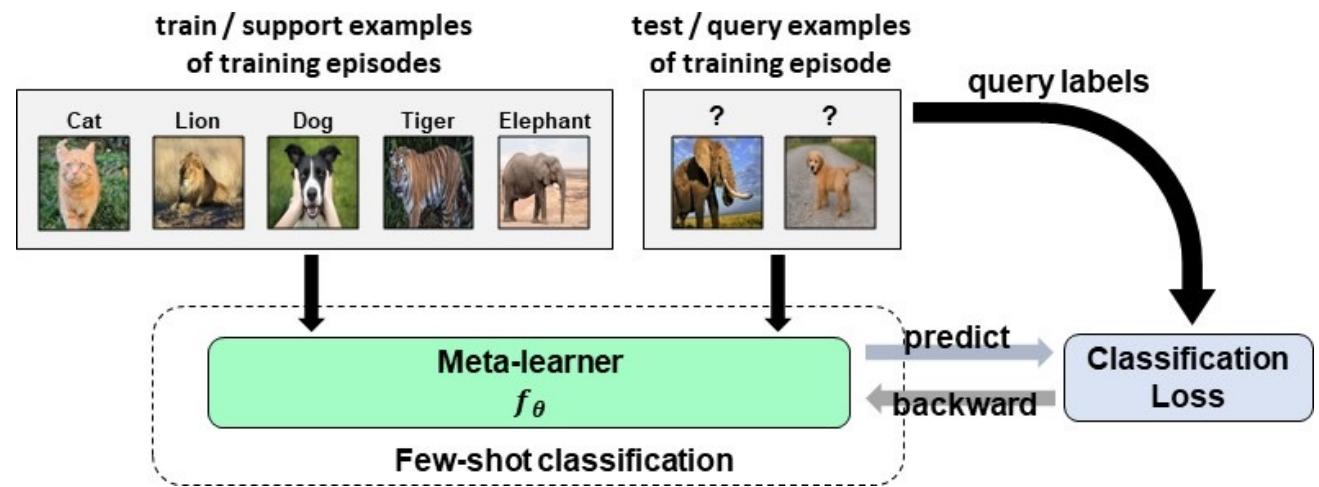
## Meta-learning with memory modules



### Meta-training routine:

1. Sample training episode  $(S, Q)$
2. Generate classification model  $m_\varphi = f_\theta(S)$
3. Predict classification scores  $p_m = m_\varphi(x_m^Q)$
4. Optimize  $\theta$  w.r.t. the query classification loss  $L(f_\theta(\cdot | S), Q) = \sigma_m - \log(p_m[y_m^Q])$

# Meta-learning with memory modules



## Meta-training routine:

1. Sample training episode  $(S, Q)$
  2. Generate classification model  $m_\phi = f_\theta(S)$
  3. Predict classification scores  $\{p_m\}_m = f_\theta(\{x_m^Q\}_m | S)$  for all queries
  4. Optimize  $\theta$  w.r.t. the query classification loss  $L(f_\theta(\cdot | S), Q) = \sigma_m - \log(p_m[y_m^Q])$
- store & access support data with a memory module

## Meta-learning with memory modules

- More generic than metric learning methods
  - applicable to other learning problems: regression, RL, ...
- More data hungry (for training the meta-learner)
- More computational expensive

# Agenda

- Introduction
- Main types of few-shot learning algorithms
  - Metric learning
  - Meta-learning with memory modules
  - Learn to predict model parameters
- Few-shot learning without forgetting

# Agenda

- Introduction
- Main types of few-shot learning algorithms
  - Metric learning
  - Meta-learning with memory modules
  - **Learn to predict model parameters**
- Few-shot learning without forgetting

## Learn to predict model parameters

**Key idea:** train the meta-learner to predict task-specific model parameters

Usually, a **small subset of model parameters**:

- Predict diagonal of factorized weights:
  - “Learning feed-forward one-shot learners”, Bertinetto et al. 16
- Predict weights of classification head
  - “Learning to model the tail”, Wang et al. 17

## Learn to predict model parameters

**Key idea:** train the meta-learner to predict task-specific model parameters

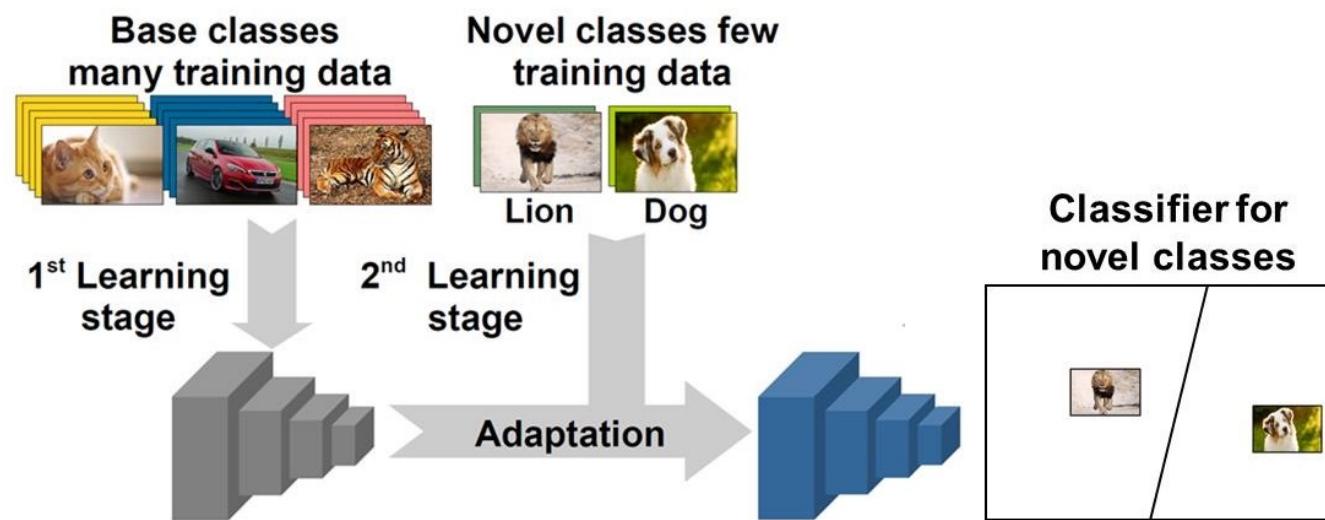
Here focus:

- **predicting the weights of the classification head**
- in the context of the “**few-shot learning without forgetting**” problem

# Agenda

- Introduction
- Main types of few-shot learning algorithms
  - Metric learning
  - Meta-learning with memory modules
  - Optimization based meta-learning
  - Learn to predict model parameters
- **Few-shot learning without forgetting**

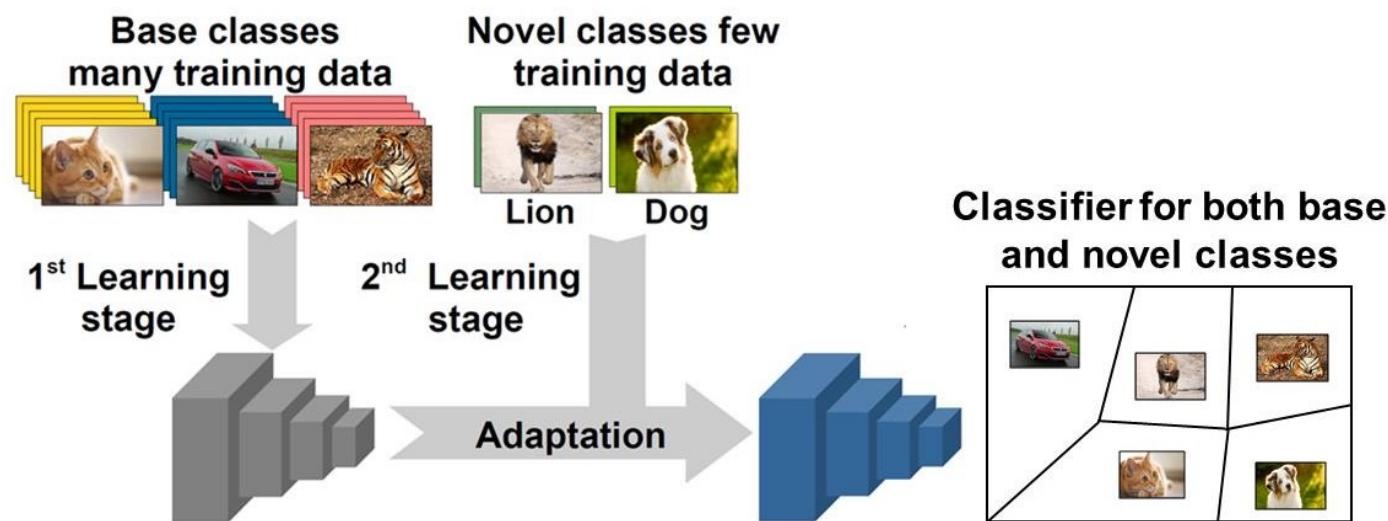
## Few-shot learning without forgetting



### Typical few-shot models:

- focus on learning novel classes with limited data
- but “forget” the initial base classes ☹
  - “forget”: worse than base class models or unable to recognize base classes

