# Chapter 16 Meta Learning

# Introduction
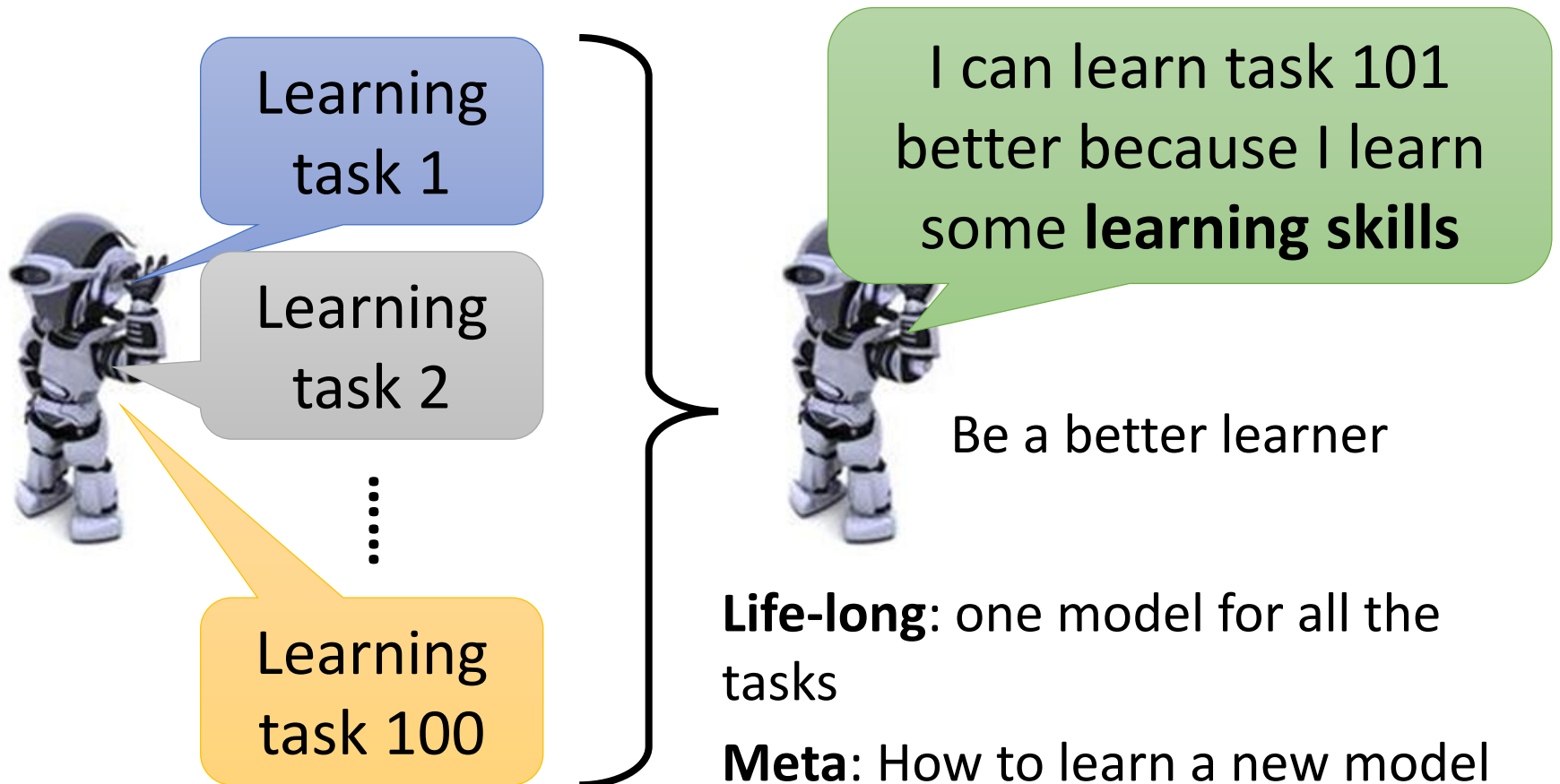
Task 1: speech recognition

Task 2: image recognition

⋮

Task 100: text classification
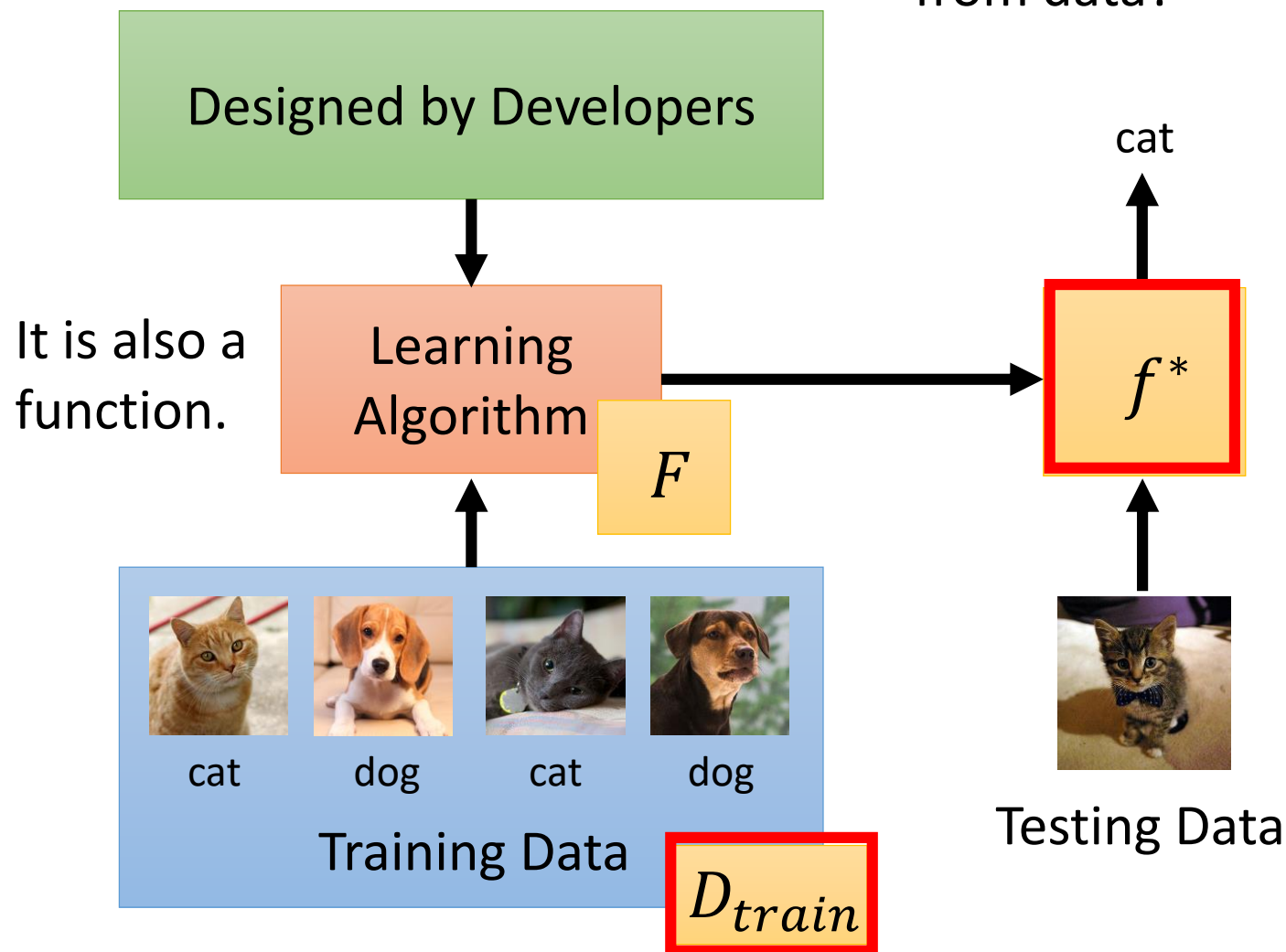
- Meta learning = Learn to learn

