# DEEP LEARNING

## Lecture 13: Meta-Learning

Dr. Yang Lu

Department of Computer Science and Technology

luyang@xmu.edu.cn

- When meta is used as a single word, it refers to itself or to the conventions of its genre, self-referential.

  - The chinese translation is 元.

厦門大學信息學院（特色化示范性软件学院）
School of Informatics Xiamen University (National Characteristic Demonstration Software School)

厦门大学 计算机科学与技术系
Department of Computer Science and Technology, Xiamen University

Image source: https://technews.tw/2021/10/29/a-social-technology-company-meta/
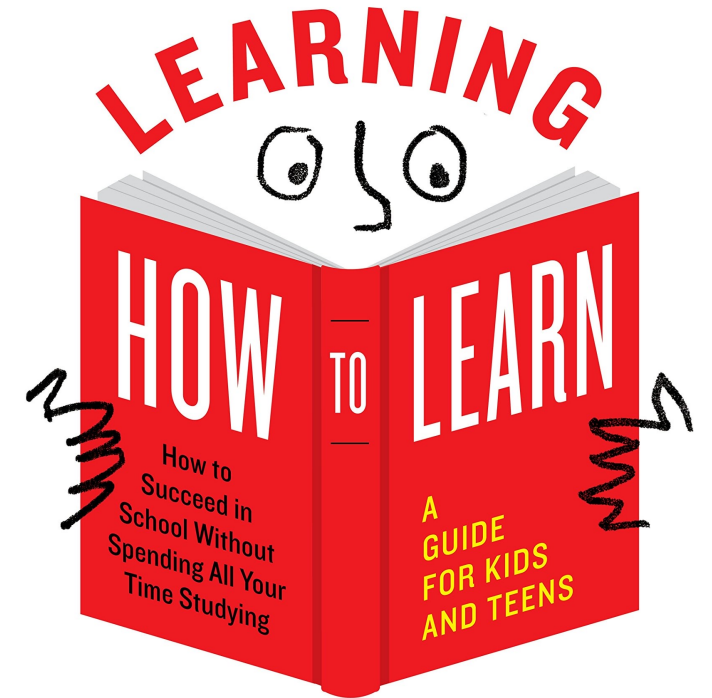
# What is Meta?

- When meta is used as a prefix, meta-X means "beyond-X," "after-X," or "X about X".

- Examples:

  - Metadata: data that describes other data.

  - Metafile: in computer graphics, define objects and images using a list of coordinates.

  - Metaphysics (形而上学): a branch of philosophy that examines the fundamental nature of reality.

  - Meta-analysis: a statistical analysis that combines the results of multiple scientific studies.

  - Metaverse: a virtual world supporting persistent online 3-D virtual environments.

From the bestselling author of *A Mind for Numbers* and the creators of the popular online course Learning How to Learn

- Meta-learning （元学习） also means "beyond learning", "above learning" or "learning about learning".

- It has another name:

  Learning to learn

LEARNING
HOW TO LEARN

How to Succeed in School Without Spending All Your Time Studying

A GUIDE FOR KIDS AND TEENS

BARBARA OAKLEY, PhD, AND
TERRENCE SEJNOWSKI, PhD,
WITH ALISTAIR McCONVILLE

厦門大學信息学院（特色化示范性软件学院）
School of Informatics Xiamen University (National Characteristic Demonstration Software School)

厦门大学 计算机科学与技术系
Department of Computer Science and Technology, Xiamen University

Image source: https://images-na.ssl-images-amazon.com/images/I/81j+JH7WshL.jpg

# Meta-Learning

- Historically, the success of machine learning was driven by the choice of hand-engineered features with model learning.

  - Problem: how to select feature…

- Deep learning realized the promise of joint feature and model learning, providing a huge improvement in performance for many tasks.

  - Problem: how to select algorithm/hyperparameter…

- Meta-learning in neural networks can be seen as aiming to provide the next step of integrating joint feature, model, and algorithm learning.

厦門大学信息学院（特色化示范性软件学院）
School of Informatics Xiamen University (National Characteristic Demonstration Software School)

厦门大学 计算机科学与技术系
Department of Computer Science and Technology, Xiamen University

- Learning aspects:

  - Model learning: select the best model for the task.

  - Feature learning: select the best feature for the task.

  - Algorithm learning: select the best algorithm for the task.

- How do we select an algorithm to train a model?

  - Manually try different algorithms.

  - Manually try an algorithm with different hyperparameters.

# Algorithm Selection vs. Algorithm Learning

Industry

Academia





Using 1000 GPUs to try 1000 sets of hyperparameters

Randomly choose a set by imagination and claim that it is the best!

Meta-learning automates this procedure by end-to-end neural networks.

# Relation to AutoML

- AutoML aims to automate parts of the machin[e] learning process that are typically manual.

  - Such as data preparation, algorithm selection, hype[r] parameter tuning, and architecture search.

- AutoML sometimes makes use of end-to-en[d] optimization.

  - Meta-learning can be seen as a specialization of AutoML.

- <span style="color:red">Meta-learning is about algorithm learning, rather than algorithm tuning.</span>

Google Cloud
AutoML

# From Machine Learning to Meta-Learning

■In machine learning, we select an algorithm $F$, train it by optimizing the parameter $\theta$, and obtain model $f_\theta$.

Training data

$$F( \text{[cat] [dog]} ) \rightarrow f_\theta$$

cat    dog

Test data

$$f_\theta( \text{[cat image]} ) \rightarrow \text{cat}$$

By the idea of learning to learn, can we learn $F$?

Image source: Meta Learning: Learn to Learn, Hung-Yi Lee

■ Now, our goal is to find the best algorithm $F$ for the task, just like the best model $f$ for the data.



$A($ cat dog $) \rightarrow F_\omega$

Training data

$F_\omega($ cat dog $) \rightarrow f_\theta$

Test data

$f_\theta($ $) \rightarrow$ cat

$\omega$ is learnable algorithm parameters. It is usually called meta-knowledge.

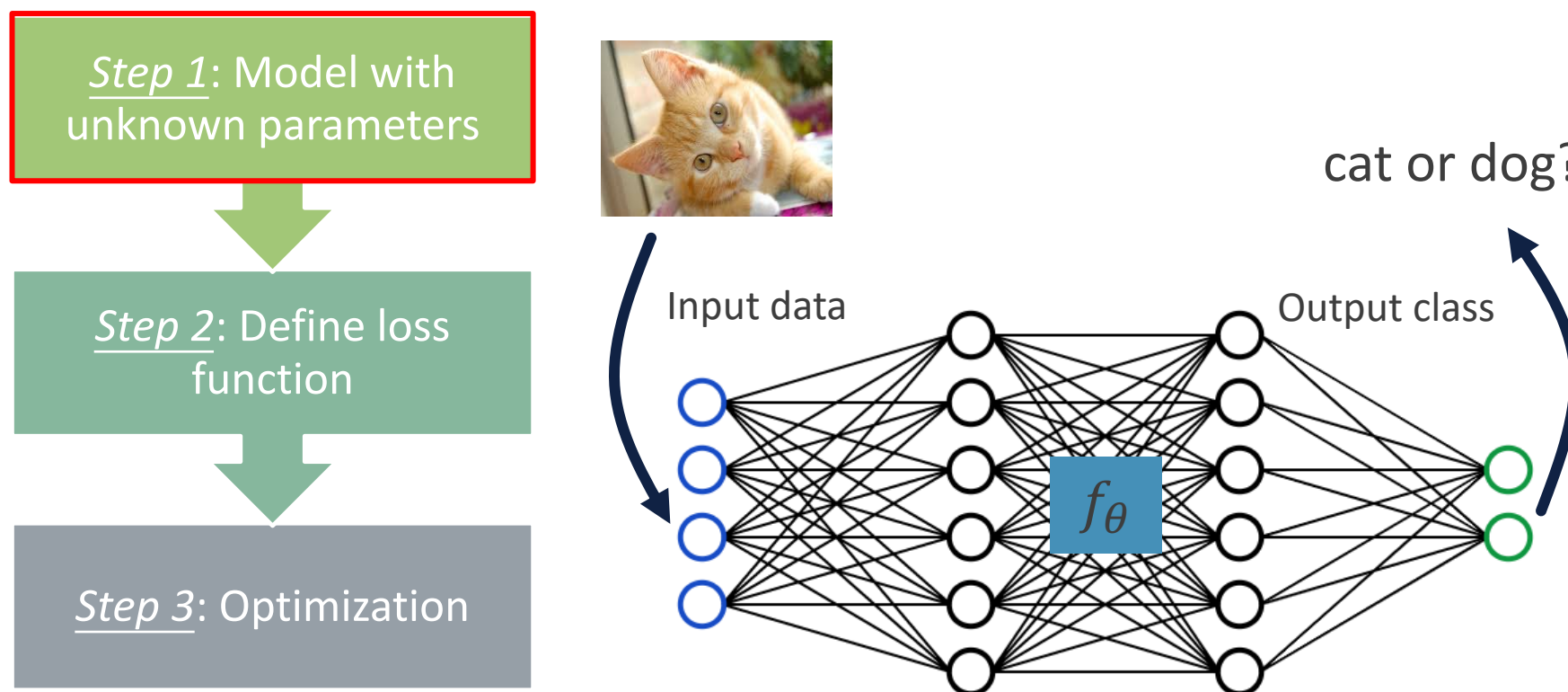- When we deal with multiple tasks, we can also train $F_\omega$ to be good for all tasks. In this way, it is capable of generalizing new task.



Training tasks

Training data of new task

$$A(\ \ ) \to F_\omega$$

$$F_\omega(\ \ ) \to f_\theta$$

Test data of new task

$$f_\theta(\ \ ) \to \text{laptop}$$

e

■ Single-task and multi-task meta-learning actually deals with different problems.

■ Single-task meta-learning aims at learning the most suitable algorithm for this task.

■ Multi-task meta-learning aims at learning the most suitable algorithm for all tasks, and be capable of dealing with new task.
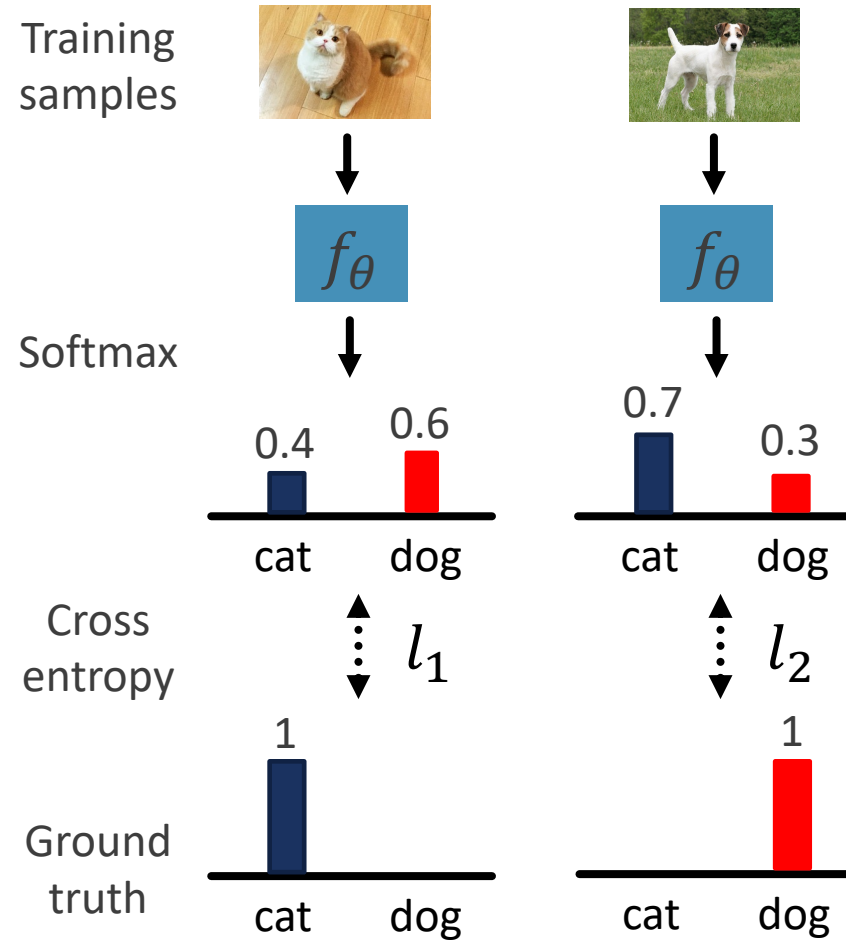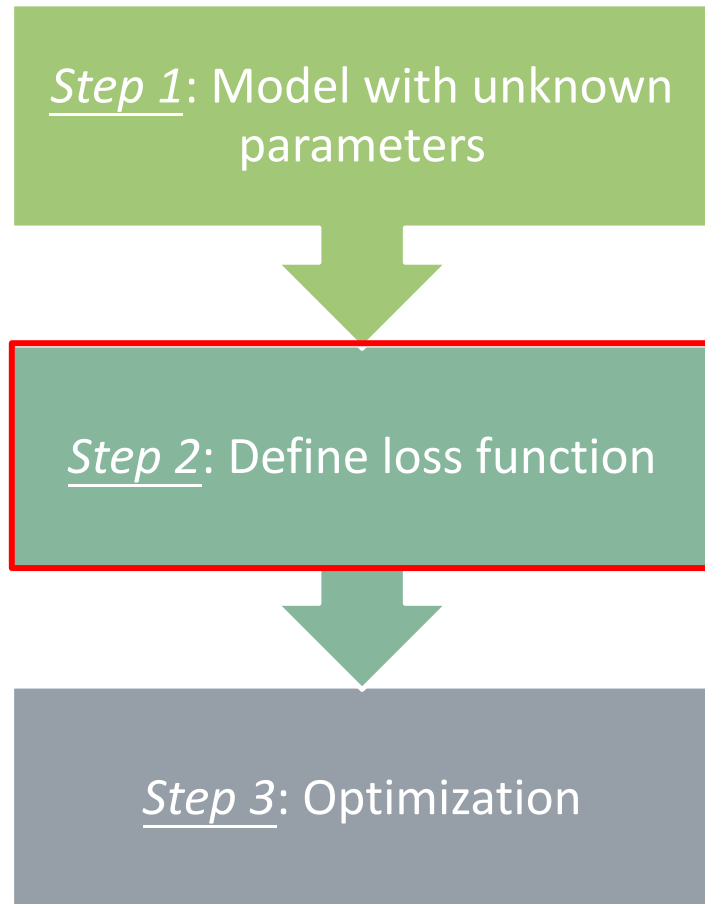
Notice the difference between multi-task meta-learning and multi-task learning.