## Few-shot learning without forgetting



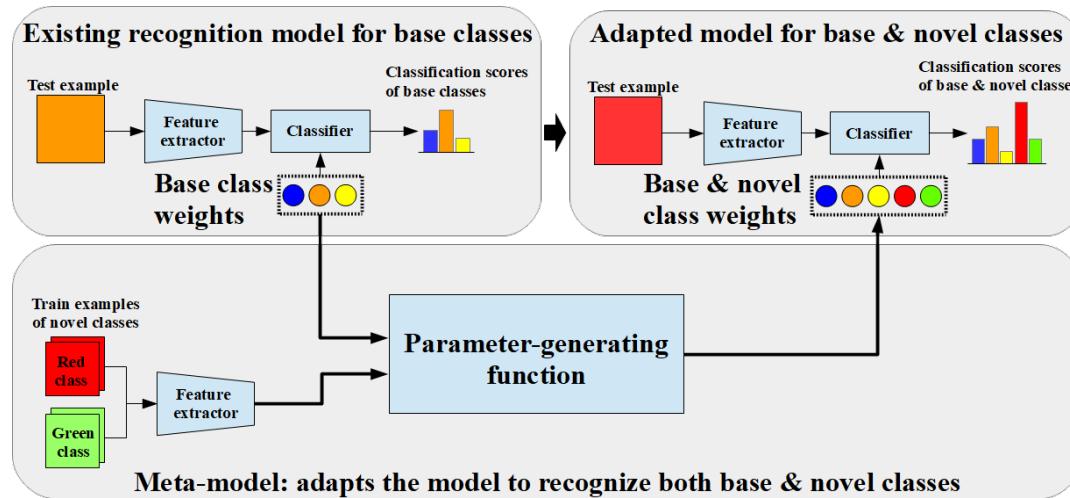
- In contrast, **practical applications** often require:
  - to **extend base classes with novel ones** using few training data
  - and **without re-training** on the full dataset (base+novel)
- “**Few-shot learning without forgetting**” targets this problem
  - combines elements from both incremental and few-shot learning

## Agenda

- Introduction
- Main types of few-shot learning algorithms
  - Metric learning
  - Meta-learning with memory modules
  - Optimization-based meta-learning
    - **Learn to predict classification weights**
- **Few-shot learning without forgetting**

The description of the “Learn to predict classification weights” methods is in the context of the “few-shot learning without forgetting” setting.

# Learn to generate classification weights



- **Pre-trained network:** feature extractor + cosine classification head
- **Extend with parameter-generating function:**
  - **outputs:** new weights for the novel classes

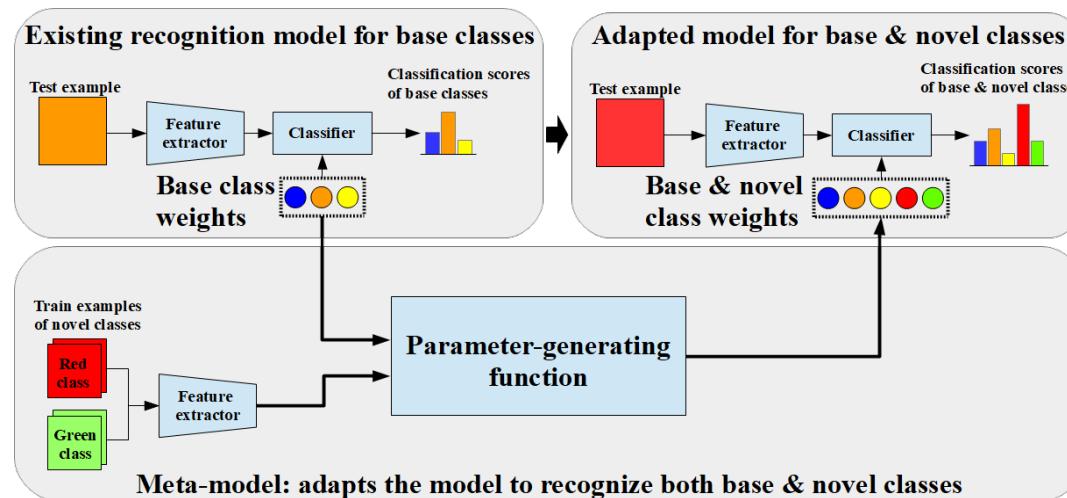
“Dynamic Few-Shot Visual Learning without Forgetting”, Gidaris et al. 18

“Low-Shot Learning with Imprinted Weights”, Qi et al. 18

“Few-Shot Image Recognition by Predicting Parameters from Activations”, Qiao et al. 18

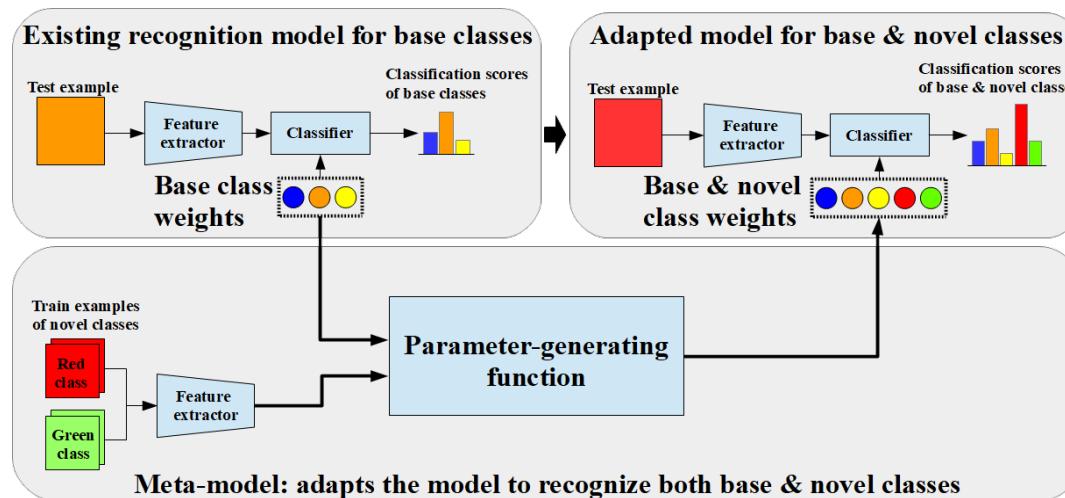
“Learning to model the tail”, Wang et al. 17

# Learn to generate classification weights



- **Important to use cosine classification head:**
  - **$L_2$ -normalize weights:** all classes have same  $L_2$  norms
    - **avoids class imbalance:** biasing towards classes with bigger  $L_2$ norms
  - Easier to add novel weights → unified recognition of both type of classes

# Learn to generate classification weights



- **Important to use cosine classification head:**
- Beneficial in the traditional incremental learning setting as well:
  - “Learning a unified classifier incrementally via rebalancing”, Hou et al. 19
  - “Memory efficient incremental through feature adaptation” Iscen et al. 20