Learning task 1

Learning task 2

⋮

Learning task 100

I can learn task 101 better because I learn some **learning skills**

Be a better learner

**Life-long**: one model for all the tasks

**Meta**: How to learn a new model

# Meta Learning

$$f^* = F(D_{train})$$

Can machine find $F$ from data?

Designed by Developers

It is also a function.

Learning Algorithm

$F$

$f^*$

cat



cat    dog    cat    dog

Training Data

$D_{train}$

Testing Data

# Meta Learning

**Machine Learning** ≈ 根据资料找一個函数 f 的能力

$$f(\quad \text{🐱} \quad) = \text{"Cat"}$$

**Meta Learning**

≈ 根据资料找一個找一個函數 f 的函数 F 的能力

$$F(\quad) = f^*$$

cat    dog    cat    dog

**Training Data**

# ~~Machine~~ Learning is Simple
## Meta

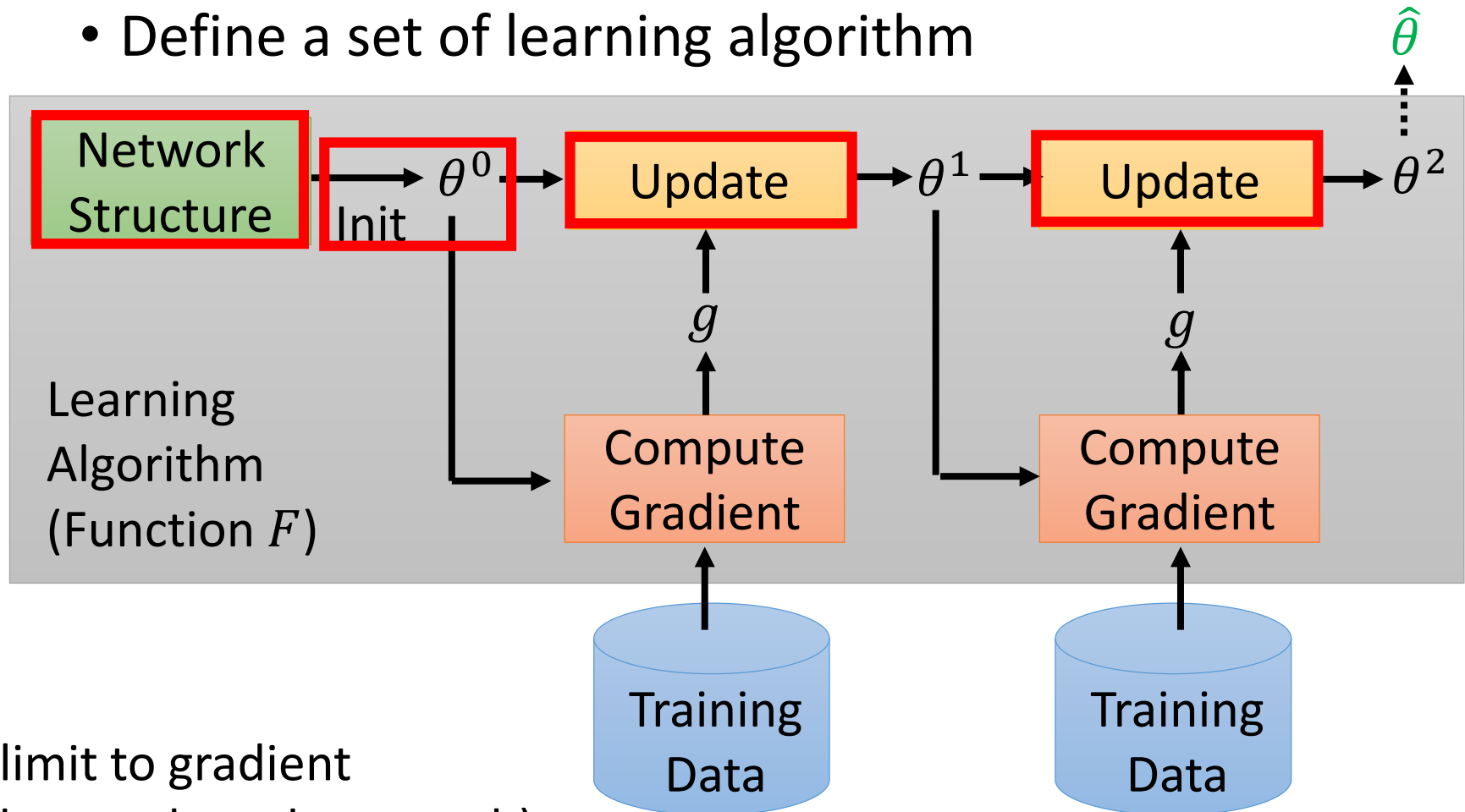| Step 1: define a set of ~~function~~ | → | Step 2: goodness of ~~function~~ | → | Step 3: pick the best ~~function~~ |
|---|---|---|---|---|