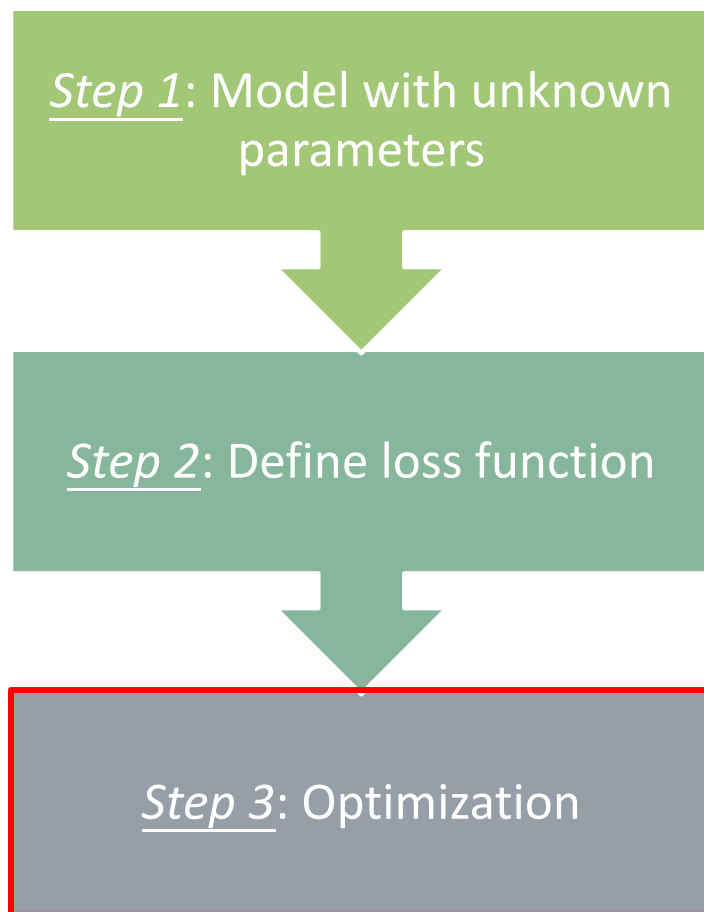School of Informatics Xiamen University (National Characteristic Demonstration Software School)

Department of Computer Science and Technology, Xiamen University

11

- Now, the problem is <span style="color:red">how to learn algorithm parameter $\omega$?</span>
- Recall that how we learn model parameter $\theta$.



*Step 1*: Model with unknown parameters

*Step 2*: Define loss function

*Step 3*: Optimization

Input data

$f_\theta$

Output class

cat or dog?

Content reproduced from "Meta Learning: Learn to Learn" by Hung-Yi Lee

# How to Learn Model Parameter?



*Step 1*: Model with unknown parameters

*Step 2*: Define loss function

*Step 3*: Optimization

Training samples

$f_\theta$   $f_\theta$

Softmax

0.4   0.6        0.7   0.3

cat   dog        cat   dog

Cross entropy   $l_1$        $l_2$

1                    1

Ground truth

cat   dog        cat   dog

# How to Learn Model Parameter?

**Step 1**: Model with unknown parameters

⬇

**Step 2**: Define loss function

⬇

**Step 3**: Optimization

loss:

$$L(\theta) = \sum_{i=1}^{|D|} l_i$$

sum over training examples

$$\theta^* = \underset{\theta}{\arg\min} \, L(\theta)$$

done by gradient descent

$f_{\theta^*}$ is the model learned by a learning algorithm from data.

- Learning algorithm parameter is exactly same as learning model parameter.

- Instead of generalizing over data, it generalizes over tasks.

# How to Learn Algorithm Parameter?

**Task 1**

Training samples

Testing samples

$F_\omega$

Optimize

$f_{\theta_1^*}$

Evaluation

$l^1$

**Task 2**

Training samples

Testing samples

$F_\omega$

Optimize

$f_{\theta_2^*}$

Evaluation

$l^2$

Total loss: $L(\phi) = \sum_{n=1}^{N} l^n$   ($N$ is the number of the training tasks)

Content reproduced from "Meta Learning: Learn to Learn" by Hung-Yi Lee

# How to Learn Algorithm Parameter?

- In typical machine learning, the loss is computed based on training examples.

- In meta-learning, the loss is computed based on testing examples.

- Is there any problem here?

You dare use testing
data during training!

# How to Learn Algorithm Parameter?

- However, we have to use "testing data" to evaluate how the algorithm parameter performs.

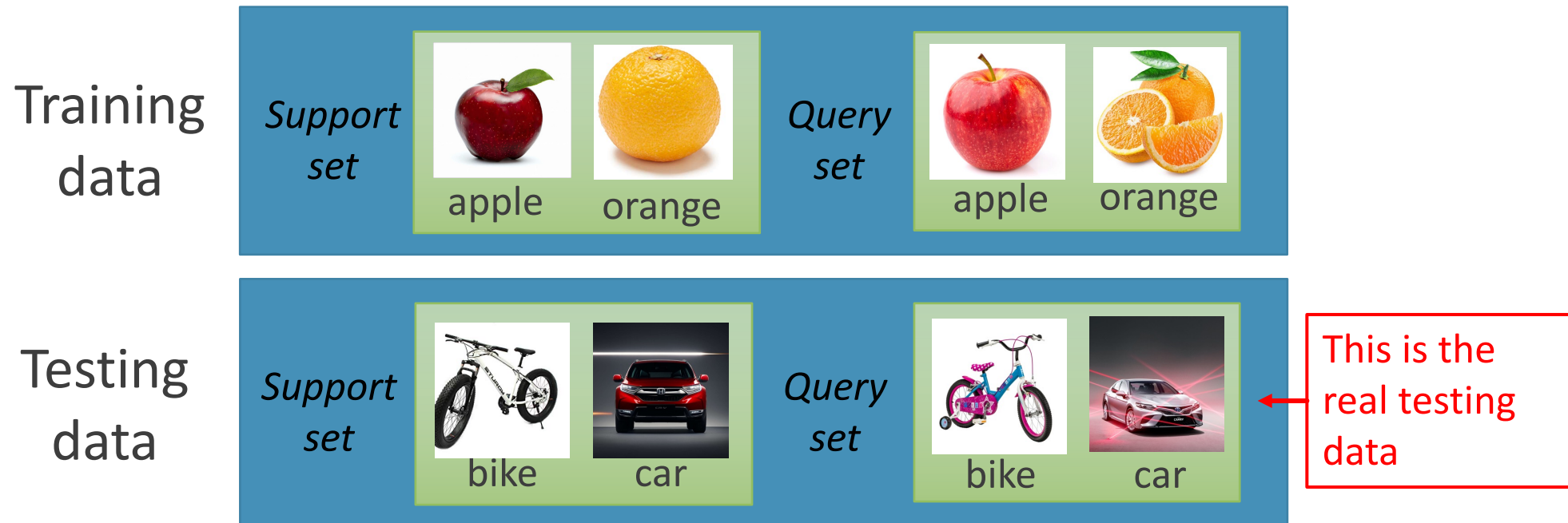- We can cut a part from the training data, just like validation data.

Training data

*Support set*

apple    orange

*Query set*

apple    orange

Used for optimizing model parameter $\theta$

Used for optimizing algorithm parameter $\omega$

Query set is nothing but "validation data for training"

# How to Learn Algorithm Parameter?

- Similarly, we also split testing data into support and query set.

- The label of testing support set is available during testing.

- The testing query set is the real testing data to evaluate the performance of meta-learning algorithm.

- Training and testing phases in meta-learning is called meta-training and meta-testing.

  - Meta-training can be formalized as bi-level optimization:

  Loss to evaluate algorithm

  Outer loop

  $$\omega^* = \underset{\omega}{\arg\min} \sum_{i=1}^{M} L^{meta}\left(\theta^{*(i)}(\omega), D_{train}^{query(i)}\right)$$

  Inner loop

  $$s.t.\ \theta^{*(i)}(\omega) = \underset{\theta}{\arg\min} L^{task}\left(\theta, \omega, D_{train}^{support(i)}\right)$$

  Loss to evaluate model

  - Meta-testing for new task $j$ can be formalized as:

  $$\theta^* = \underset{\theta}{\arg\min} L^{task}\left(\theta, \omega^*, D_{test}^{support(j)}\right)$$

  And finally we use $f_{\theta^*}$ to predict samples in $D_{test}^{query(j)}$.

# Few-Shot Learning

- One direct application of multi-task meta-learning is few-shot learning (小样本学习).

- Do we human beings need great amount of training data to recognize image category?

- No! We have the learning ability. We know how to learn!

  - When we deal with new task, our experience help us learn with only a few samples.

  - We are experts of "learning to learn".

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

Image source: https://cdn.britannica.com/80/146380-050-9144819E/Okapi.jpg

# Few-Shot Learning



A typical setting for few-shot learning: $n$-way-$k$-shot
$k$ is usually set at 1 or 5

- Few-shot learning has its own benchmark datasets.

  - MiniImageNet, Fewshot-CIFAR100, Omniglot, etc.

- A typical dataset split for MiniImageNet is: 64 training classes, 12 validation classes, and 24 test classes.

  - No class overlap among training, validating and test.

- We run a lot of episodes for training. In each episode, we randomly select $N$ classes with $k + 1$ samples.

  - $k$ support samples and 1 query sample.

# Few-Shot Classification Leaderboard

## *mini*ImageNet Leaderboard (5-class)
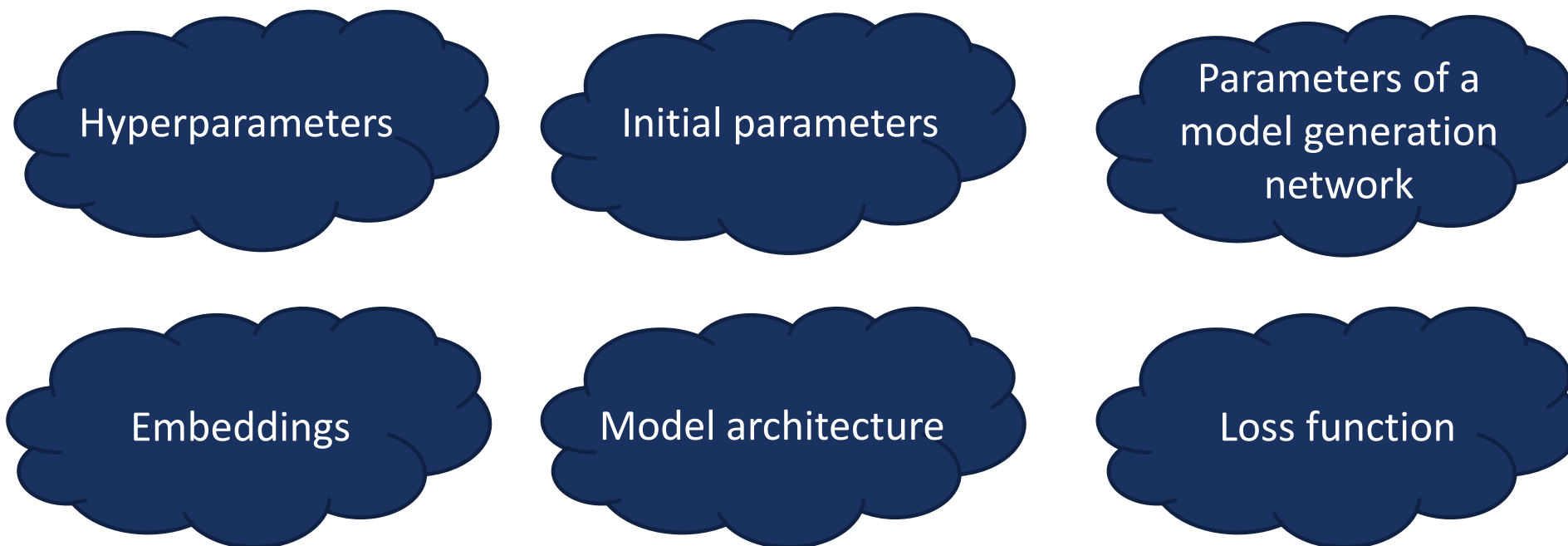
Edit this leaderboard

Search: [_____]

| Method | Venue | Year | Backbone | Setting | 1-shot | 5-shot | Code | Reported by |
|--------|-------|------|----------|---------|--------|--------|------|-------------|
| EASY | arXiv | 2022 | 3xResNet012 | Transductive | 84.04 ± 0.23 | 89.14 ± 0.11 | [PyTorch] | [Source] |
| iLPC | ICCV | 2021 | WRN-28-10 | Semi-supervised | 83.58±0.79 | 89.68±0.37 | [PyTorch] | [Source] |
| iLPC | ICCV | 2021 | WRN-28-10 | Transductive | 83.05±0.79 | 88.82±0.42 | [PyTorch] | [Source] |
| PT+MAP | arXiv | 2021 | WRN | Transductive | 82.92 ± 0.26 | 88.82 ± 0.13 | [PyTorch] | [Source] |
| PTN | AAAI | 2021 | WRN-28-10 | Semi-supervised | 82.66 ± 0.97 | 88.43 ± 0.67 | None | [Source] |
| EASY | arXiv | 2022 | 2xResNet-12(1/√2) | Transductive | 82.31 ± 0.24 | 88.57 ± 0.12 | [PyTorch] | [Source] |
| Simple CNAPS | CVPR | 2020 | ResNet18 (pre-trained on ImageNet) | Inductive | 82.16 | 89.80 | [PyTorch] | [Source] |
| Oblique Manifold | ICCV | 2021 | WRN-28-10 | Transductive | 80.64±0.34 | 89.39±0.39 | [PyTorch] | [Source] |
| ICA + MSP | ECCV | 2020 | DenseNet | Semi-supervised | 80.11 ± 0.25 | 85.78 ± 0.13 | None | [Source] |
| EPNet | ECCV | 2020 | WRN-28-10 | Semi-supervised | 79.22 ± 0.92 | 88.05 ± 0.51 | [PyTorch] | [Source] |

Source: https://few-shot.yyliu.net/miniimagenet.html

# Meta-Knowledge

- Talking so much about the learnable algorithm $F_\omega$.