## Learn to generate classification weights

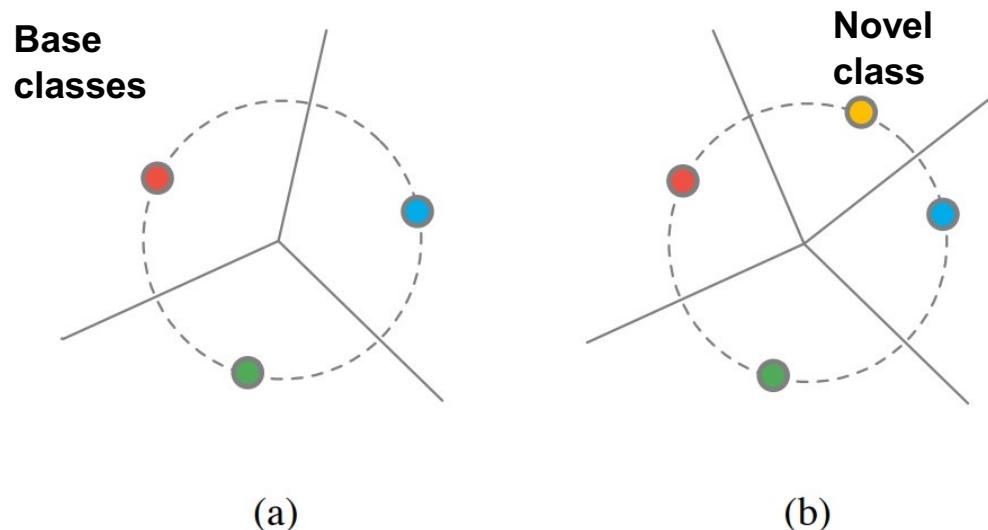


Figure 2. Illustration of imprinting in the normalized embedding space. (a) Before imprinting, the decision boundaries are determined by the trained weights. (b) With imprinting, the embedding of an example (the yellow point) from a novel class defines a new region.

Source: "Low-Shot Learning with Imprinted Weights", Qi et al. 18

## Meta-training in few-shot learning without forgetting

### Meta-training routine:

1. Sample training episode  $(S, Q)$
2. Generate classification model  $m_\varphi = f_\theta(S)$
3. Classification scores  $p_m = m_\varphi(x_m^Q)$
4. Optimize  $\theta$  w.r.t. the query classification loss  $L(f_\theta(S), Q) = \sigma_m - \log(p_m[y_m^Q])$

## Meta-training in few-shot learning without forgetting

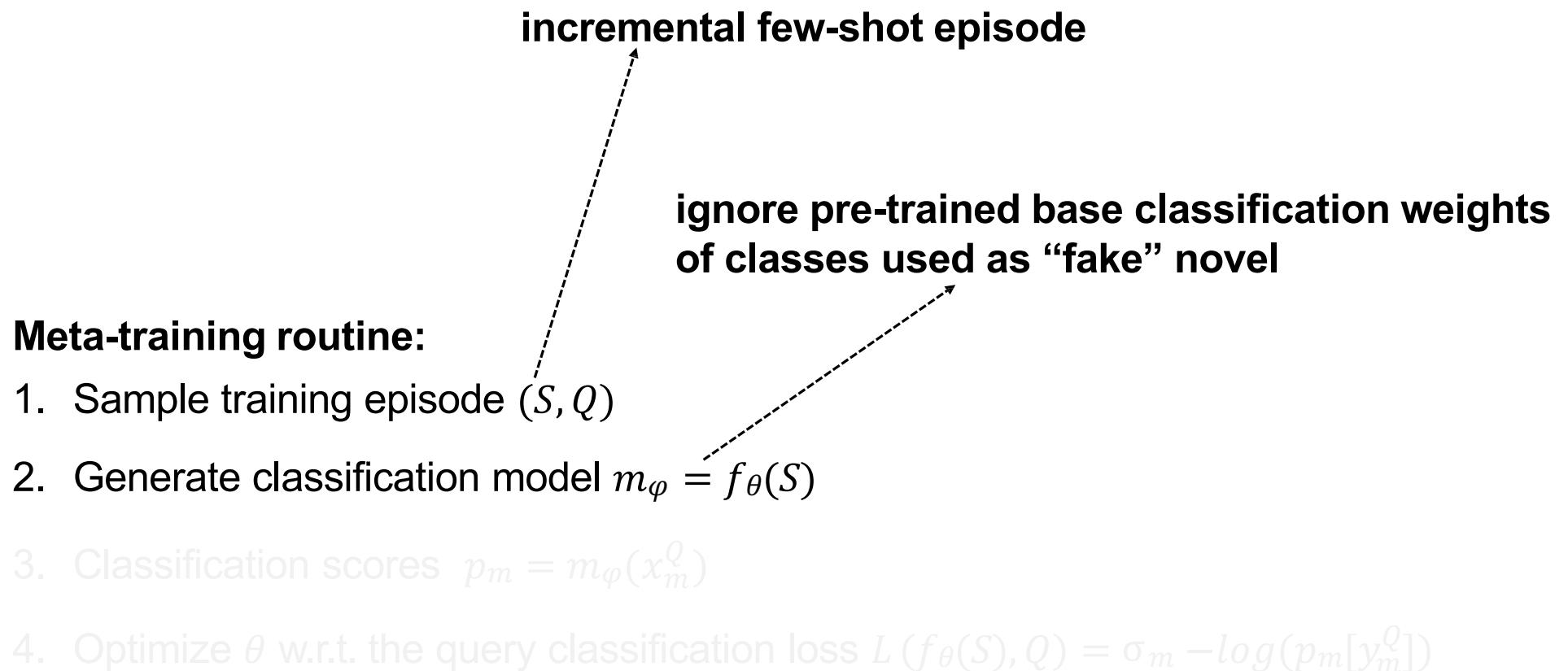
**incremental few-shot episode:**

- randomly choose some base classes as “fake” novel
- $S$ : examples from the “fake” novel classes
- $Q$ : examples form both “fake” novel and remaining base

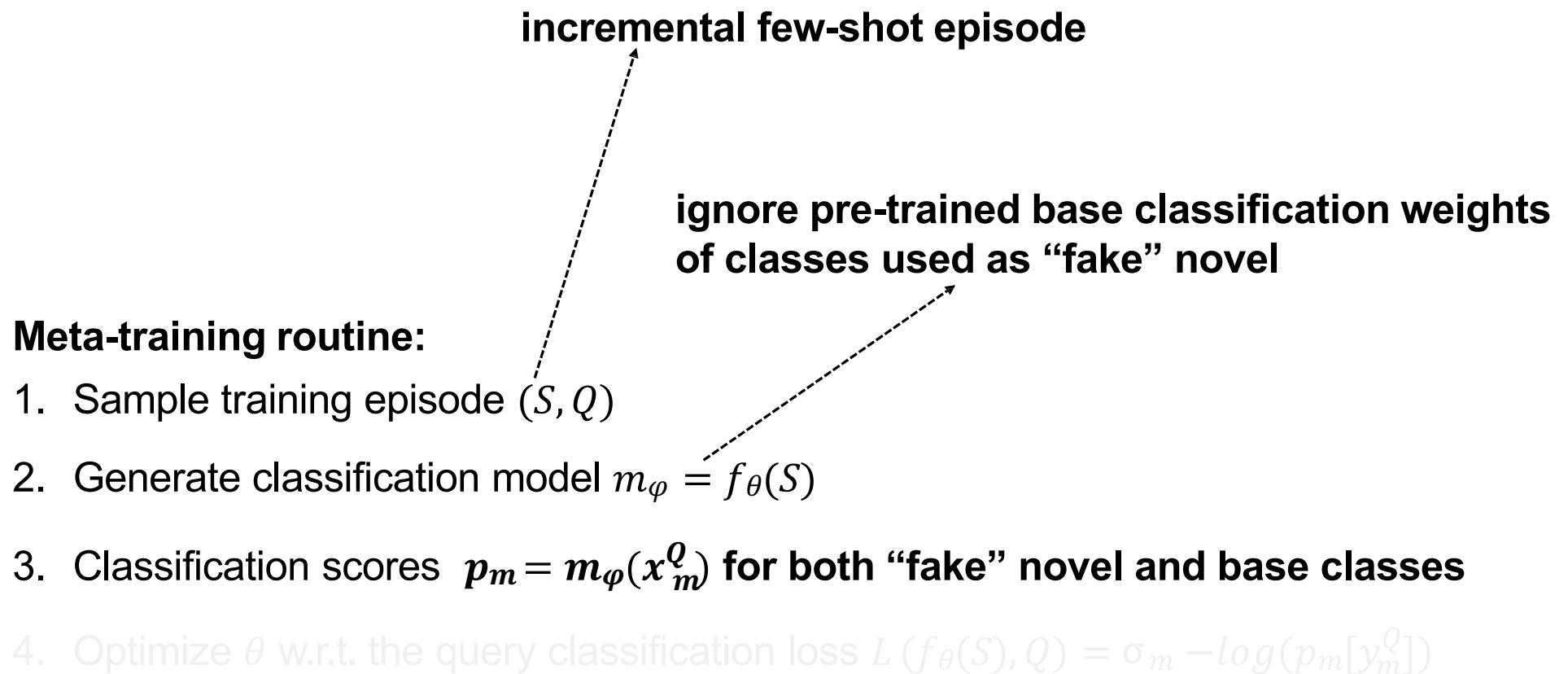
**Meta-training routine:**

1. Sample training episode  $(S, Q)$
2. Generate classification model  $m_\varphi = f_\theta(S)$
3. Classification scores  $p_m = m_\varphi(x_m^Q)$
4. Optimize  $\theta$  w.r.t. the query classification loss, e.g.:  $L(f_\theta(S), Q) = \sigma_m - \log(p_m[y_m^Q])$