Function $f$ ➡ Learning algorithm $F$

就好像把大象放进冰箱 ……

# Meta Learning

Different decisions in the red boxes lead to different algorithms. What happens in the red boxes is decided by humans until now.

- Define a set of learning algorithm



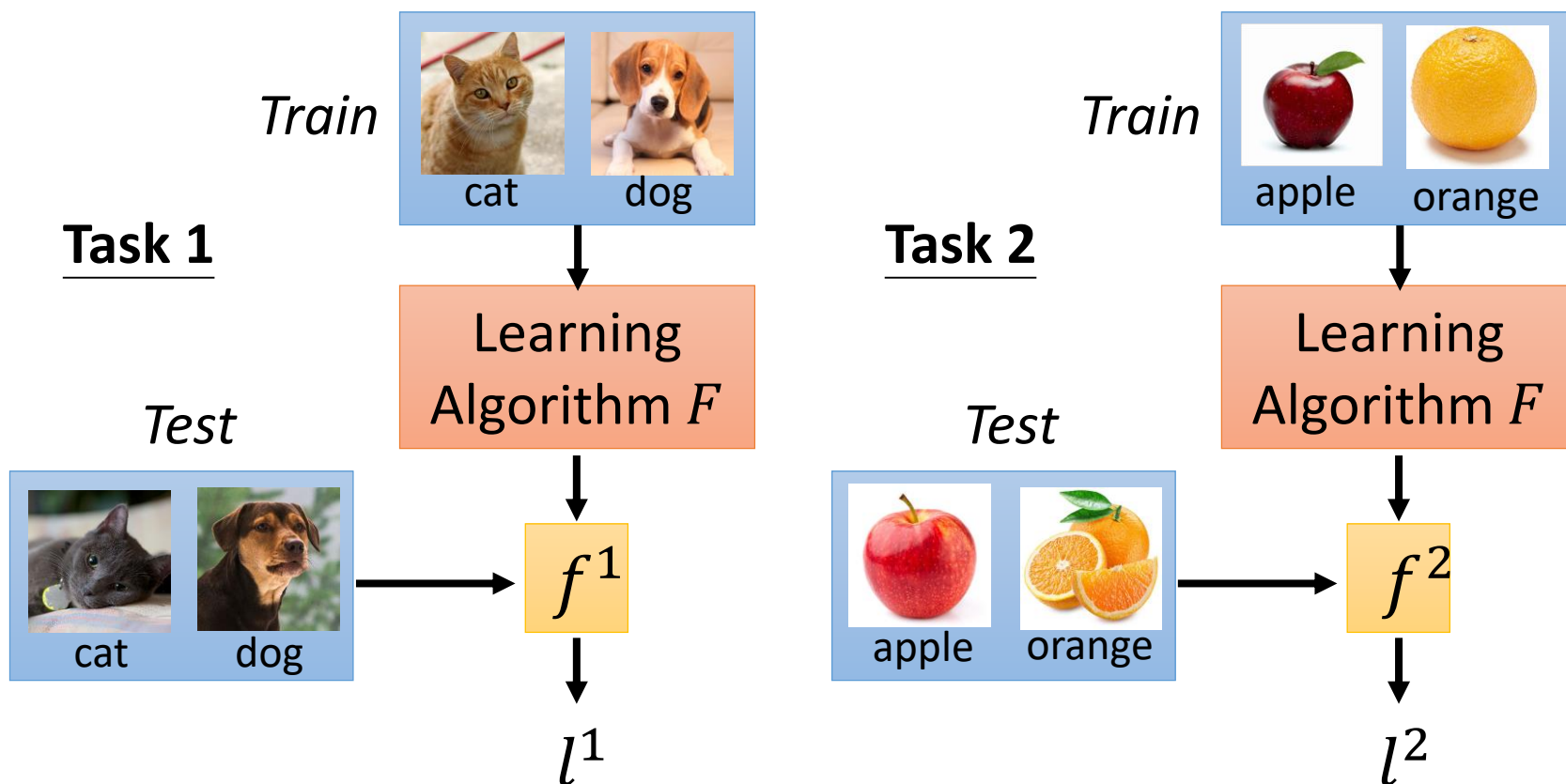Learning Algorithm (Function $F$)

(limit to gradient descent based approach)

# Meta Learning

$$L(F) = \sum_{n=1}^{N} l^n$$

N tasks

Testing loss for task n after training

- Defining the goodness of a function $F$



**Task 1**

*Train*

cat  dog

Learning Algorithm $F$

*Test*

cat  dog

$f^1$

$l^1$

**Task 2**

*Train*

apple  orange

Learning Algorithm $F$

*Test*

apple  orange

$f^2$

$l^2$

Machine Learning

Meta Learning

Widely considered in **few-shot learning**

Training Tasks

Task 1

Task 2

Sometimes you need validation tasks

Testing Tasks

Support set

Query set

Train    Test

cat    dog    cat    dog

Train    Test

apple    orange    apple    orange

Train    Test

bike    car    bike    car

# Meta Learning

- Defining the goodness of a function $F$

$$L(F) = \sum_{n=1}^{N} l^n$$

**Testing:**
**Task New**

- Find the best function $F^*$

$$F^* = arg \min_{F} L(F)$$



*Train*

bike    car

Learning Algorithm $F^*$

*Test*

bike    car

$f^*$

$l$

# Omniglot

- 1623 characters

- Each has 20 examples

Tagalog
character 1

# Omniglot
# – Few-shot Classification

- N-ways K-shot classification: In each training and test tasks, there are N classes, each has K examples.