- What exactly are those algorithm parameters $\omega$?

Hyperparameters

Initial parameters

Parameters of a model generation network

Embeddings

Model architecture

Loss function

• • •

# Meta-Knowledge Taxonomy

- Optimization-based Method

- Model-based Method

- Metric-based Method

# OPTIMIZATION-BASED METHOD

# Optimization-Based Method

■ The meta-knowledge $\omega$ is related to the optimization process.

- ■ Learning to optimize
- ■ Learning to initialize
- ■ Learning to weight
- ■ Learning to reward
- ■ Learning to augment
- ■ Dataset distillation
- ■ Neural architecture search

# Learning to learn by gradient descent by gradient descent

Marcin Andrychowicz[1], Misha Denil[1], Sergio Gómez Colmenarejo[1], Matthew W. Hoffman[1],
David Pfau[1], Tom Schaul[1], Brendan Shillingford[1,2], Nando de Freitas[1,2,3]

[1]Google DeepMind    [2]University of Oxford    [3]Canadian Institute for Advanced Research

marcin.andrychowicz@gmail.com
{mdenil,sergomez,mwhoffman,pfau,schaul}@google.com
brendan.shillingford@cs.ox.ac.uk, nandodefreitas@google.com

厦門大學信息学院（特色化示范性软件学院）
School of Informatics Xiamen University (National Characteristic Demonstration Software School)

厦门大学 计算机科学与技术系
Department of Computer Science and Technology, Xiamen University

M. Andrychowicz, M. Denil, S. G. Colmenarejo, and M. W. Hoffman, "Learning to learn by gradient descent by gradient descent," in NIPS, 2016, pp. 1–17.

- Conventionally, when we do optimization

$$\theta_{t+1} \leftarrow \theta_t + \lambda g\big(\nabla_\theta L(\theta_t)\big)$$
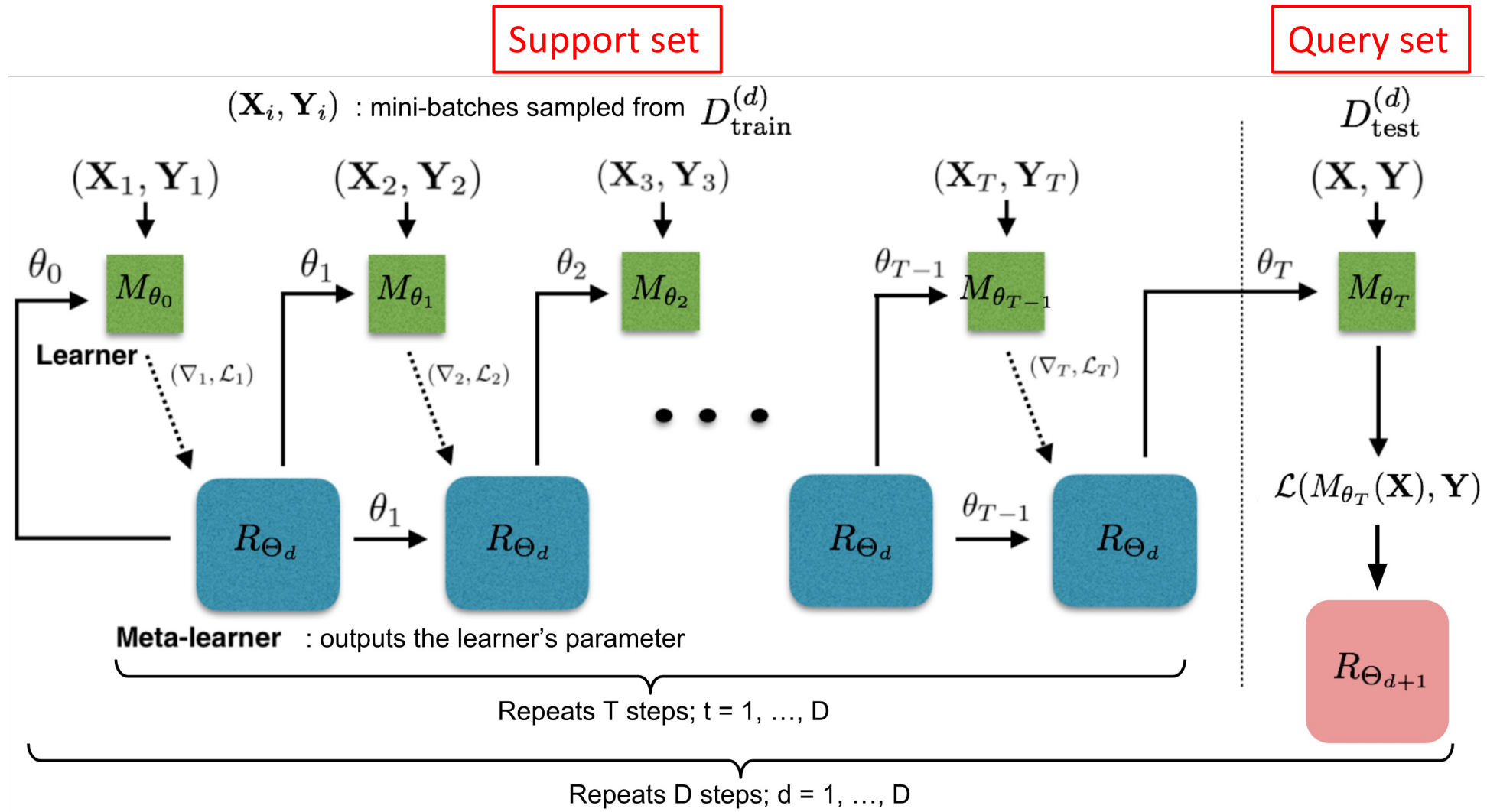
we select the optimizer $g$ such as SGD, momentum, AdaGrad, ADAM, etc.

  - $g$ can be seen as a hand-crafted function of the gradients.

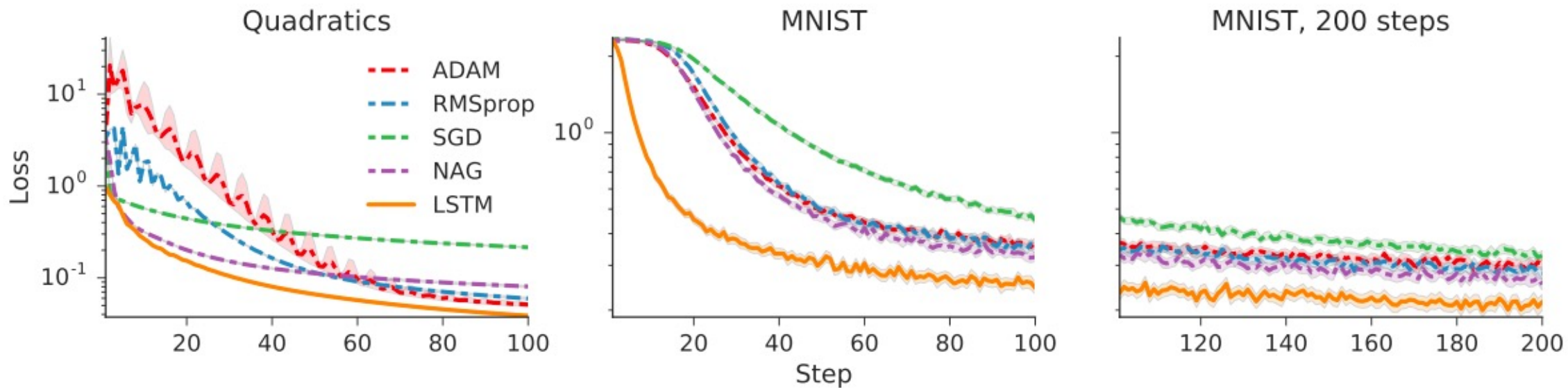- In meta-learning, we can learn proper optimization function by meta-knowledge $\omega$:

$$\theta_{t+1} \leftarrow \theta_t + g_\omega\big(\nabla_\theta L(\theta_t)\big)$$
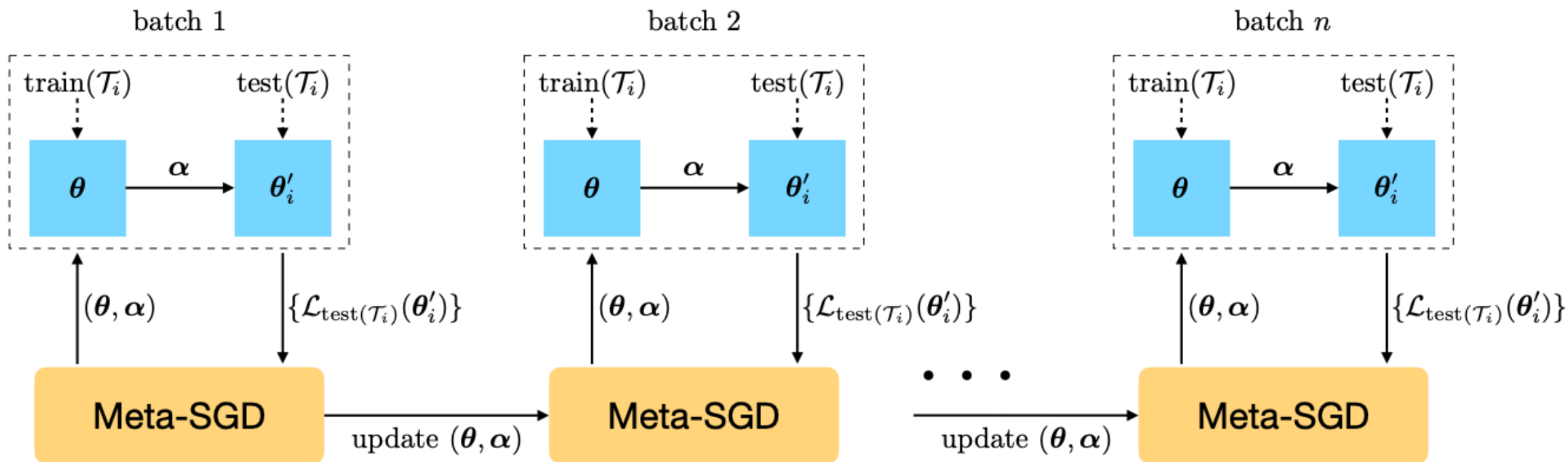
M. Andrychowicz, M. Denil, S. G. Colmenarejo, and M. W. Hoffman, "Learning to learn by gradient descent by gradient descent," in NIPS, 2016, pp. 1–17.

厦門大學信息學院（特色化示范性软件学院）
School of Informatics Xiamen University (National Characteristic Demonstration Software School)

厦门大学 计算机科学与技术系
Department of Computer Science and Technology, Xiamen University

Image source: M. Andrychowicz, M. Denil, S. G. Colmenarejo, and M. W. Hoffman, "Learning to learn by gradient descent by gradient descent," in NIPS, 2016, pp. 1–17.

Learn update direction and learning rate separately

Image source: Z. Li, F. Zhou, F. Chen, and H. Li, "Meta-SGD : Learning to Learn Quickly for Few-Shot Learning," in arXiv, 2017.

# Learning to Initialize

- In Model-Agnostic Meta-Learning (MAML), $\omega$ is the initialized model parameter $\theta$.

- The goal is to find a good $\theta$, such that only a few steps optimization can obtain good model for a task.

meta-learning
learning/adaptation

$\theta$

$\nabla \mathcal{L}_3$

$\nabla \mathcal{L}_2$

$\nabla \mathcal{L}_1$

$\theta_3^*$

$\theta_1^*$

$\theta_2^*$

厦門大學信息学院（特色化示范性软件学院）
School of Informatics Xiamen University (National Characteristic Demonstration Software School)

厦门大学 计算机科学与技术系
Department of Computer Science and Technology, Xiamen University

Image source: C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in ICML, 2017, vol. 3, pp. 1856–1868.