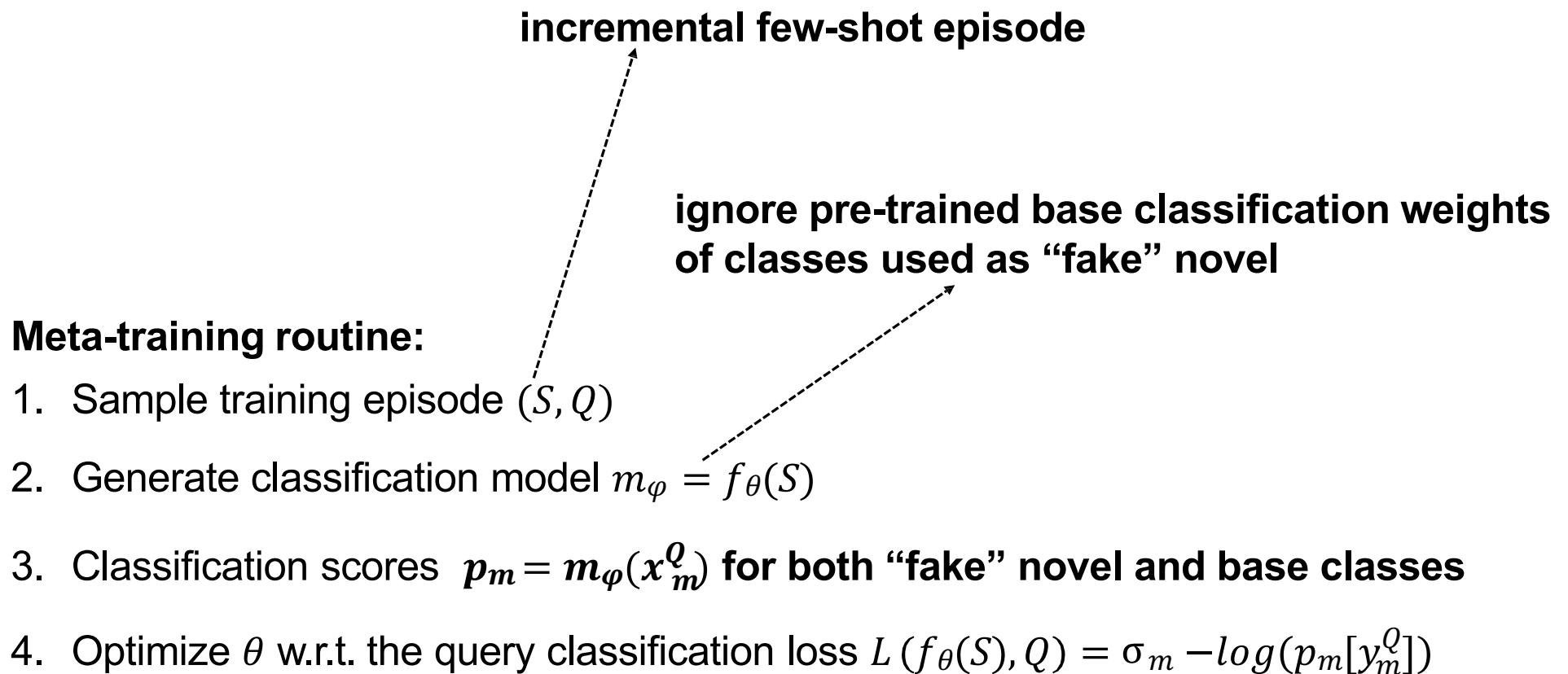
## Meta-training in few-shot learning without forgetting



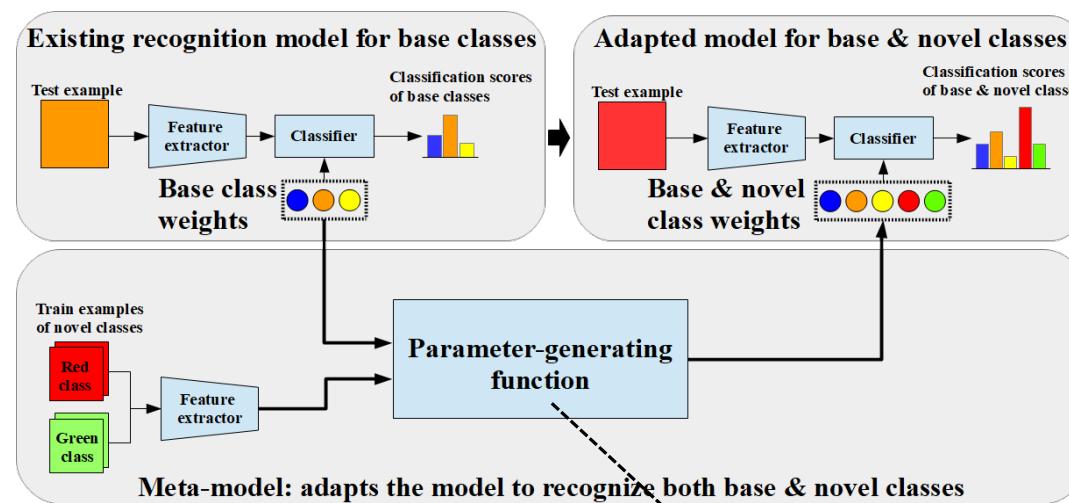
## Meta-training in few-shot learning without forgetting



## Meta-training in few-shot learning without forgetting



# Generate weights with prototypical feature averaging



## Simplest case:

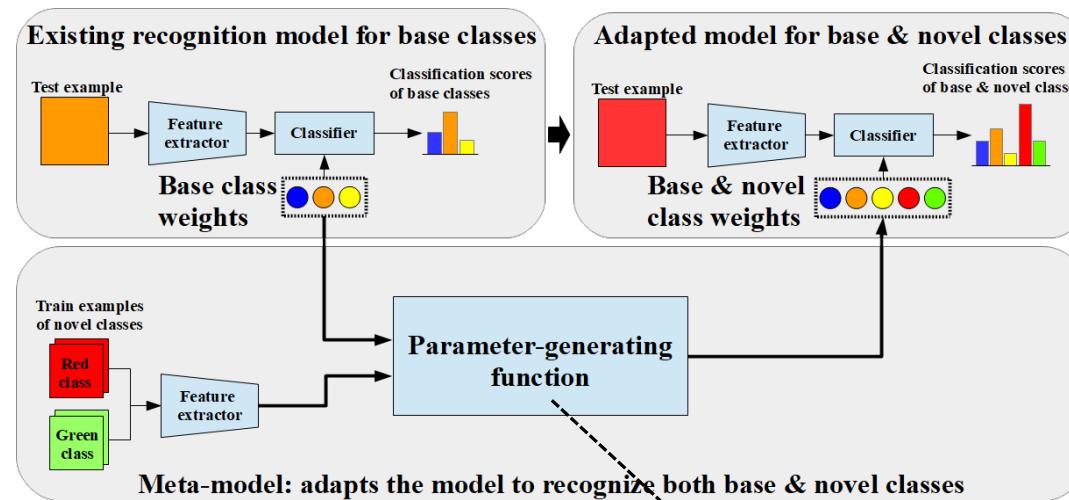
- $S_i$ : support set of  $i$ -th novel class
- **novel weight** = average feature vector of  $S_i$ :

$$w_i^{avg} = \frac{1}{|S_i|} \sigma_{(x_k^S, y_k^S) \in S_i} F_\theta(x_k^S)$$

"Low-Shot Learning with Imprinted Weights", Qi et al. 18

"Dynamic Few-Shot Visual Learning without Forgetting", Gidaris et al. 18

# Generate weights with prototypical feature averaging



## Simplest case:

- $S_i$ : support set of  $i$ -th novel class
- **novel weight** = average feature vector of  $S_i$ :