**_20 ways_**
**_1 shot_**

Each character represents a class



Testing set (Query set)

Training set (Support set)

- Split your characters into training and testing characters
    - Sample N training characters, sample K examples from each sampled characters → one training task
    - Sample N testing characters, sample K examples from each sampled characters → one testing task

# Techniques Today

- **MAML**
  - Chelsea Finn, Pieter Abbeel, and Sergey Levine, "Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks", ICML, 2017

- **Reptile**
  - Alex Nichol, Joshua Achiam, John Schulman, On First-Order Meta-Learning Algorithms, arXiv, 2018

# MAML

$\hat{\theta}^n$: model learned from task $n$

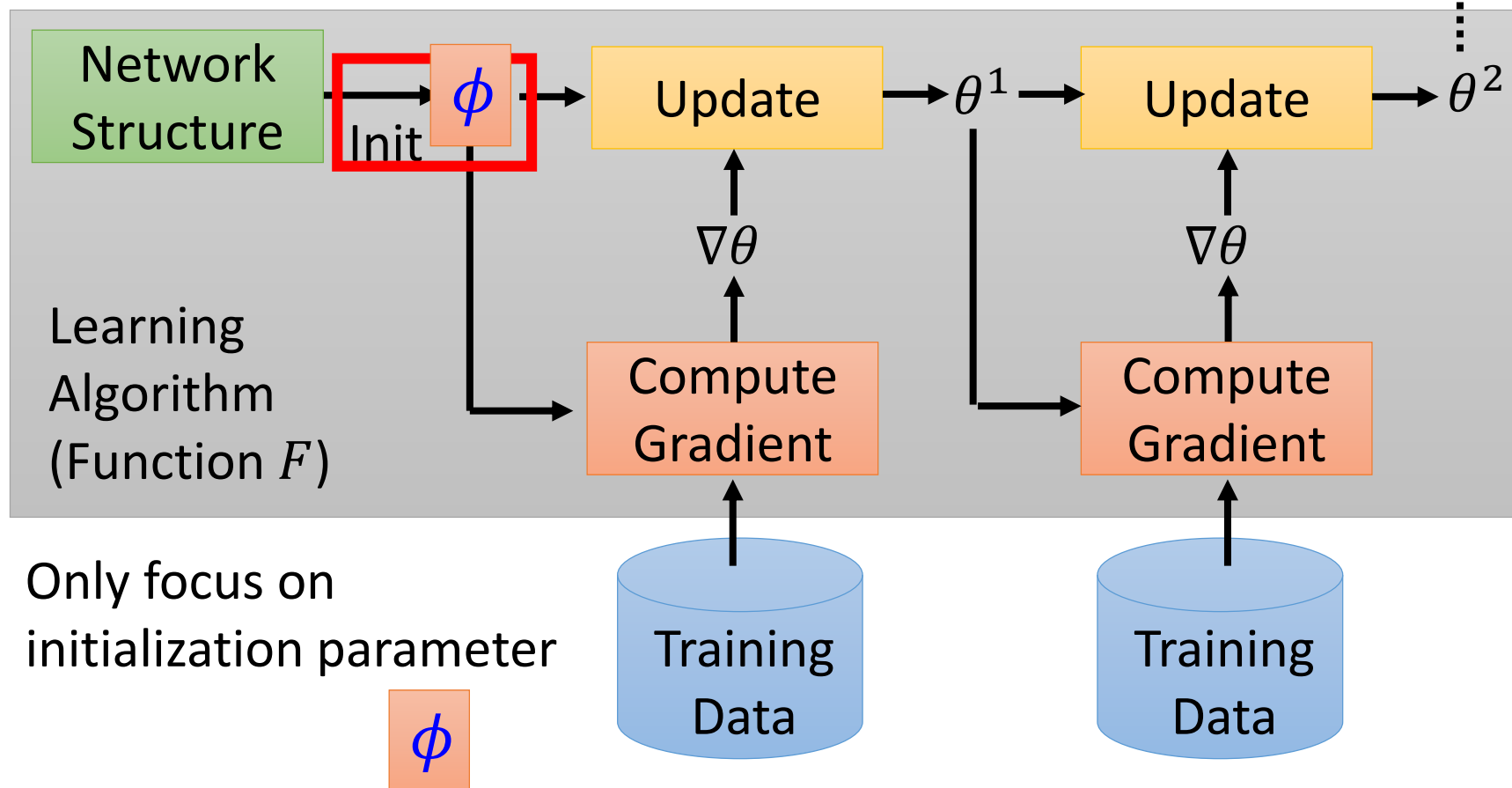$\hat{\theta}^n$ depends on $\phi$

Loss Function:

$$L(\phi) = \sum_{n=1}^{N} l^n(\hat{\theta}^n)$$

$l^n(\hat{\theta}^n)$: loss of task $n$ on the testing set of task $n$



Learning Algorithm (Function $F$)

Network Structure — Init — $\phi$ — Update — $\theta^1$ — Update — $\theta^2$

$\nabla\theta$

Compute Gradient

Training Data

$\hat{\theta}$

Only focus on initialization parameter

$\phi$

# *MAML*

$\hat{\theta}^n$: model learned from task $n$

Loss Function:

$\hat{\theta}^n$ depends on $\phi$

$$L(\phi) = \sum_{n=1}^{N} l^n(\hat{\theta}^n)$$

$l^n(\hat{\theta}^n)$: loss of task $n$ on the testing set of task $n$

How to minimize $L(\phi)$?  Gradient Descent

$$\phi \leftarrow \phi - \eta \nabla_\phi L(\phi)$$

---

# *Model Pre-training*

Widely used in transfer learning

Loss Function:

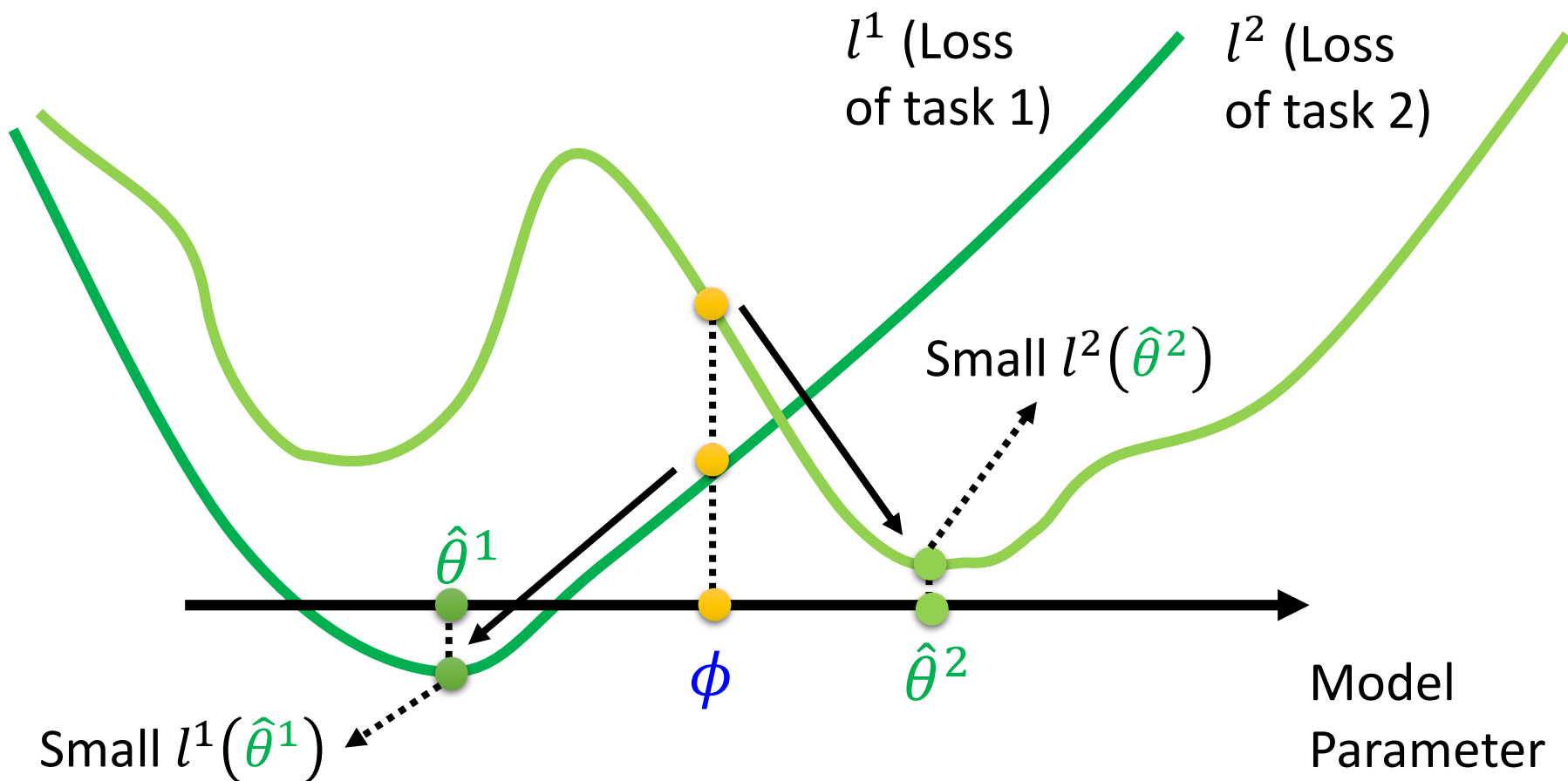$$L(\phi) = \sum_{n=1}^{N} l^n(\phi)$$

# *MAML*

$$L(\phi) = \sum_{n=1}^{N} l^n(\hat{\theta}^n)$$

$l^1$ (Loss of task 1)

$l^2$ (Loss of task 2)

Small $l^2(\hat{\theta}^2)$

$\hat{\theta}^1$

$\phi$

$\hat{\theta}^2$

Model Parameter

Small $l^1(\hat{\theta}^1)$

# Model Pre-training

$$L(\phi) = \sum_{n=1}^{N} l^n(\phi)$$

找寻在所有 task 都最好的 $\phi$

并不保证拿 $\phi$ 去训练以后会得到好的 $\hat{\theta}^n$

$l^2(\hat{\theta}^2)$

$l^1$ (Loss of task 1)

$l^2$ (Loss of task 2)

$\phi$

Model Parameter

## _MAML_

$\hat{\theta}^n$: model learned from task $n$

$\hat{\theta}^n$ depends on $\phi$

Loss Function:

$$L(\phi) = \sum_{n=1}^{N} l^n(\hat{\theta}^n)$$

$l^n(\hat{\theta}^n)$: loss of task $n$ on the testing set of task $n$

How to minimize $L(\phi)$?   Gradient Descent

$$\phi \leftarrow \phi - \eta \nabla_\phi L(\phi)$$

Find $\phi$ achieving good performance **after training** | 潜力

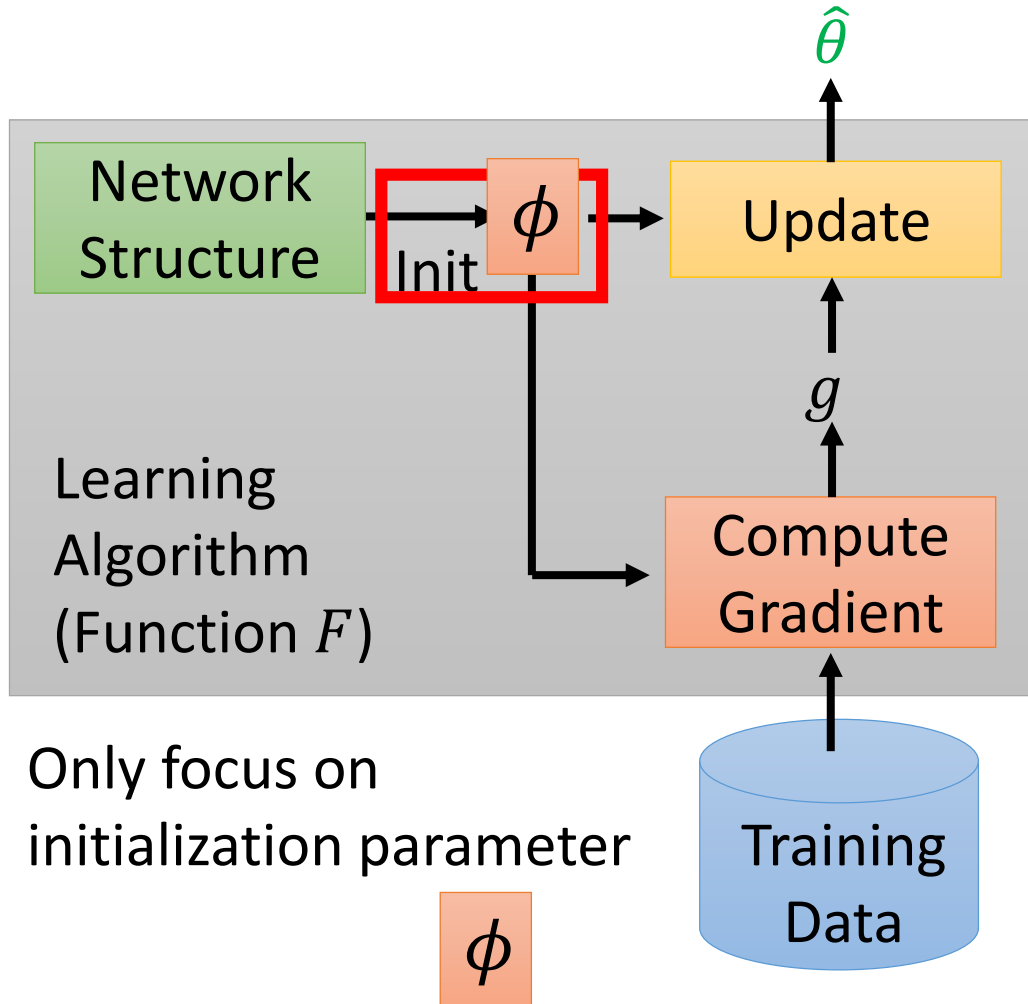## _Model Pre-training_

Widely used in transfer learning

Loss Function:

$$L(\phi) = \sum_{n=1}^{N} l^n(\phi)$$

Find $\phi$ achieving good performance | 现在表现如何

# MAML



Only focus on initialization parameter $\phi$

$$L(\phi) = \sum_{n=1}^{N} l^n(\hat{\theta}^n)$$

$$\phi \leftarrow \phi - \eta \nabla_{\phi} L(\phi)$$

Considering one-step training:

$$\hat{\theta} = \phi - \varepsilon \nabla_{\phi} l(\phi)$$

# Warning of Math

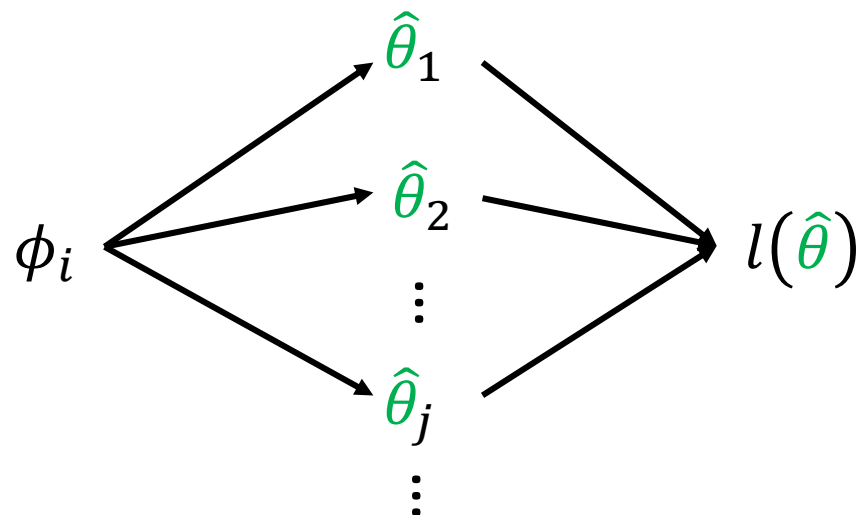$$\phi \leftarrow \phi - \eta \nabla_\phi L(\phi)$$

$$L(\phi) = \sum_{n=1}^{N} l^n(\hat{\theta}^n)$$

$$\hat{\theta} = \phi - \varepsilon \nabla_\phi l(\phi)$$

$$\nabla_\phi L(\phi) = \nabla_\phi \sum_{n=1}^{N} l^n(\hat{\theta}^n) = \sum_{n=1}^{N} \underline{\nabla_\phi l^n(\hat{\theta}^n)}$$

$$\frac{\partial l(\hat{\theta})}{\partial \phi_i} = \sum_j \frac{\partial l(\hat{\theta})}{\partial \hat{\theta}_j} \frac{\partial \hat{\theta}_j}{\partial \phi_i}$$

$$\nabla_\phi l(\hat{\theta}) = \begin{bmatrix} \partial l(\hat{\theta})/\partial \phi_1 \\ \partial l(\hat{\theta})/\partial \phi_2 \\ \vdots \\ \boxed{\partial l(\hat{\theta})/\partial \phi_i} \\ \vdots \end{bmatrix}$$

$$\phi \leftarrow \phi - \eta \nabla_\phi L(\phi)$$

$$L(\phi) = \sum_{n=1}^{N} l^n(\hat{\theta}^n)$$

$$\boxed{\hat{\theta} = \phi - \varepsilon \nabla_\phi l(\phi)}$$

$$\nabla_\phi L(\phi) = \nabla_\phi \sum_{n=1}^{N} l^n(\hat{\theta}^n) = \sum_{n=1}^{N} \underline{\nabla_\phi l^n(\hat{\theta}^n)}$$