# MAML

- The optimization of MAML follows:

$$\theta_i' \leftarrow \theta - \beta \nabla_\theta L\left(\theta, D_{train}^{support(i)}\right)$$

$$\theta \leftarrow \theta - \alpha \nabla_\theta \sum_{i=1}^{M} L\left(\theta_i', D_{train}^{query(i)}\right)$$

- $\omega$ is the model parameter $\theta$ itself.

- The loss function is same: $L = L^{meta} = L^{task}$.

- It is call model-agnostic because there is no specified meta-learning model for $\omega$.

- Any model can apply MAML.

**Algorithm 1** Model-Agnostic Meta-Learning

**Require:** $p(\mathcal{T})$: distribution over tasks
**Require:** $\alpha, \beta$: step size hyperparameters
1: randomly initialize $\theta$
2: **while** not done **do**       *Outer loop*
3:    Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
4:    **for all** $\mathcal{T}_i$ **do**       *Inner loop*
5:       Evaluate $\nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$ with respect to $K$ examples
6:       Compute adapted parameters with gradient descent: $\theta'_i = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$
7:    **end for**
8:    Update $\theta \leftarrow \theta - \beta \nabla_\theta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$
9: **end while**

Inner loop updates model parameter for each task $i$.

Outer loop updates $\theta$ by evaluating each $\theta_i'$ on query set.

$\theta_i'$ is obtained from $\theta$. Therefore evaluating $\theta_i'$ implicitly evaluates $\theta$.

MAML vs. pre-trained model

- MAML looks for a good initialization to generalize new task.

- Pre-trained model transfers knowledge from a well-learned model on source tasks to a target task by finetuning.

<span style="color:red">What is the difference here?</span>

- MAML doesn't require the initialized model $\theta$ perform well on each task, but the one-step optimized $\theta_i{'}$:

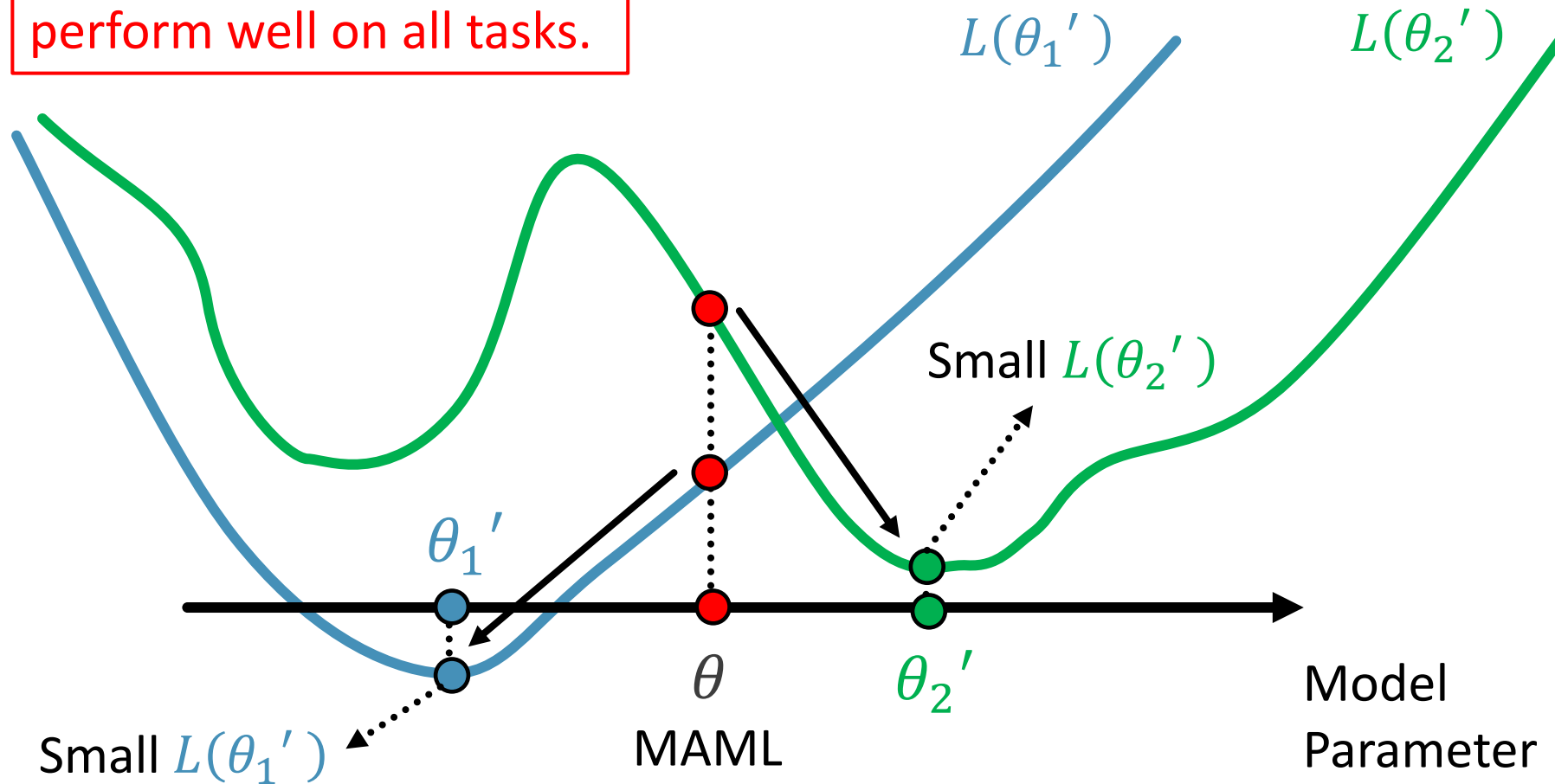$$L(\theta) = \sum_{i=1}^{M} L\left(\theta_i', D_{train}^{query(i)}\right)$$

- Pre-trained model usually require the initialized model $\theta$ perform well on each task:

$$L(\theta) = \sum_{i=1}^{M} L\left(\theta, D_{train}^{(i)}\right)$$

MAML doesn't require $\theta$ perform well on all tasks.

$L(\theta_1')$    $L(\theta_2')$

Small $L(\theta_2')$

$\theta_1'$

$\theta$
MAML

$\theta_2'$

Small $L(\theta_1')$

Model Parameter

Pre-trained model requires $\theta$ perform well on all tasks, but $\theta_i{}'$ may not perform well after finetuning.

$L(\theta_1{}')$   $L(\theta_2{}')$

Large $L(\theta_2{}')$

Small $L(\theta_1{}')$

$\theta_2{}'$   $\theta_1{}'$

$\theta$

Pre-trained Model

Model Parameter

- Toy example: try to learn sine function $y = a\sin(x + b)$.

- Each combination of $a$ and $b$ is a task.

- The goal is to fit a new sin function based on only a few points.



3 random tasks



1000 random tasks

- MAML is able to quickly adapt with only a few datapoints.
- MAML trained model $f_\theta$ has learned to model the periodic nature of the sine wave!

Pre-trained model learns a straight line because it has to minimize the error for all tasks.



Pre-trained model

MAML

Image source: https://towardsdatascience.com/paper-repro-deep-metalearning-using-maml-and-reptile-fd1df1cc81b0

# First-Order MAML

- MAML takes partial derivatives on $\theta$ at both inner and outer loop.
- Therefore, the outer loop actually calculates the second-order derivatives, i.e. Hessian.

$$\nabla_\theta \sum_{i=1}^{M} L(\theta_i')$$

$$= \sum_{i=1}^{M} \nabla_\theta L(\theta_i')$$

$$= \sum_{i=1}^{M} \begin{bmatrix} \partial L(\theta_i')/\partial\theta_1 \\ \partial L(\theta_i')/\partial\theta_2 \\ \vdots \\ \boxed{\partial L(\theta_i')/\partial\theta_j} \\ \vdots \end{bmatrix}$$

$$\theta_k' = \theta - \beta\nabla_\theta L(\theta)$$

$$\partial L(\theta_i')/\partial\theta_j$$

$$= \sum_k \frac{\partial L(\theta_i')}{\partial\theta_k'} \boxed{\frac{\partial\theta_k'}{\partial\theta_j}}$$

$$\frac{\partial\theta_k'}{\partial\theta_j} = \begin{cases} 1 - \beta\dfrac{\partial L(\theta)}{\partial\theta_k\partial\theta_j}, & k = j \\ -\beta\dfrac{\partial L(\theta)}{\partial\theta_k\partial\theta_j} & k \neq j \end{cases}$$

- The second order derivative $\dfrac{\partial L(\theta)}{\partial\theta_k \partial\theta_j}$ is the element of Hessian matrix $H_\theta(L)$.

- We can rewrite the outer loop gradient as:

$$\nabla_\theta L(\theta_i{}') = \big(I - \beta H_\theta(L)\big)\nabla_{\theta_i'} L(\theta_i')$$

- First-Order MAML (FOMAML) calculates the approximation by simply setting the Hessian matrix at 0:

$$\nabla_\theta L(\theta_i{}') \approx \nabla_{\theta_i'} L(\theta_i')$$

# First-Order MAML

$$8: \quad \text{Update } \theta \leftarrow \theta - \beta \nabla_\theta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta_i'})$$

$$9: \textbf{ end while}$$

Replacing $\theta$ by $\theta_i'$ highly improves the efficiency without loss of much accuracy.

| MiniImagenet (Ravi & Larochelle, 2017) | 5-way Accuracy | |
|---|---|---|
| | 1-shot | 5-shot |
| fine-tuning baseline | $28.86 \pm 0.54\%$ | $49.79 \pm 0.79\%$ |
| nearest neighbor baseline | $41.08 \pm 0.70\%$ | $51.04 \pm 0.65\%$ |
| matching nets (Vinyals et al., 2016) | $43.56 \pm 0.84\%$ | $55.31 \pm 0.73\%$ |
| meta-learner LSTM (Ravi & Larochelle, 2017) | $43.44 \pm 0.77\%$ | $60.60 \pm 0.71\%$ |
| **MAML, first order approx. (ours)** | $\mathbf{48.07 \pm 1.75\%}$ | $\mathbf{63.15 \pm 0.91\%}$ |
| **MAML (ours)** | $48.70 \pm 1.84\%$ | $63.11 \pm 0.92\%$ |

■Reptile further simplify FOMAML.

$k$ time update instead of 1

**Algorithm 1** Reptile (serial version)

Initialize $\phi$, the vector of initial parameters
**for** iteration $= 1, 2, \dots$ **do**
  Sample task $\tau$, corresponding to loss $L_\tau$ on weight vectors $\widetilde{\phi}$
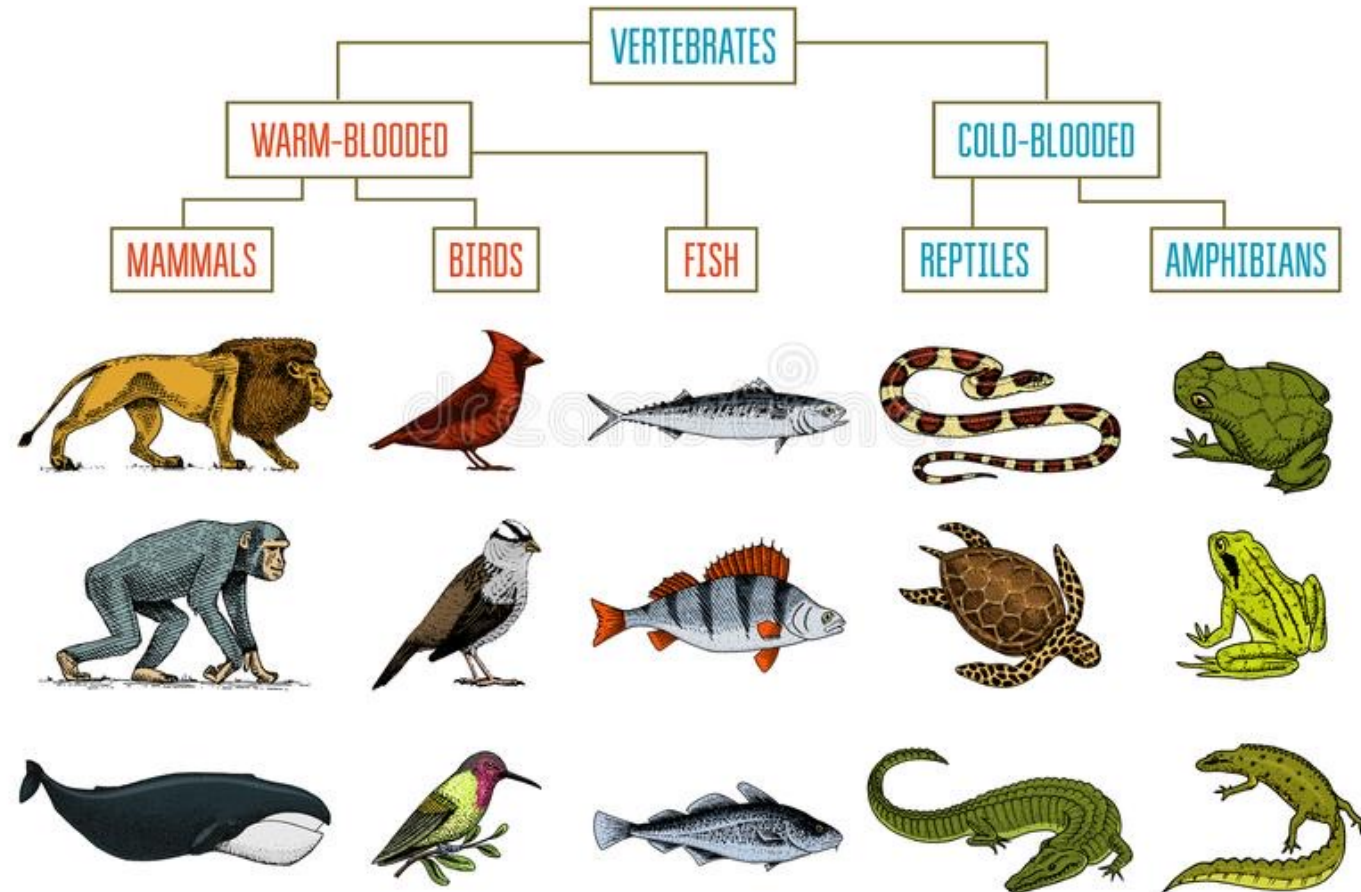  Compute $\widetilde{\phi} = U_\tau^k(\phi)$, denoting $k$ steps of SGD or Adam
  Update $\phi \leftarrow \phi + \epsilon(\widetilde{\phi} - \phi)$
**end for**

Simply use the direction instead of calculating gradient

46

# MAML and Reptile

# Learning to Weight

- During optimization, we may assign different weights to different training samples, according to its learning difficulty.

$$\theta \leftarrow \theta - \alpha \nabla_\theta \sum_{i=1}^{N} w(\boldsymbol{x}_i) l(\theta, \boldsymbol{x}_i)$$

- In this manner, difficult (frequently misclassified) samples are assigned higher weights.