No meta-learning here

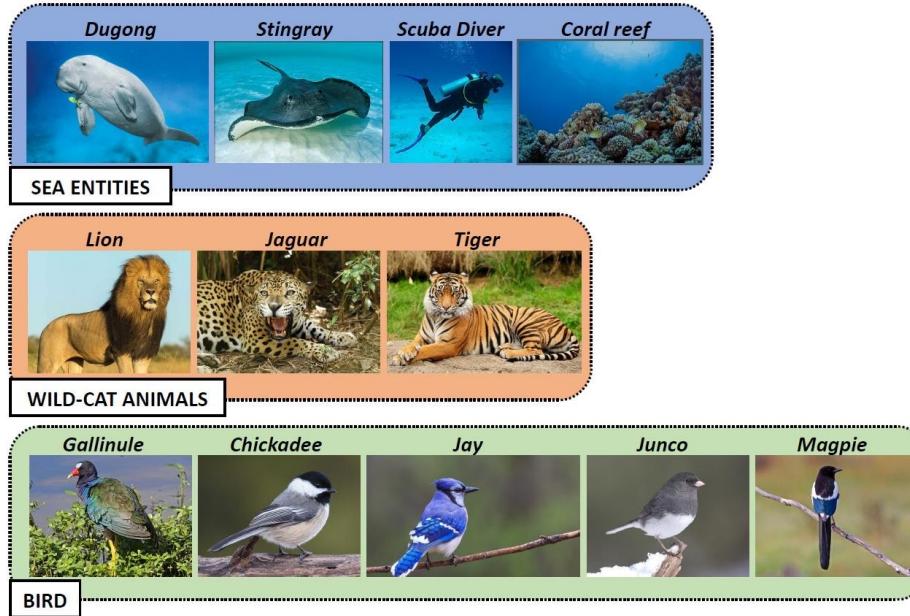
$$w_i^{avg} = \frac{1}{|S_i|} \sigma_{(x_k^S, y_k^S) \in S_i} \theta(x_k^S)$$

"Low-Shot Learning with Imprinted Weights", Qi et al. 18

"Dynamic Few-Shot Visual Learning without Forgetting", Gidaris et al. 18

$F$

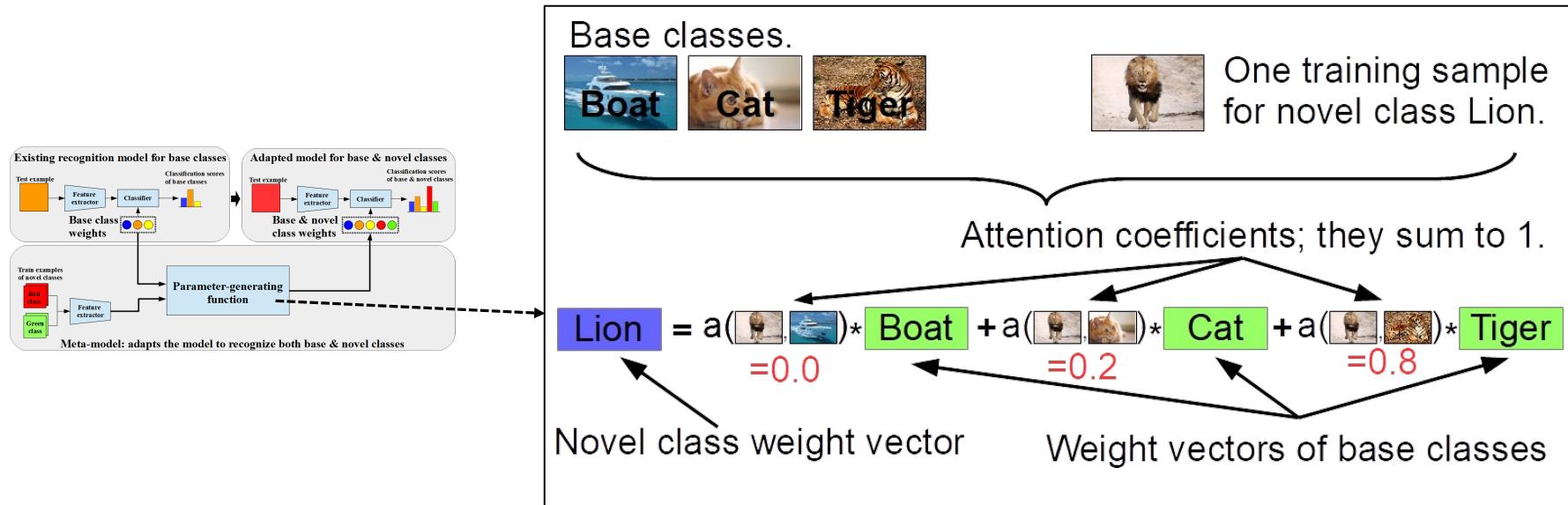
# Correlations between classification weights



Many classes are semantically / visually related:

- **Correlations between their classification weights**
- Exploit those correlations for generating novel class weights?

# Generate classification weights with attention module



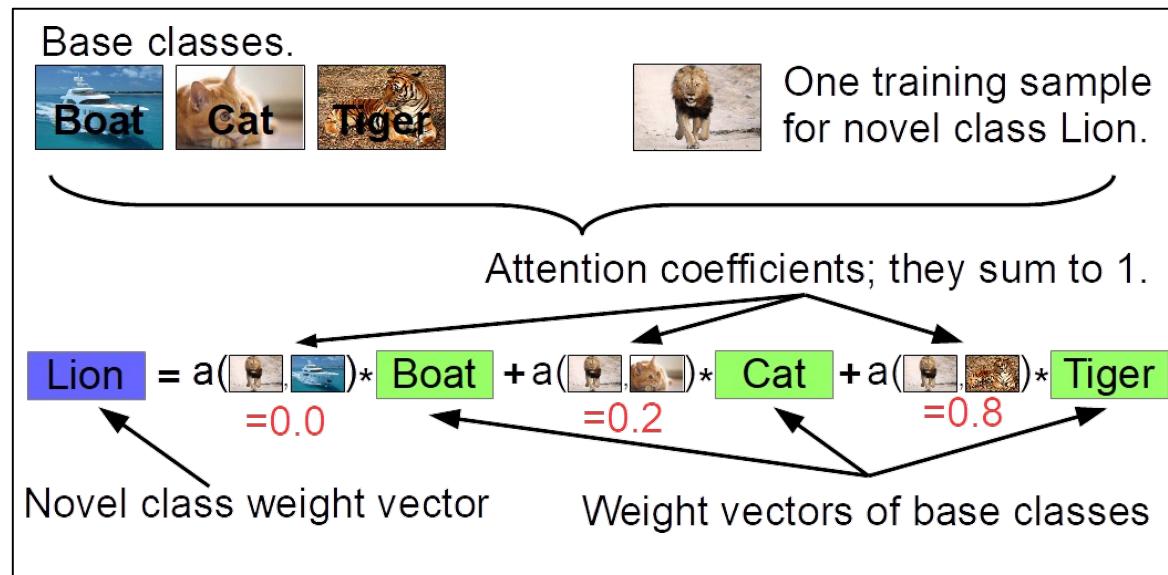
- Novel weight using attention over base weights  $w_b$ :