$$\frac{\partial l(\hat{\theta})}{\partial \phi_i} = \sum_j \frac{\partial l(\hat{\theta})}{\partial \hat{\theta}_j} \boxed{\frac{\partial \hat{\theta}_j}{\partial \phi_i}} \approx \frac{\partial l(\hat{\theta})}{\partial \hat{\theta}_i}$$

$$\hat{\theta}_j = \phi_j - \varepsilon \frac{\partial l(\phi)}{\partial \phi_j}$$

$$\nabla_\phi l(\hat{\theta}) = \begin{bmatrix} \partial l(\hat{\theta})/\partial \phi_1 \\ \partial l(\hat{\theta})/\partial \phi_2 \\ \vdots \\ \boxed{\partial l(\hat{\theta})/\partial \phi_i} \\ \vdots \end{bmatrix}$$

$i \neq j:$

$$\frac{\partial \hat{\theta}_j}{\partial \phi_i} = -\varepsilon \frac{\partial l(\phi)}{\partial \phi_i \partial \phi_j} \approx 0$$
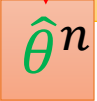
$i = j:$

$$\frac{\partial \hat{\theta}_j}{\partial \phi_i} = 1 - \varepsilon \frac{\partial l(\phi)}{\partial \phi_i \partial \phi_j} \approx 1$$

$$\phi \leftarrow \phi - \eta \nabla_\phi L(\phi)$$

$$L(\phi) = \sum_{n=1}^{N} l^n(\hat{\theta}^n)$$

$$\hat{\theta} = \phi - \varepsilon \nabla_\phi l(\phi)$$

$$\nabla_\phi L(\phi) = \nabla_\phi \sum_{n=1}^{N} l^n(\hat{\theta}^n) = \sum_{n=1}^{N} \nabla_\phi l^n(\hat{\theta}^n)$$
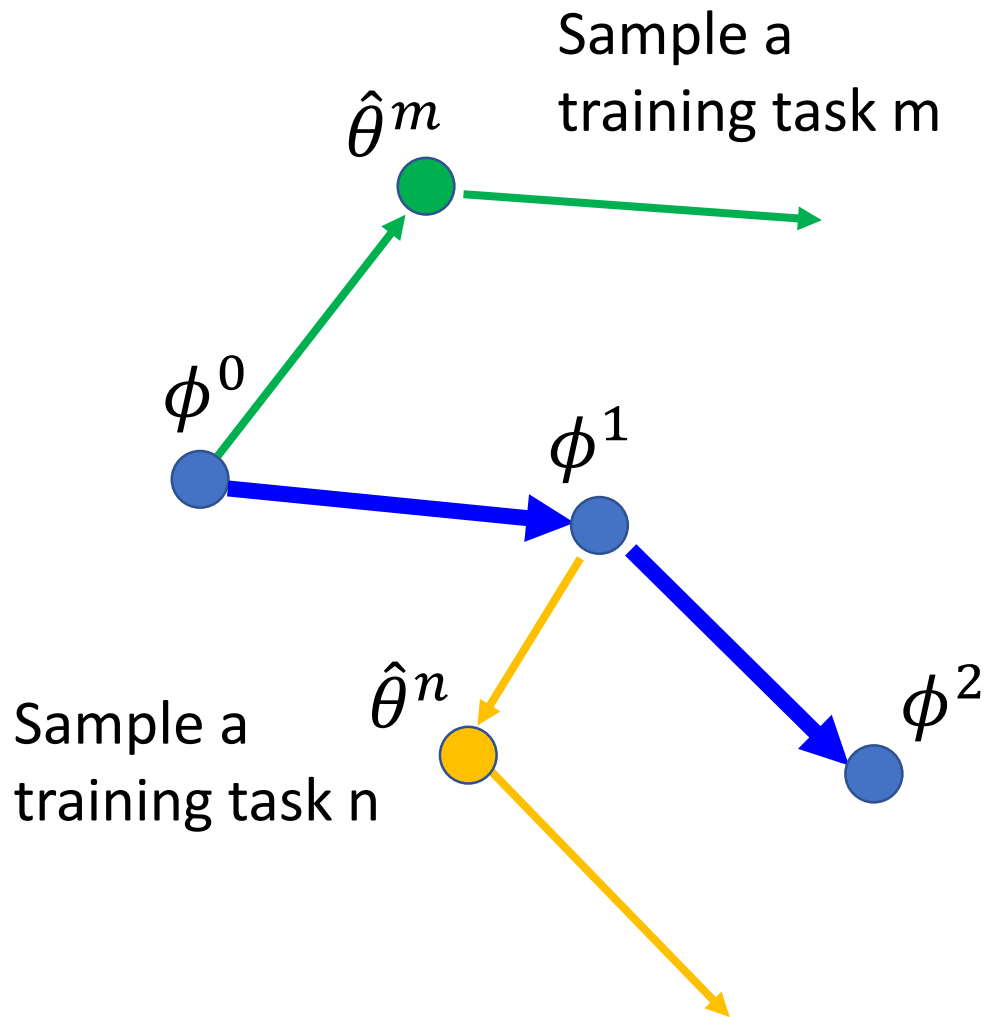
$$\hat{\theta}^n$$

$$\frac{\partial l(\hat{\theta})}{\partial \phi_i} = \sum_j \frac{\partial l(\hat{\theta})}{\partial \hat{\theta}_j} \frac{\partial \hat{\theta}_j}{\partial \phi_i} \approx \frac{\partial l(\hat{\theta})}{\partial \hat{\theta}_i}$$