- For example, Focal loss assigns weight by:

Hand-crafted design!

$$w(\boldsymbol{x}_i) = \left(1 - p_{y_i}\right)^{\gamma}$$

where $p_{y_i}$ is the probability belonging to its ground truth $y_i$.

- Can we learn a mapping function from the sample $\boldsymbol{x}_i$ to its weight $w(\boldsymbol{x}_i)$?

- Of course! Simply train an MLP to learn the relationship:

$$\theta \leftarrow \theta - \alpha \nabla_\theta \sum_{i=1}^{N} F_\omega\big(l(\theta, \boldsymbol{x}_i)\big) l(\theta, \boldsymbol{x}_i)$$

where $F_\omega\big(l(\theta, \boldsymbol{x}_i)\big)$ takes the training loss as input and output the corresponding weight.
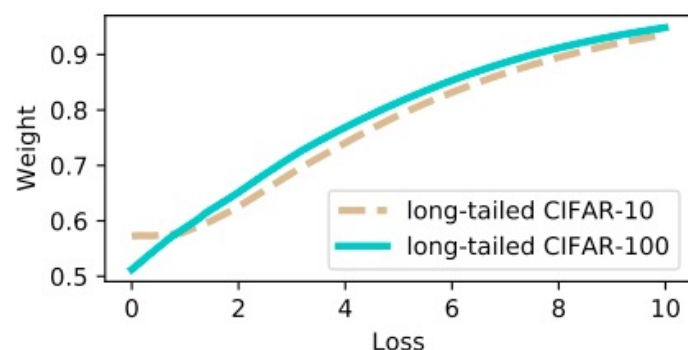
# Learning to Weight

$(\boldsymbol{w}^{(t)}$ in the figure)     $(\Theta^{(t)}$ in the figure)

- Again, we use $D_{train}^{support}$ to optimize $\theta$, and use $D_{train}^{query}$ to optimize $\omega$.

Update meta-knowledge



Temporarily update model parameter
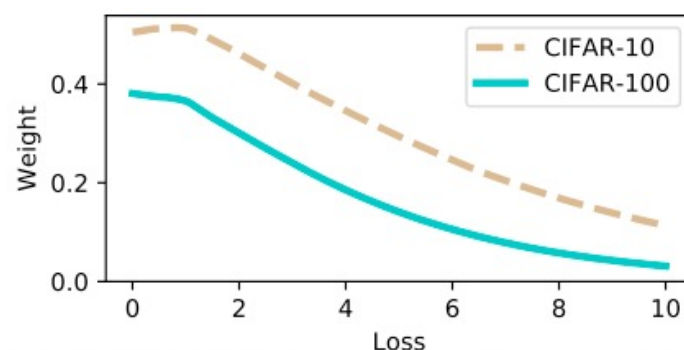
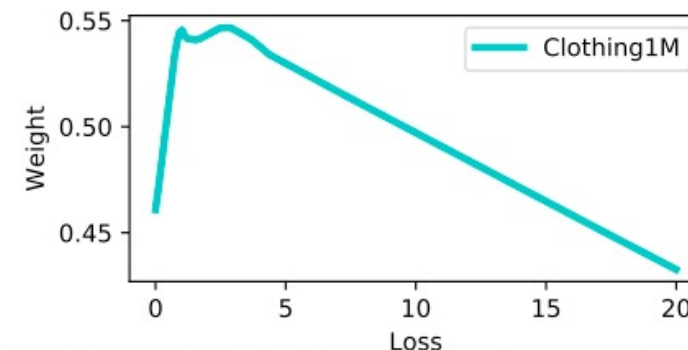Formally update model parameter using updated meta-knowledge

- The weighting function is learned based on the distribution of the dataset.



(d) MW-Net function learned in class imbalance case

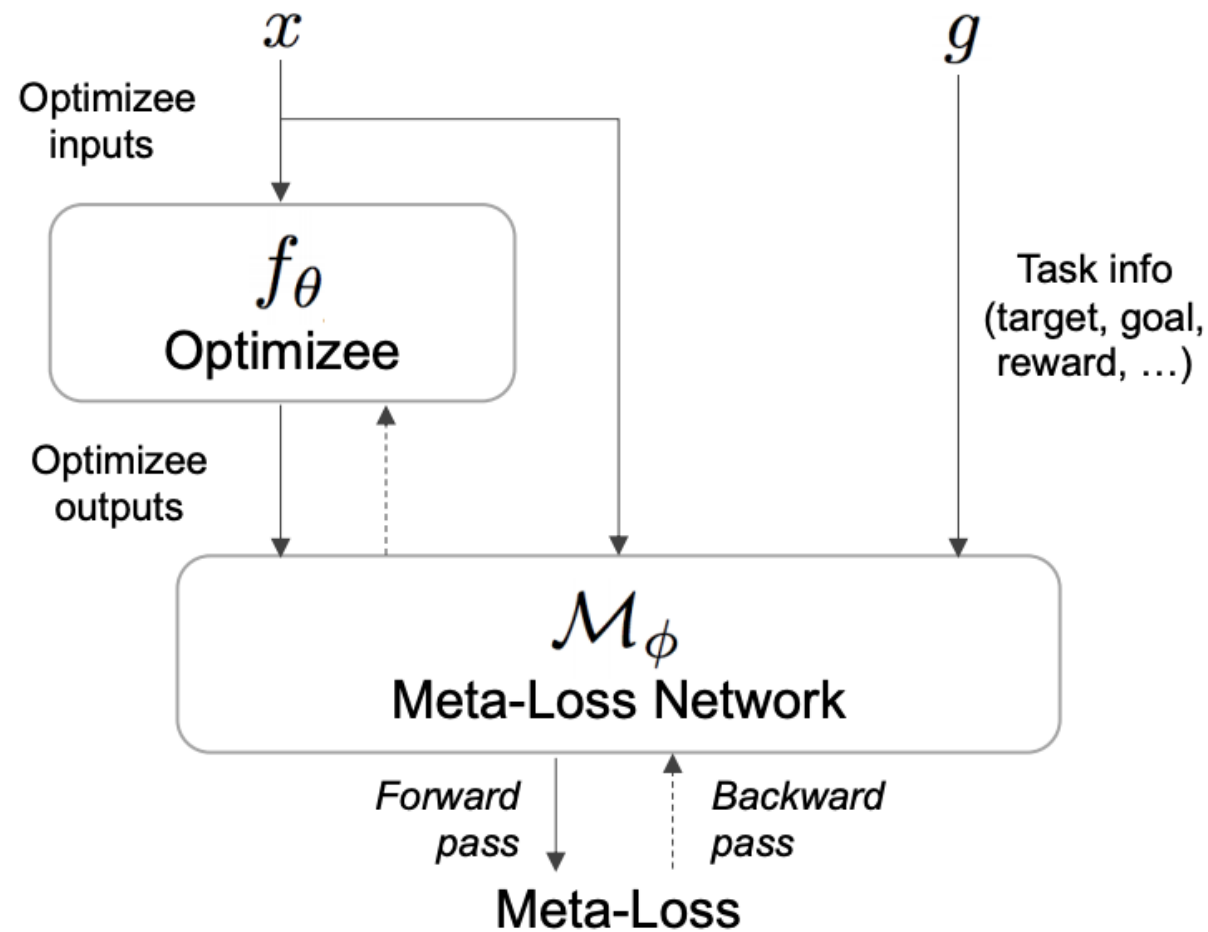(e) MW-Net function learned in corrupter labels case

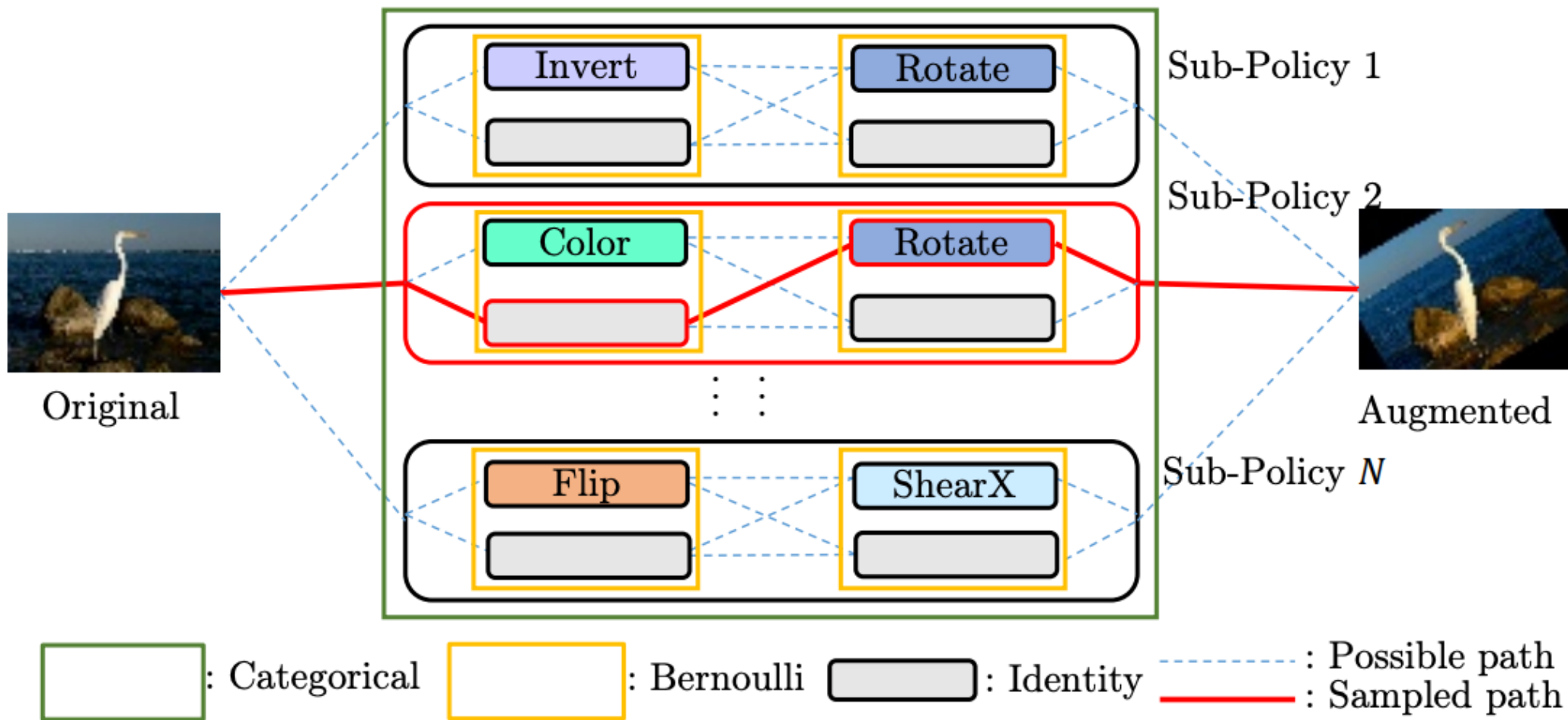(f) MW-Net function learned in real Clothing1M dataset

Image source: J. Shu et al., "Meta-Weight-Net: Learning an Explicit Mapping For Sample Weighting," in NIPS, 2019, pp. 1–12.

Image source: S. Bechtle, "Meta Learning via Learned Loss," in arXiv.

- The data augmentation operation is wrapped up in inner optimization, which is conventionally hand-designed.

  - E.g. crop, zoom, flip, rotate, etc.