$$w_i^{att} = \sigma_{b=1}^{N_b} a(S_i)[b] \cdot w_b$$

- $N_b$ : number of base classes

“Dynamic Few-Shot Visual Learning without Forgetting”, Gidaris et al. 18

## Generate classification weights with attention module



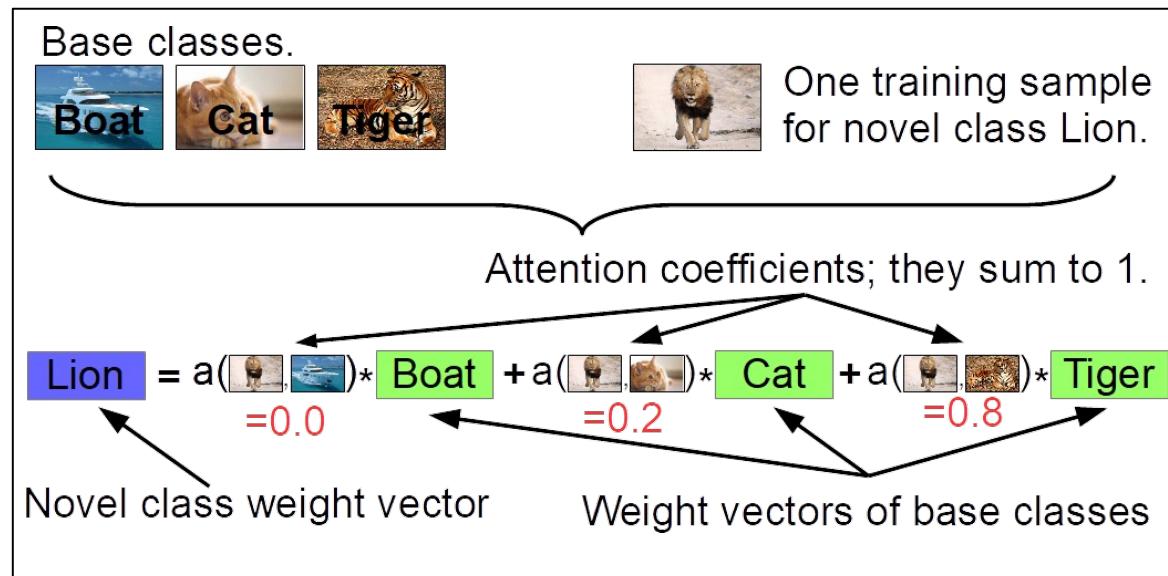
- Novel weight using attention over base weights  $w_b$ :

$$w_i^{att} = \sigma_{b=1}^{N_b} a(S_i)[b] \cdot w_b$$

- $N_b$ : number of base classes

“Dynamic Few-Shot Visual Learning without Forgetting”, S. Gidaris et al. 18

## Generate classification weights with attention module



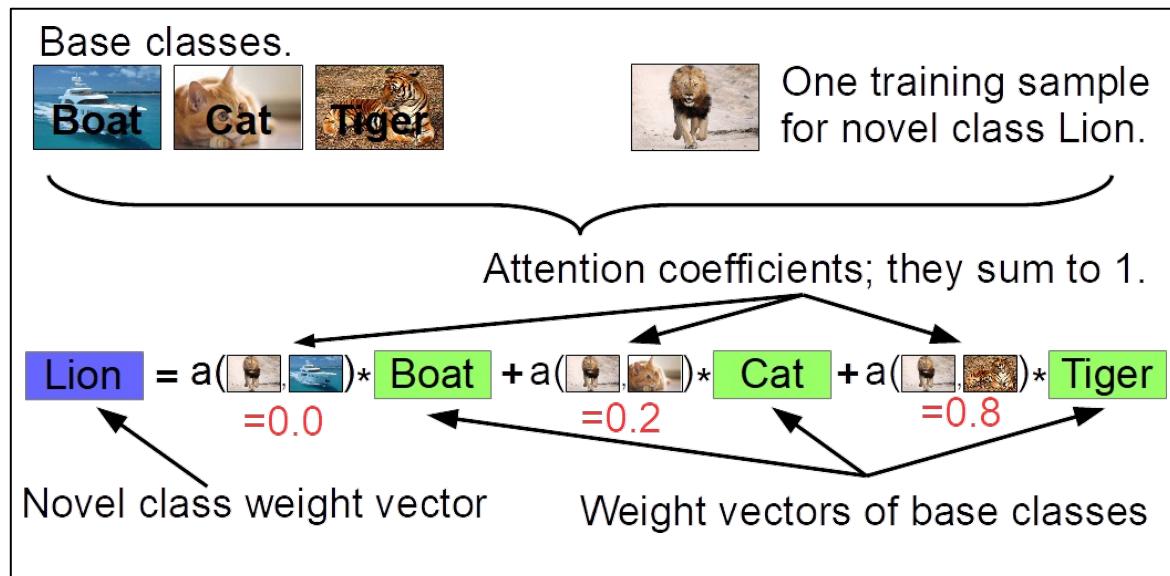
- Novel weight using attention over base weights  $w_b$ :

$$w_i^{att} = \sigma_{b=1}^{N_b} a(S_i)[b] \cdot w_b$$

- $a(S_i)[b]$ : average similarity of support features with base class weight  $w_b$ 
  - Computed with cosine + softmax

"Dynamic Few-Shot Visual Learning without Forgetting", Gidaris et al. 18

## Generate classification weights with attention module



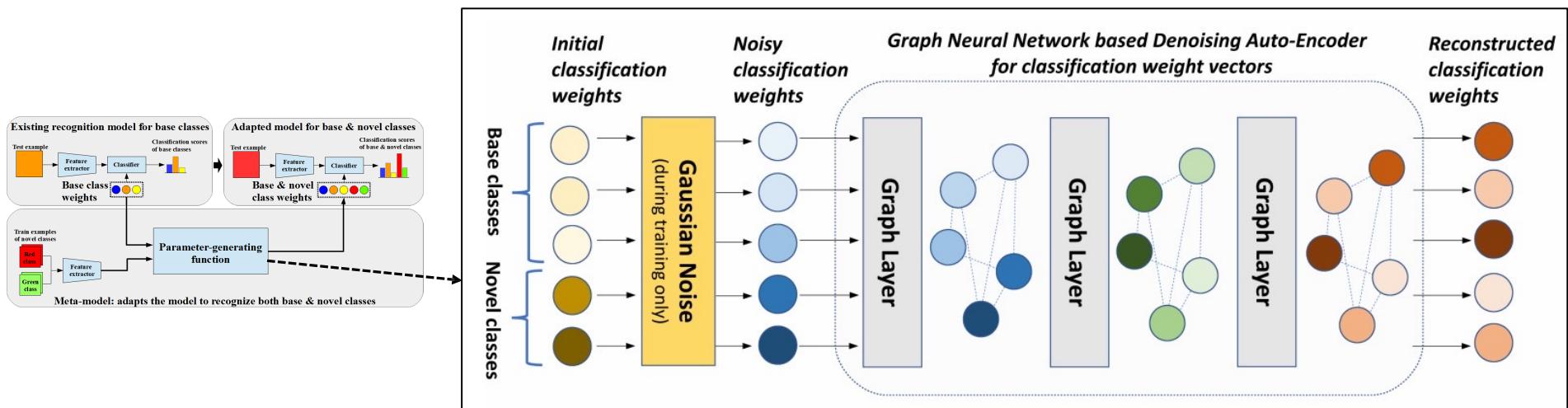
- Novel weight using attention over base weights  $w_b$ :

$$w_i^{att} = \sigma_{b=1}^{N_b} a(S_i)[b] \cdot w_b$$

- Final novel weight:  $w_i^{att}$  combined with prototypical averaging weight  $w_i^{avg}$

“Dynamic Few-Shot Visual Learning without Forgetting”, Gidaris et al. 18

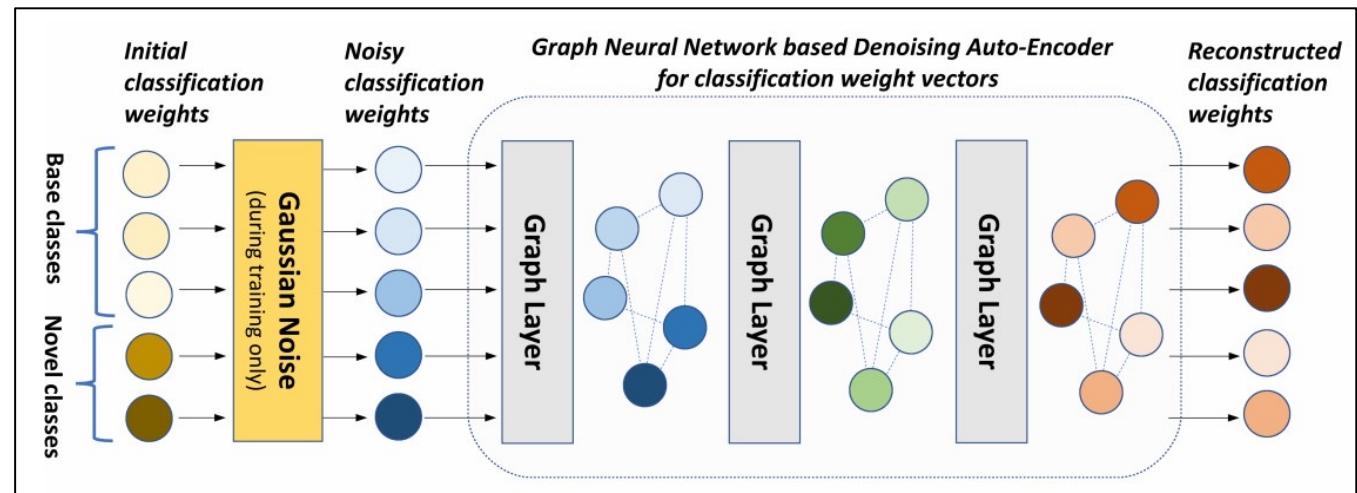
# Generate weights with a GNN Denoising AutoEncoder



- **Learning inter-class correlations with GNN based Denoising AutoEncoders**
  - Nodes = classes
  - Edges = each class connected to top most similar classes
  - **More expressive than a single layer attention mechanism**

"Generating classification weights with GNN denoising auto encoders for few-shot learning ", Gidaris et al. 19

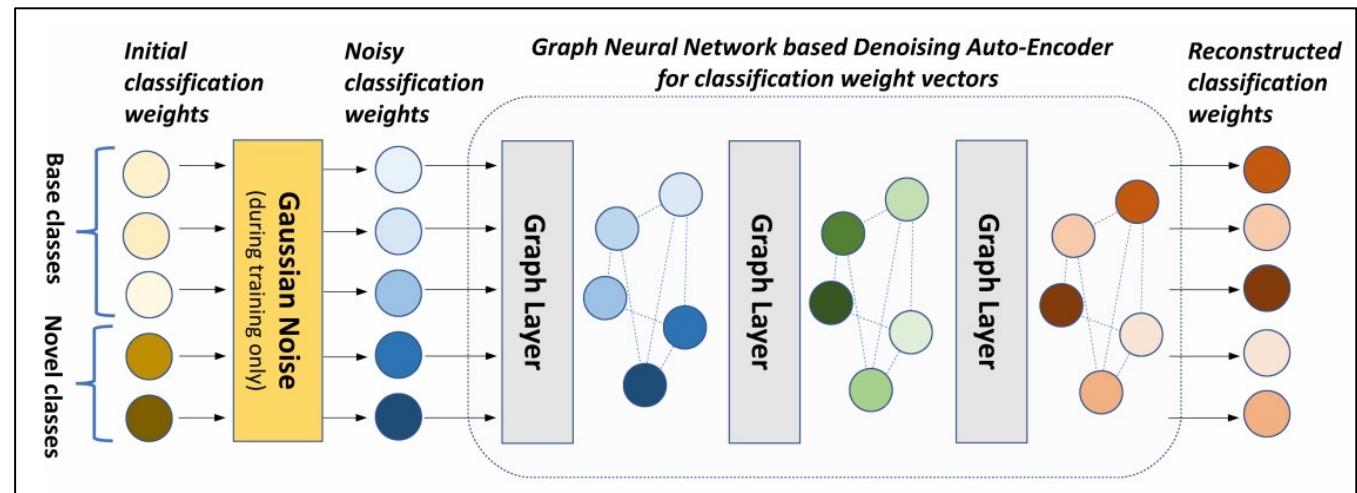
## Generate weights with a GNN Denoising AutoEncoder



- **Learning inter-class correlations with GNN based Denoising AutoEncoders**
  - Nodes = classes
  - Edges = each class connected to top most similar classes
  - **More expressive than a single layer attention mechanism**

“Generating classification weights with GNN denoising auto encoders for few-shot learning ”, Gidaris et al. 19

## Generate weights with a GNN Denoising AutoEncoder



- **Learning inter-class correlations with GNN based Denoising AutoEncoders**
- **DAE:** reconstructs initial (noisy) prototypical averaging weights
  - Meta-training here can be data hungry
  - injecting noise during meta-training → **regularize meta-training**

“Generating classification weights with GNN denoising auto encoders for few-shot learning ”, Gidaris et al. 19

## Few-shot learning without forgetting

Approach	Novel classes					All classes				
	K=1	2	5	10	20	K=1	2	5	10	20
<b>Prior work</b>										
Prototypical Networks	39.3	54.4	66.3	71.2	73.9	49.5	61.0	69.7	72.9	74.6
Matching Networks	43.6	54.0	66.0	72.5	76.9	54.4	61.0	69.0	73.7	76.5
Logistic regression [Hariharan <i>et al.</i> 16]	38.4	51.1	64.8	71.6	76.6	40.8	49.9	64.2	71.9	76.9
Logistic regression w/ H [Hariharan <i>et al.</i> 16]	40.7	50.8	62.0	69.3	76.5	52.2	59.4	67.6	72.8	76.9
Prototype Matching Nets w/ H [Wang <i>et al.</i> 18]	45.8	57.8	69.0	74.3	77.4	57.6	64.7	71.9	75.2	77.5
<b>Cosine Classifier with few-shot classification weight generation</b>										
Feature Averaging [Gidaris <i>et al.</i> 18]	45.4	56.9	68.9	74.5	77.7	57.0	64.3	72.3	75.6	77.3
Attention Mechanism [Gidaris <i>et al.</i> 18]	46.2	57.5	69.2	74.8	<b>78.1</b>	58.2	65.2	72.7	<b>76.5</b>	<b>78.7</b>
GNN Denoising AutoEncoder [Gidaris <i>et al.</i> 19]	<b>48.0</b>	<b>59.7</b>	<b>70.3</b>	<b>75.0</b>	77.8	<b>59.1</b>	<b>66.3</b>	<b>73.2</b>	76.1	77.5

Table 2: Top-5 accuracies on the novel and on all classes for the ImageNet-FS benchmark [13]. To report results we use 100 test episodes.

**Exploiting inter-class correlations (attention, GNN) leads to better performance**

Source: "Generating classification weights with GNN denoising auto encoders for few-shot learning ", Gidaris et al. 18

## Learn to generate classification weights

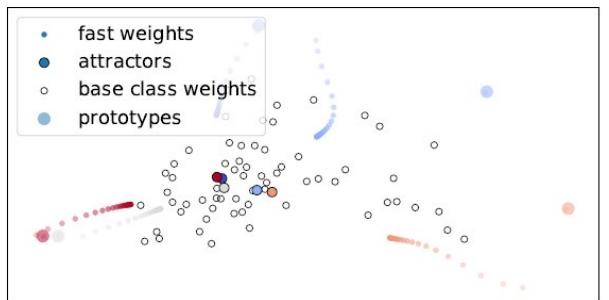
- **(Almost) simple training:**
    - Single classification network, standard supervised pre-training
    - Meta-training: only for the parameter generating module
  - More flexible: unified recognition of both base and novel classes
  - Same test speed as typical classification networks
- 
- The parameter generating module might be data hungry
  - Constrained by quality of pre-trained representations
    - Similar to metric learning based methods

# Learning Weights with Attention Attractor Networks

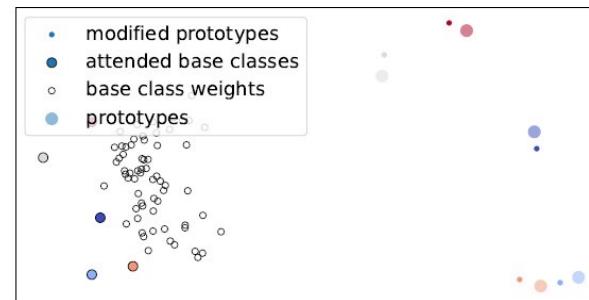
**Optimization-based meta-learning with dynamic regularization:**

$$W = \min_W \left( \text{CrossEntropyLoss}(S, W) + \sum_i R(w_i - w_i^{att}) \right)$$

- The meta-learner is trained to predict (using  $S_i$ ) priors  $w_i^{att}$  so that the optimized weights  $W$  would minimize the classification loss on the query set  $Q$



(a) Ours



(b) LwoF [9]

Figure 3: Visualization of a 5-shot 64+5-way episode using PCA. **Left:** Our attractor model learns to “pull” prototypes (large colored circles) towards base class weights (white circles). We visualize the trajectories during episodic training; **Right:** Dynamic few-shot learning without forgetting [9].