- When $\omega$ defines the data augmentation strategy, it can be learned by the outer optimization, in order to maximize validation performance.
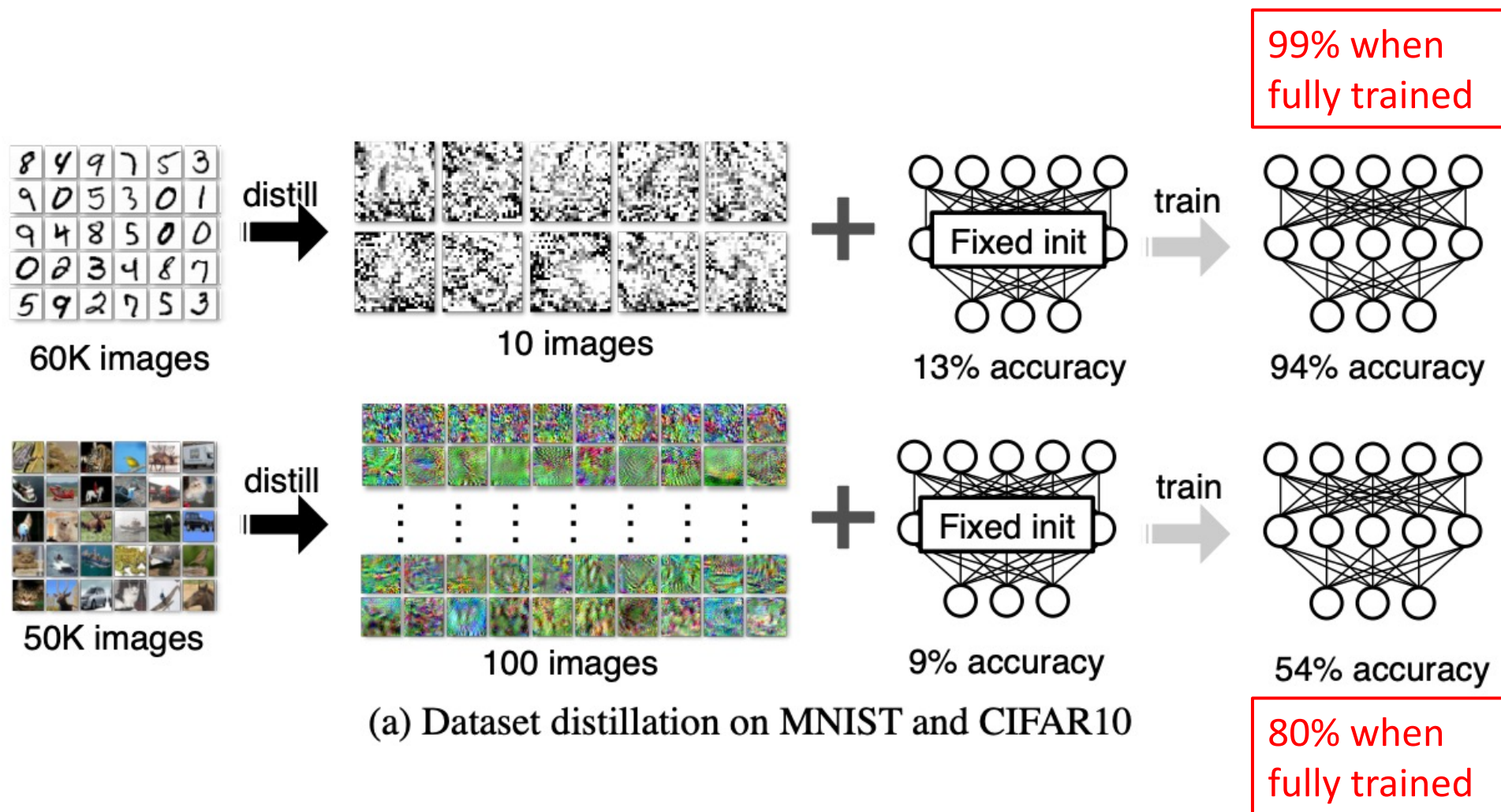
Image source: Y. Li, G. Hu, Y. Wang, T. Hospedales, N. M. Robertson, and Y. Yang, "DADA: Differentiable Automatic Data Augmentation," in arXiv, 2020, pp. 1–16.

厦門大學信息學院（特色化示范性软件学院）
School of Informatics Xiamen University (National Characteristic Demonstration Software School)

厦门大学 计算机科学与技术系
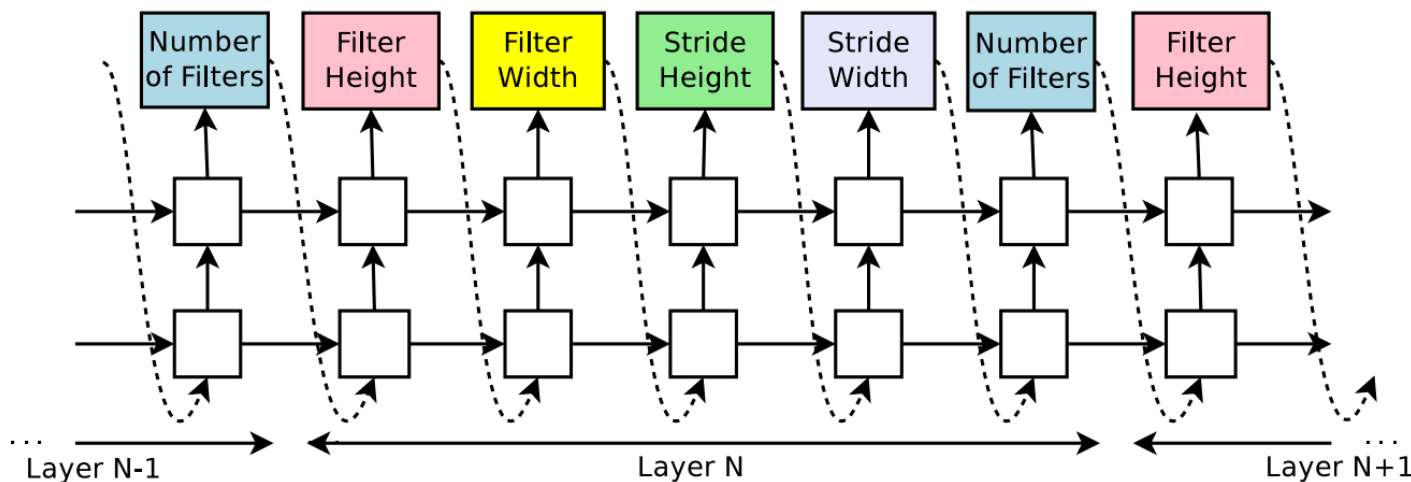Department of Computer Science and Technology, Xiamen University

- In bi-level optimization, we always use the same support data to optimize model parameter $\theta$.

- Can the support data itself be the meta-knowledge $\omega$?
  - Select the most significant samples to train the model.
  - Only a few selected samples can achieve high performance.

# Dataset Distillation



99% when fully trained

80% when fully trained

(a) Dataset distillation on MNIST and CIFAR10

60K images → distill → 10 images + Fixed init → 13% accuracy → train → 94% accuracy

50K images → distill → 100 images + Fixed init → 9% accuracy → train → 54% accuracy

# Neural Architecture Search

- $\omega$ specifies the architecture of a neural network.
  - E.g. number of filters, filter size, stride and pooling size, activation functions, shortcut connections, etc.
- The search space is usually hard to define, and optimize.
  - Most search spaces are broad, and the space of architectures is not trivially differentiable.



Reinforcement learning is usually adopted to search the spaces.

Image source: Zoph, Barret, and Quoc V. Le. "Neural architecture search with reinforcement learning." arXiv preprint arXiv:1611.01578 (2016).

# MODEL-BASED METHOD

- The optimization of model network $f_\theta$ in all optimization-based methods are still based on gradient descend.

    - In the inner loop, given $\omega$, we optimize $\theta$.

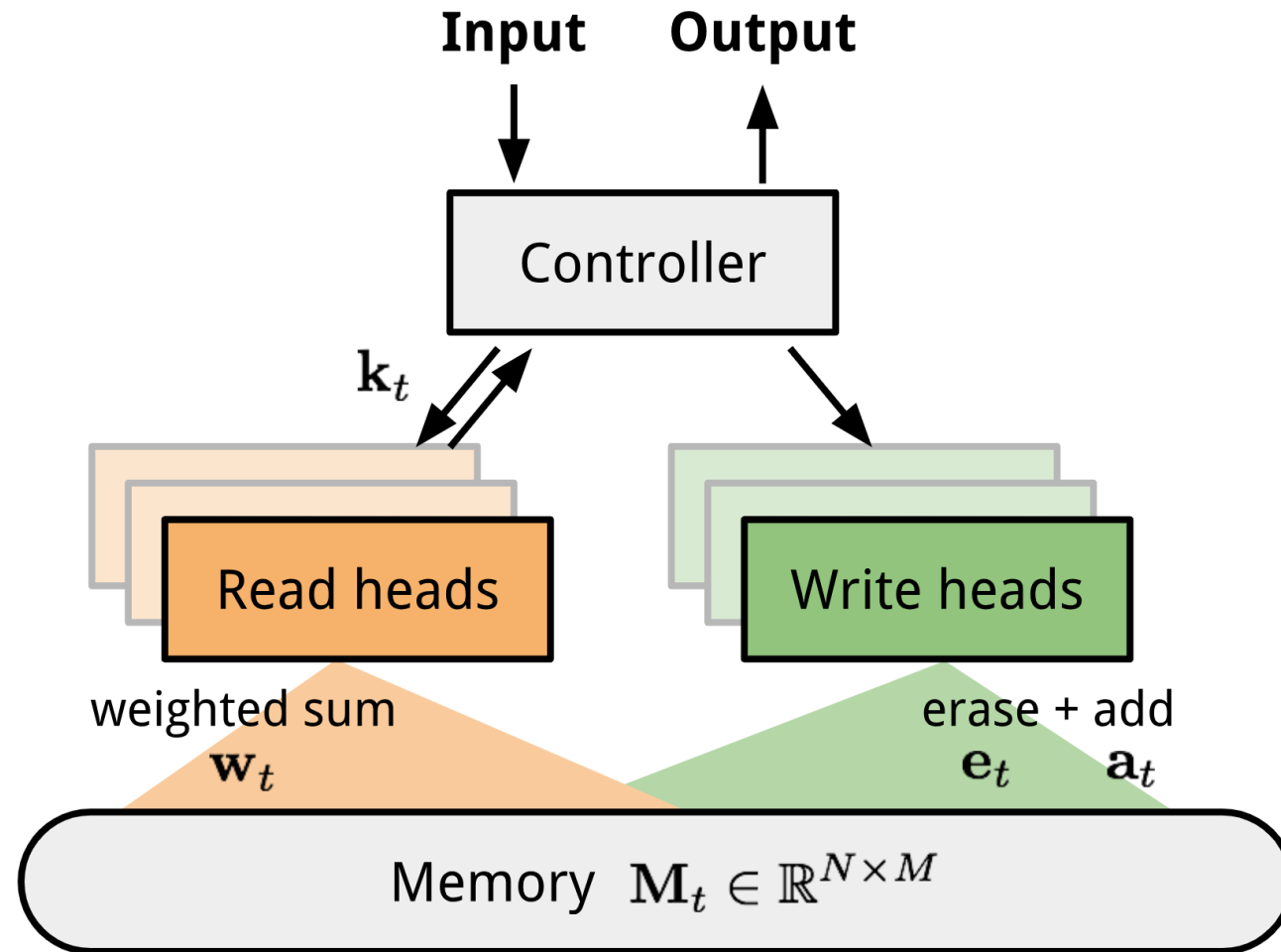$$\theta^*(\omega) = \underset{\theta}{\mathrm{argmin}}\, L^{task}\left(\theta, \omega, D_{train}^{support}\right)$$

- Can we omit this optimization step and directly obtain $\theta^*$?

$$\theta^*(\omega) = g_\omega(D_{train}^{support})$$

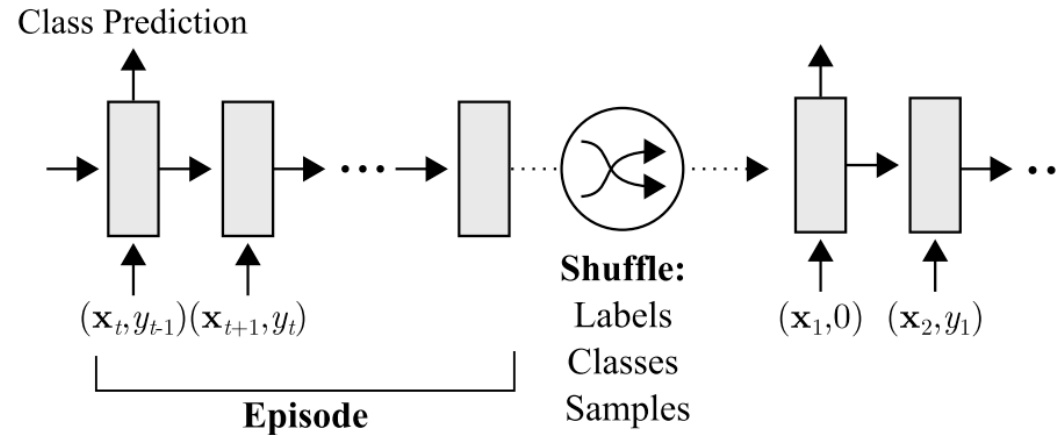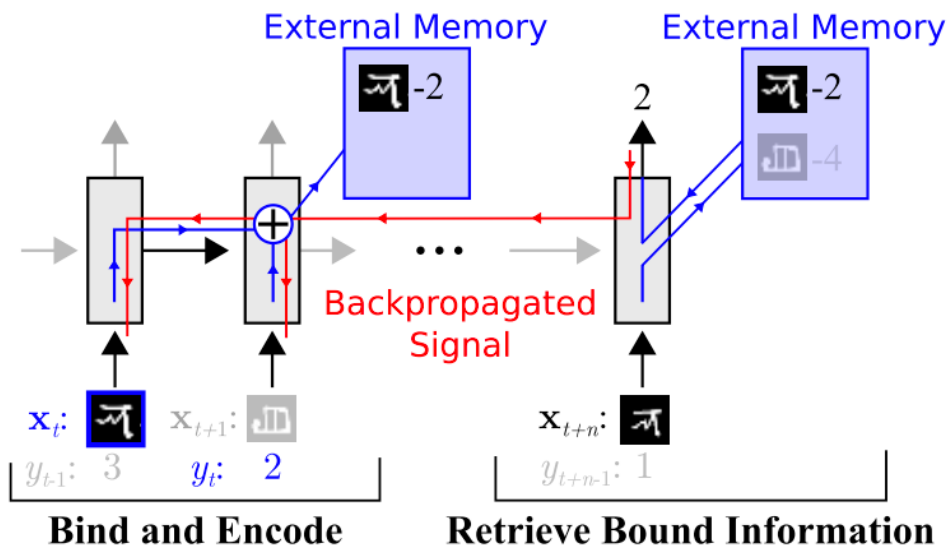Model-based methods adopt the meta-knowledge $\omega$ to directly generate a model.