# **Agenda**

- Introduction
- Main types of few-shot learning algorithms
- Few-shot learning without forgetting
- Final notes

## Final notes

- **Few-shot visual learning is important**

## Final notes

- **Few-shot visual learning is important**
- **But, common few-shot benchmarks are insufficient**
  - Omniglot: saturated
  - MinilmageNet: with proper tuning all methods achieve similar results, not realistic

## Final notes

- **Few-shot visual learning is important**
- **But, common few-shot benchmarks are insufficient**
  - Omniglot: saturated
  - MinilmageNet: with enough tuning all methods achieve similar results, not realistic setting
- **More realistic benchmarks:**
  - “Low-shot Visual Recognition by Shrinking and Hallucinating Features”, Hariharan et al. 17
  - “Few-Shot Learning with Localization in Realistic Settings”, Wertheimer et al. 19
  - “Large-Scale Long-Tailed Recognition in an Open World”, Liu et al. 19
  - “Meta-Dataset: A dataset for datasets for learning to learn from few examples”, Triantafillou et al. 19

# A Closer Look to Few-Shot Classification

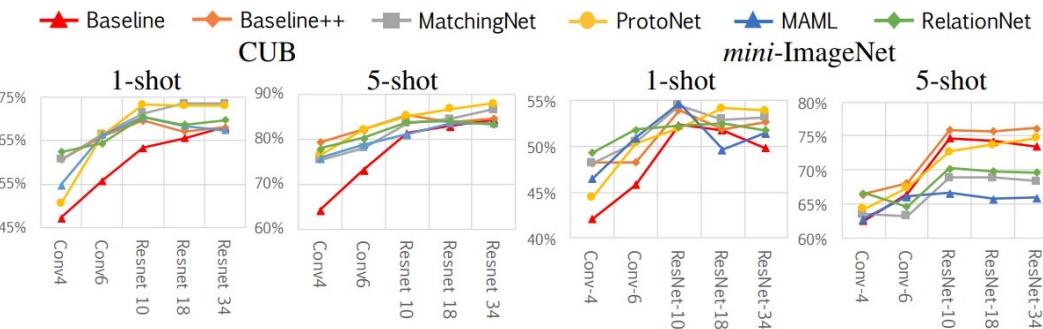


Figure 3: **Few-shot classification accuracy vs. backbone depth.** In the CUB dataset, gaps among different methods diminish as the backbone gets deeper. In mini-ImageNet 5-shot, some meta-learning methods are even beaten by Baseline with a deeper backbone. (Please refer to

"A Closer Look to Few-shot classification", Chen et al. 19

Method	CUB		mini-ImageNet	
	1-shot	5-shot	1-shot	5-shot
Baseline	47.12 ± 0.74	64.16 ± 0.71	42.11 ± 0.71	62.53 ± 0.69
Baseline++	60.53 ± 0.83	79.34 ± 0.61	48.24 ± 0.75	66.43 ± 0.63
MatchingNet Vinyals et al. (2016)	60.52 ± 0.88	75.29 ± 0.75	48.14 ± 0.78	63.48 ± 0.66
ProtoNet Snell et al. (2017)	50.46 ± 0.88	76.39 ± 0.64	44.42 ± 0.84	64.24 ± 0.72
MAML Finn et al. (2017)	54.73 ± 0.97	75.75 ± 0.76	46.47 ± 0.82	62.71 ± 0.71
RelationNet Sung et al. (2018)	62.34 ± 0.94	77.84 ± 0.68	49.31 ± 0.85	66.60 ± 0.69

**Baseline:**  
pre-training + fine-tuning last layer

**Baseline++:**  
cosine classifier

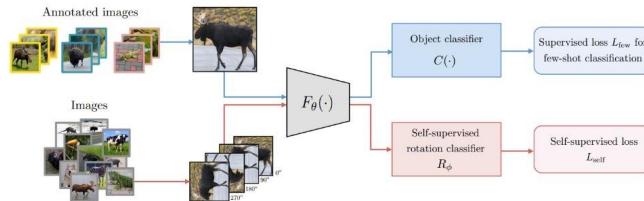
- Meta-learning algorithms and network designs of **growing complexity**, but
- **Well-tuned baselines: often on par / better than SoTA meta-learning methods**
- **Baselines: scale better with deeper backbones**

# A different direction

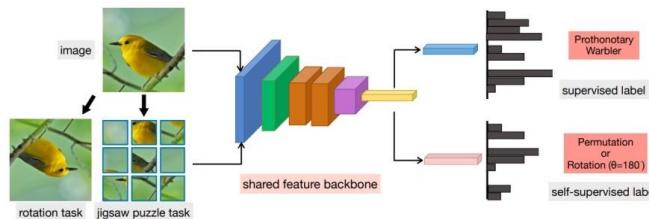
## Focus on pre-training richer representations

- Representations that know more about the world can adapt better
- **Leveraging self-supervision** (see next talk by Relja and Andrey)

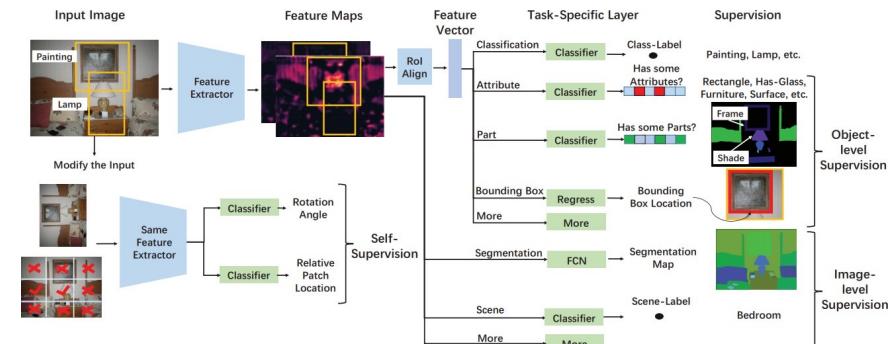
“Boosting few-shot visual learning with self-supervision”,  
Gidaris et al. 19



“When does self-supervision improve few-shot learning?”,  
Su et al. 19



“Learning generalizable representations via diverse supervision”, Pang et al. 19



Also:

“Charting the right manifold: manifold mixup for few-shot learning”, Mangla et al. 20

“Rethinking few-shot image classification: a good embedding is all you need?”, Tian et al. 20

## Not covered because of time constraints

- **Semi-supervised few-shot / meta learning:**
  - “Low-shot learning with large-scale diffusion”, Douze et al. 18
  - “Meta-learning for semi-supervised few-shot classification”, Triantafillou et al. 18
- **Few-shot / meta learning with noise labels:**
  - “Graph convolutional networks for learning with few clean and many noisy labels”, Iscen et al 20
- **Learning with imbalanced datasets (many-shot and few-shot classes):**
  - “Learning to model the tail”, Wang et al. 17
  - “Large-scale long-tailed recognition in an open world”, Liu et al. 19
  - “Decoupling representation and classifier for long-tailed recognition”, Kang et al. 20
- **Few-shot learning beyond image classification:**
  - “Few-shot object detection via feature reweighting”, Kang et al. 19
  - “Meta-learning to detect rare objects”, Wang et al. 19
  - “Few-shot semantic segmentation with prototype learning”, Dong et al. 18
  - “PANet: Few-shot image semantic segmentation with prototype alignment”, Wang et al. 19
  - “Tracking by Instance Detection: A Meta-Learning Approach”, Wang et al. 20
  - ...

### Pytorch learn2learn



### Facebook higher



### Amazon Xfer



### Opne AI GPT-3



- Pytorch backend
- Models
- Datasets
- Algorithms
- Optimizers
- TaskDataset
- Utilities

- Pytorch backend
- Optimizers
- Utilities
- Not support DataParallel

- MxNet backend
- 2019-2020 paper for meta learning/ transfer learning
- Method
- Utilities

- GPT3 Inference for NLP
- Datasets
- API
- Service

## 相關資源

# THANKS!

Q&A