# Memory-Augmented Neural Networks

厦門大學信息學院（特色化示范性软件学院）
School of Informatics Xiamen University (National Characteristic Demonstration Software School)

厦门大学 计算机科学与技术系
Department of Computer Science and Technology, Xiamen University

Image source: https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html#neural-turing-machines

# Memory-Augmented Neural Networks



Meta-knowledge $\omega$ learns proper representation to read and write with memory.

Image source: A. Santoro, M. Botvinick, T. Lillicrap, G. Deepmind, and W. G. Com, "Meta-Learning with Memory-Augmented Neural Networks," in ICML, 2016, vol. 48.

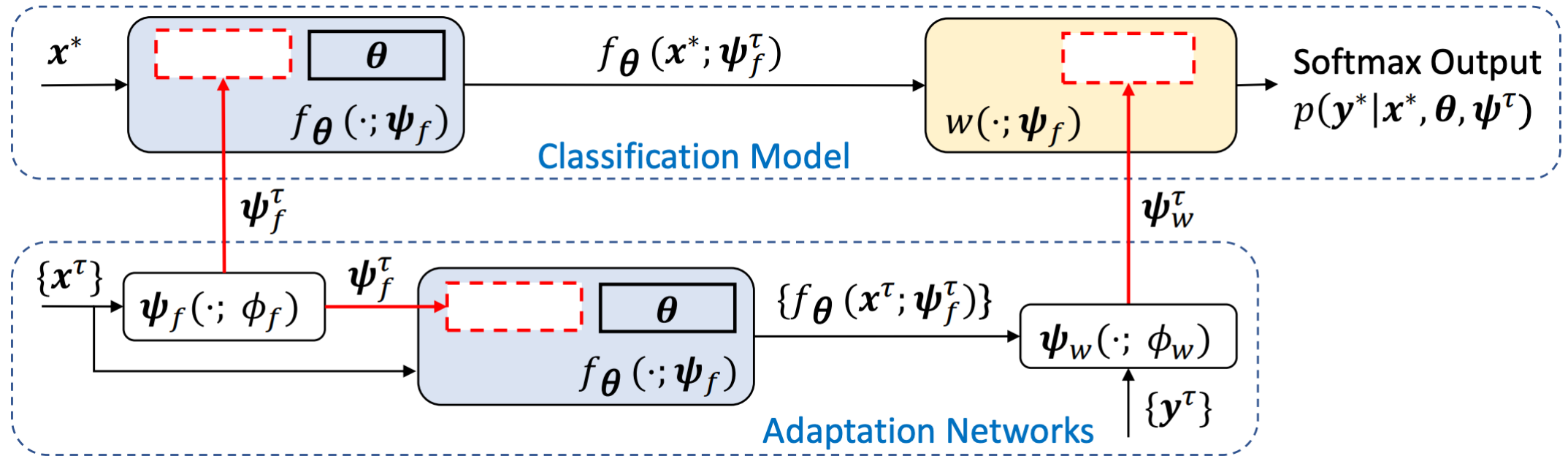# Meta Networks

- Slow weights: weights that are learned from an optimization process like SGD.

- Fast weights: weights that are directly generated by another network.

- In MetaNet, loss gradients are used as meta information to populate models that learn fast weights.

  - Slow and fast weights are combined to make predictions in neural networks.

Image source: T. Munkhdalai and H. Yu, "Meta Networks," in ICML, 2017.

# CNAPs



The model is not generated from scratch. We generate the adaptation model to new tasks, instead of optimize the model to new tasks.
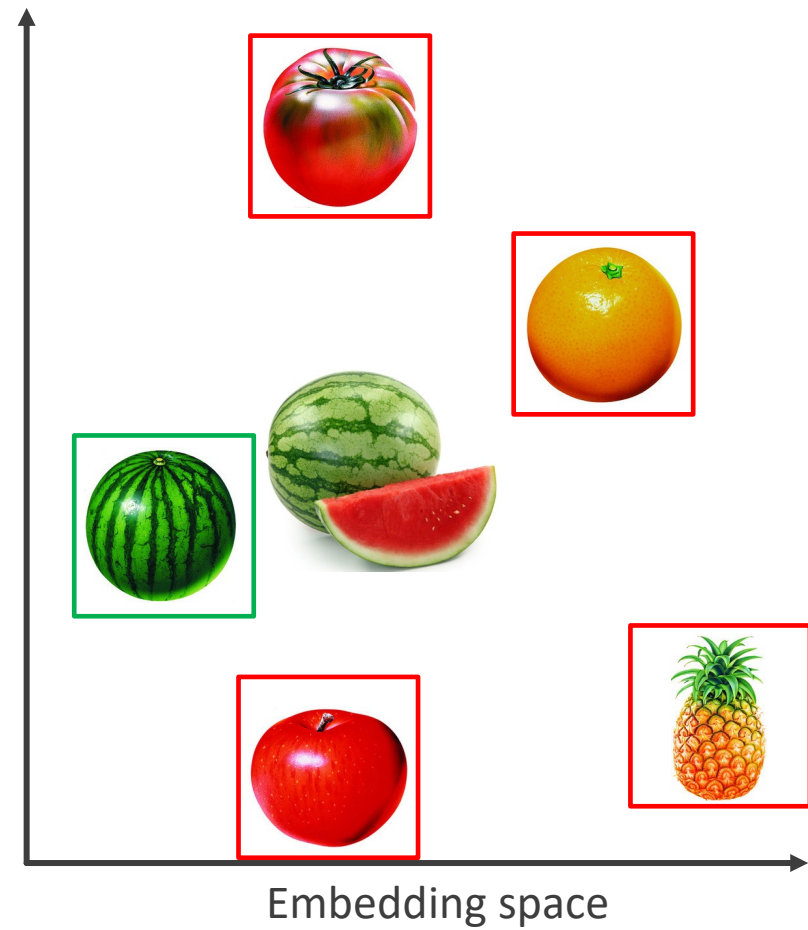
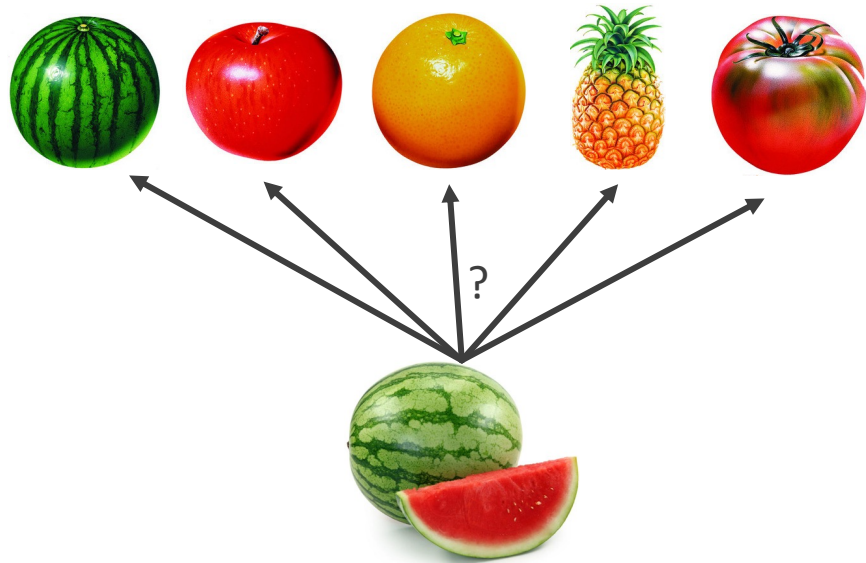# METRIC-BASED METHOD

■Previously, we always have a model $f_\theta$ to output class score $f_\theta(\boldsymbol{x})$, no matter $f_\theta$ is optimized by gradient descend with meta knowledge or directly generated by meta model.

■Do we have to use a model to do prediction? Is there any machine learning method that doesn't have a model?

$$k\text{NN}$$

# Metric-Based Method



Embedding space

- Metric-based methods learn an embedding network $F_\omega$.

- The learned representation is suitable for recognition by simple similarity comparison between query and support instances.

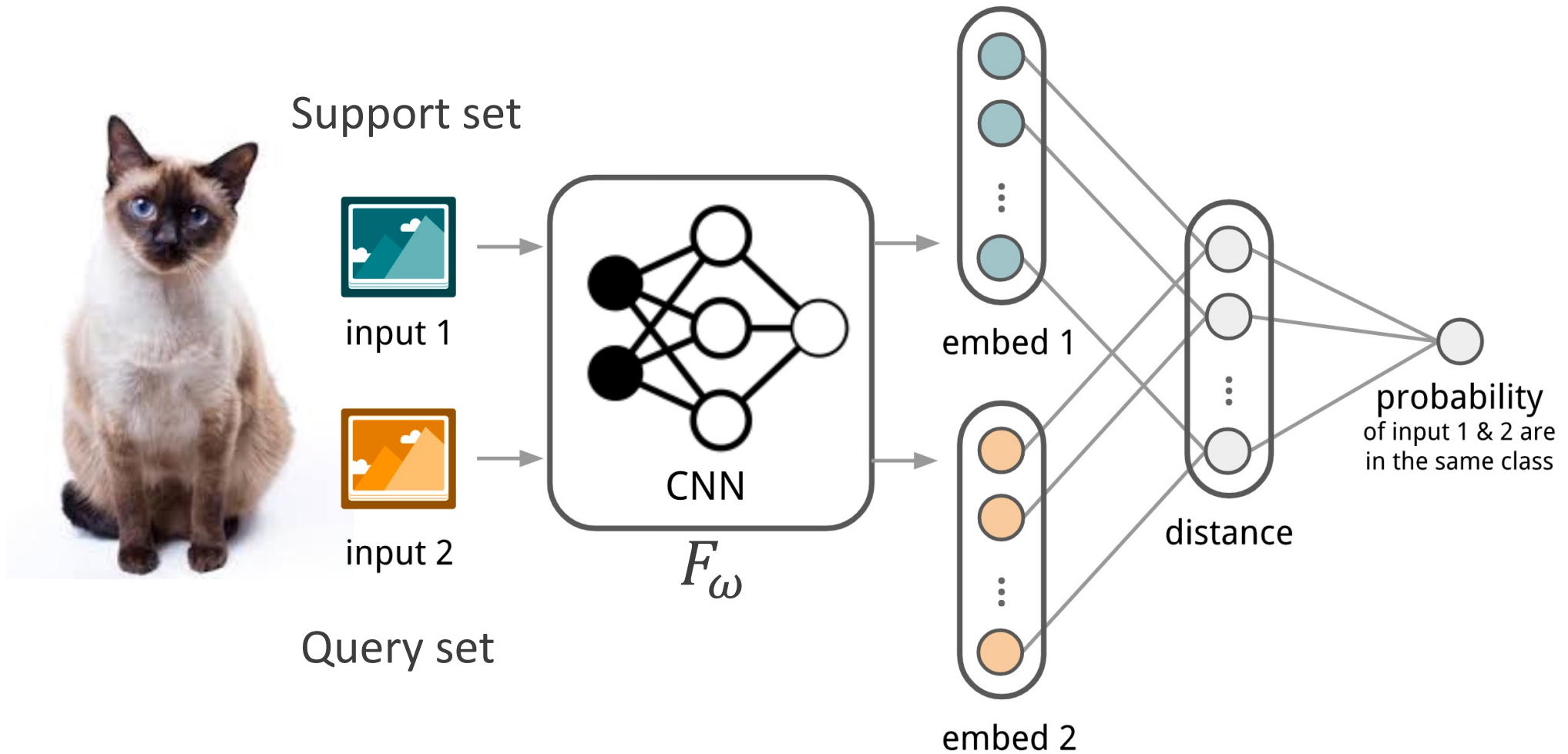- Take one-shot learning as an example, at testing phase, we simply calculating the similarity between:

$$F_\omega(\boldsymbol{x}_{test}^{query}) \quad \text{and} \quad F_\omega(\boldsymbol{x}_{test}^{support(j)})$$

- It can also be treated as model-based method with only one linear layer.

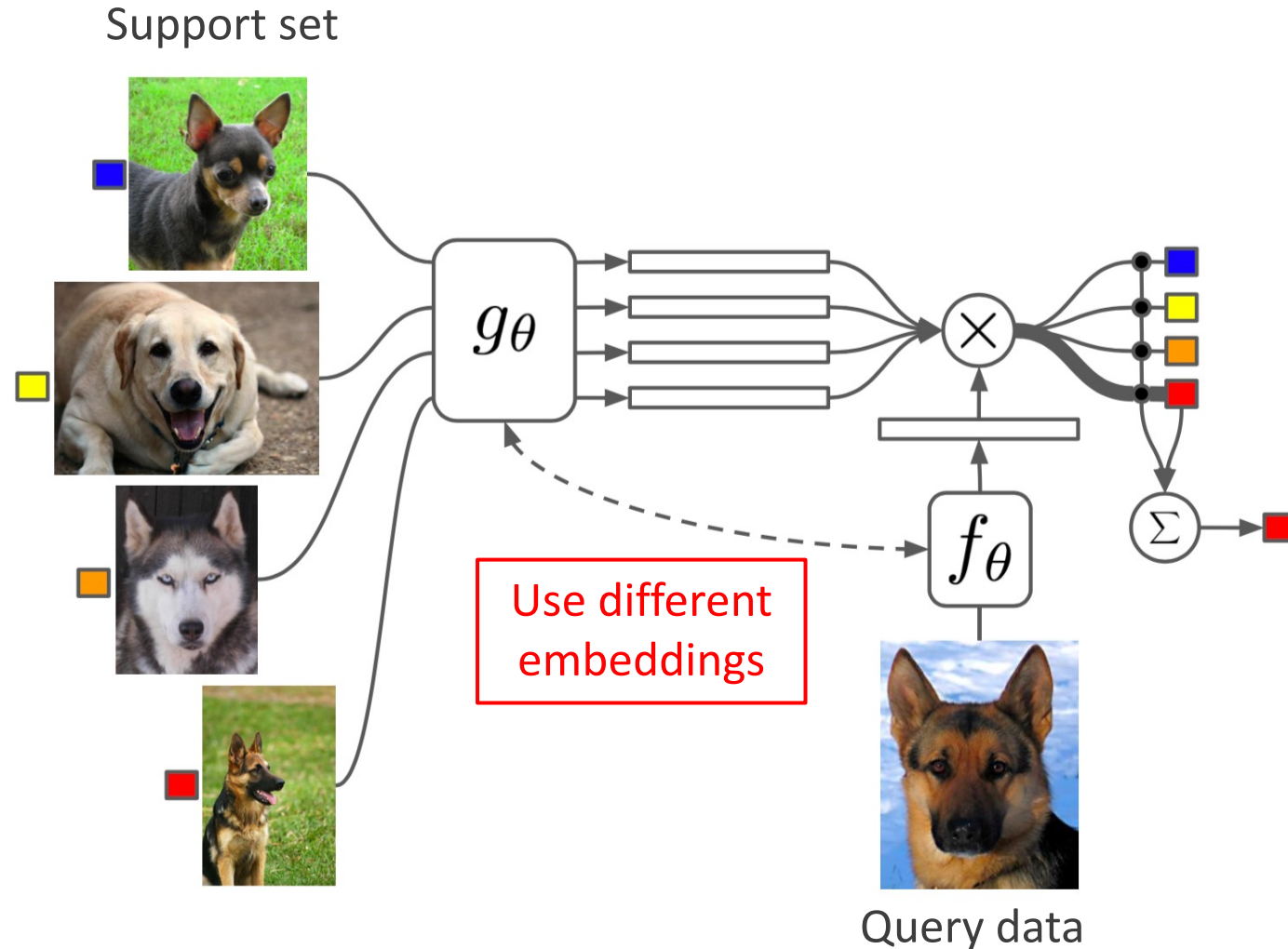# Metric-Based Method

- Siamese networks

- Matching networks

- Prototypical networks

- Relation networks

- Graph networks

Support set

Use different embeddings

$g_\theta$

$f_\theta$

Query data

厦門大學信息学院（特色化示范性软件学院）
School of Informatics Xiamen University (National Characteristic Demonstration Software School)

厦门大学 计算机科学与技术系
Department of Computer Science and Technology, Xiamen University

Image source: O. Vinyals, C. Blundell, and T. Lillicrap, "Matching Networks for One Shot Learning," in NIPS, 2016.

# Matching Networks

- Simple version:

$$a(\widehat{\boldsymbol{x}}, \boldsymbol{x}_i) = \exp\left(cos\left(F_\omega(\widehat{\boldsymbol{x}}), G_\omega(\boldsymbol{x}_i)\right)\right) / \sum_{j=1}^{k} \exp\left(cos\left(F_\omega(\widehat{\boldsymbol{x}}), G_\omega(\boldsymbol{x}_j)\right)\right)$$

- Full context version:

$$\widehat{\boldsymbol{h}}_k, \boldsymbol{c}_k = \text{LSTM}(F(\widehat{\boldsymbol{x}}), [\boldsymbol{h}_{k-1}, \boldsymbol{r}_{k-1}], \boldsymbol{c}_{k-1})$$

$$\boldsymbol{h}_k = \widehat{\boldsymbol{h}}_k + F(\widehat{\boldsymbol{x}})$$

$$\boldsymbol{r}_{k-1} = \sum_{i=1}^{|S|} a(\boldsymbol{h}_{k-1}, G(\boldsymbol{x}_i)) G(\boldsymbol{x}_i)$$

$$a(\boldsymbol{h}_{k-1}, G(\boldsymbol{x}_i)) = \exp\left(\boldsymbol{h}_{k-1}^T G(\boldsymbol{x}_i)\right) / \sum_{j=1}^{|S|} \exp\left(\boldsymbol{h}_{k-1}^T G(\boldsymbol{x}_j)\right)$$

Image source: F. Sung, Y. Yang, and L. Zhang, "Learning to Compare : Relation Network for Few-Shot Learning Queen Mary University of London," in CVPR, 2018, pp. 1199–1208.

厦門大學信息学院（特色化示范性软件学院）
School of Informatics Xiamen University (National Characteristic Demonstration Software School)

厦门大学 计算机科学与技术系
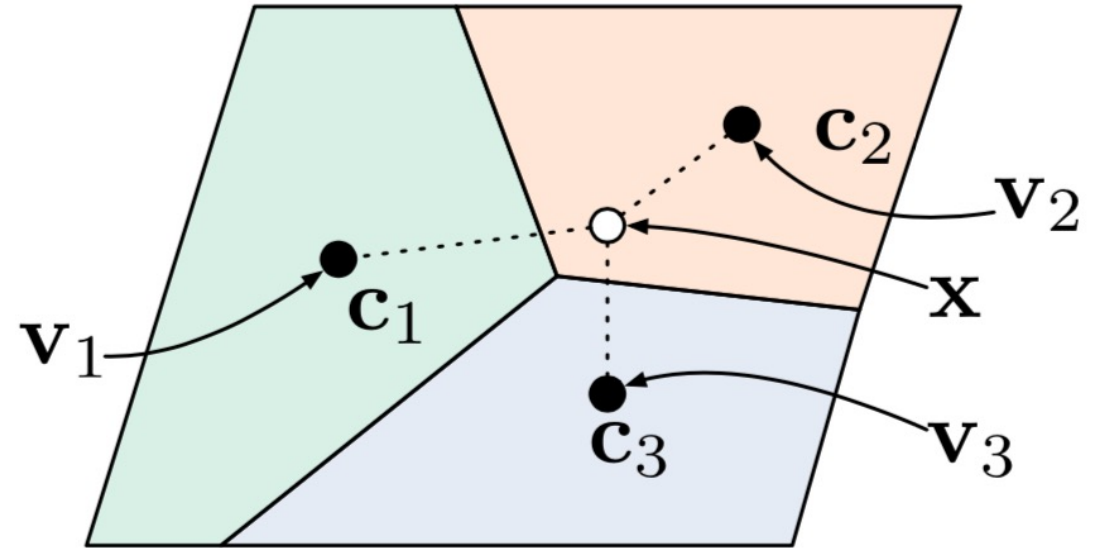Department of Computer Science and Technology, Xiamen University

# Prototypical Networks
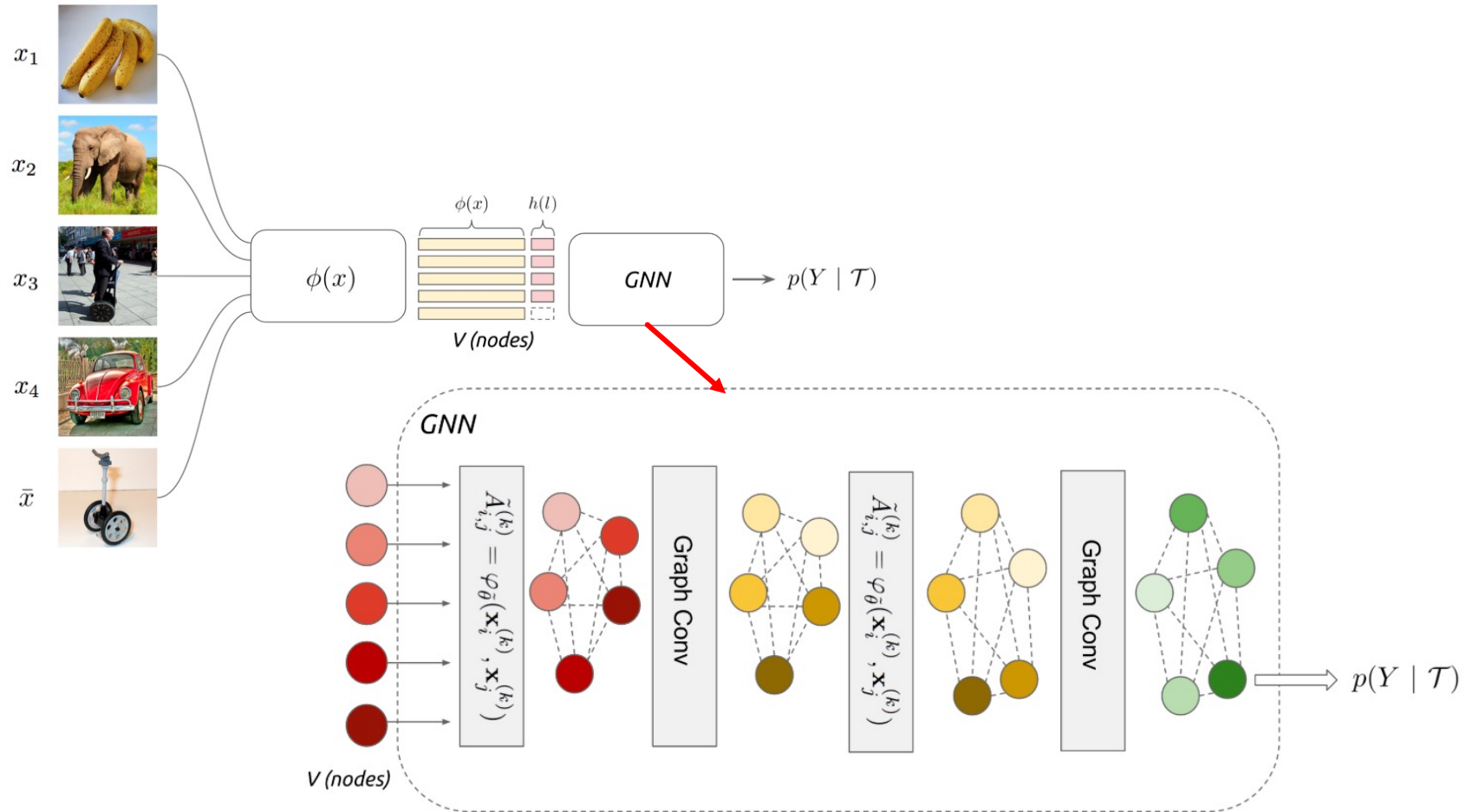


(a) Few-shot

(b) Zero-shot

# Zero-Shot Learning

- One-shot is difficult enough. How can zero-shot work?

- Zero-shot learning differs from few-shot learning in that instead of being given a support set of training points, we are given a class meta-data vector $v_k$ for each class.

    - For example, $v_k$ can be a sentence embedding for text description of the image.

    - The "zero" in zero-shot is for the labelled support set, but we can utilize other information instead.

- We can simply build a mapping between meta-data vector $v_k$ to its prototype:

$$c_k = g_\omega(v_k).$$

厦門大學信息学院（特色化示范性软件学院）
School of Informatics Xiamen University (National Characteristic Demonstration Software School)

厦门大学 计算机科学与技术系
Department of Computer Science and Technology, Xiamen University

# Graph Networks

Generally, two steps:

- Build an embedding network $F_\omega$.

- Model the relation between $F_\omega(\hat{x})$ and $F_\omega(x_i)$, and output the probabilities.

After this lecture, you should know:

- What is meta-learning and learning to learn.

- How can we utilize meta-knowledge.

- What is few-shot learning and $N$-way-$k$-shot.

- What is the difference between optimization-based, model-based and metric-based meta-learning.

# Reference & Suggested Reading

- Meta Learning: Learn to Learn

- CS 330: Deep Multi-Task and Meta Learning

- Meta-Learning: Learning to Learn Fast

- T. Hospedales, A. Antoniou, P. Micaelli, and A. Storkey, "Meta-Learning in Neural Networks: A Survey," in *arXiv*, 2020.

- Any question?

- Don't hesitate to send email to me for asking questions and discussion. ☺

厦門大學信息学院（特色化示范性软件学院）
School of Informatics Xiamen University (National Characteristic Demonstration Software School)

厦门大学 计算机科学与技术系
Department of Computer Science and Technology, Xiamen University