products, which is supported by standard deep learning libraries such as TensorFlow (Abadi et al., 2016). In our experiments, we also include a comparison to dropping this backward pass and using a first-order approximation, which we discuss in Section 5.2.
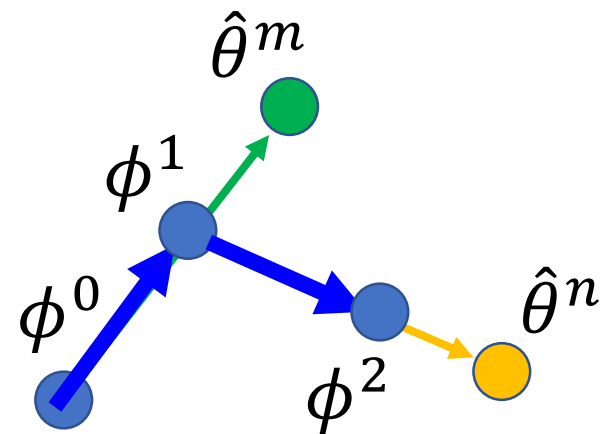
$$\nabla_\phi l(\hat{\theta}) = \begin{bmatrix} \partial l(\hat{\theta})/\partial \phi_1 \\ \partial l(\hat{\theta})/\partial \phi_2 \\ \vdots \\ \partial l(\hat{\theta})/\partial \phi_i \\ \vdots \end{bmatrix} = \begin{bmatrix} \partial l(\hat{\theta})/\partial \hat{\theta}_1 \\ \partial l(\hat{\theta})/\partial \hat{\theta}_2 \\ \vdots \\ \partial l(\hat{\theta})/\partial \hat{\theta}_i \\ \vdots \end{bmatrix} = \nabla_{\hat{\theta}} l(\hat{\theta})$$

# End of Warning
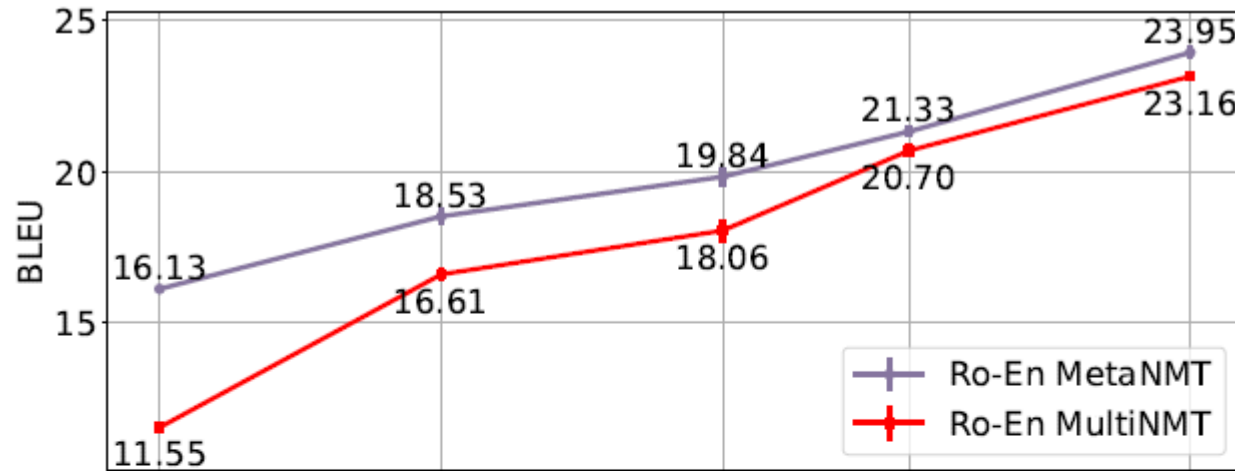
# MAML – Real Implementation



$\hat{\theta}^m$

Sample a training task m

$\phi^0$

$\phi^1$

$\phi^2$

$\hat{\theta}^n$

Sample a training task n

**_Model Pre-training_**

$\hat{\theta}^m$

$\phi^1$
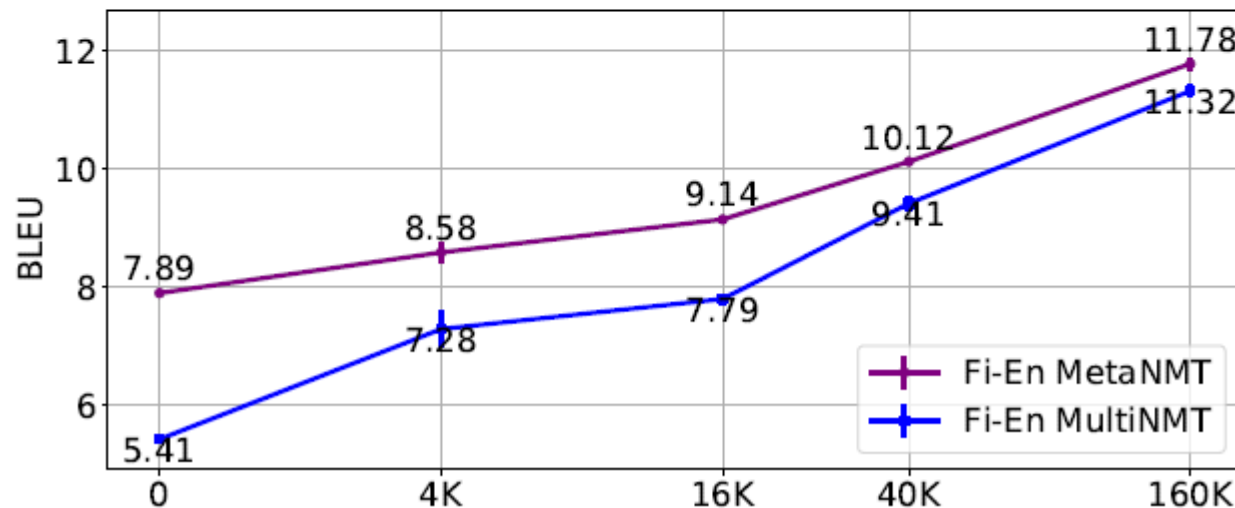
$\phi^0$

$\phi^2$

$\hat{\theta}^n$

# Translation

18 training tasks: 18 different languages translating to English

2 validation tasks: 2 different languages translating to English



Ro = Romanian

Fi = Finnish

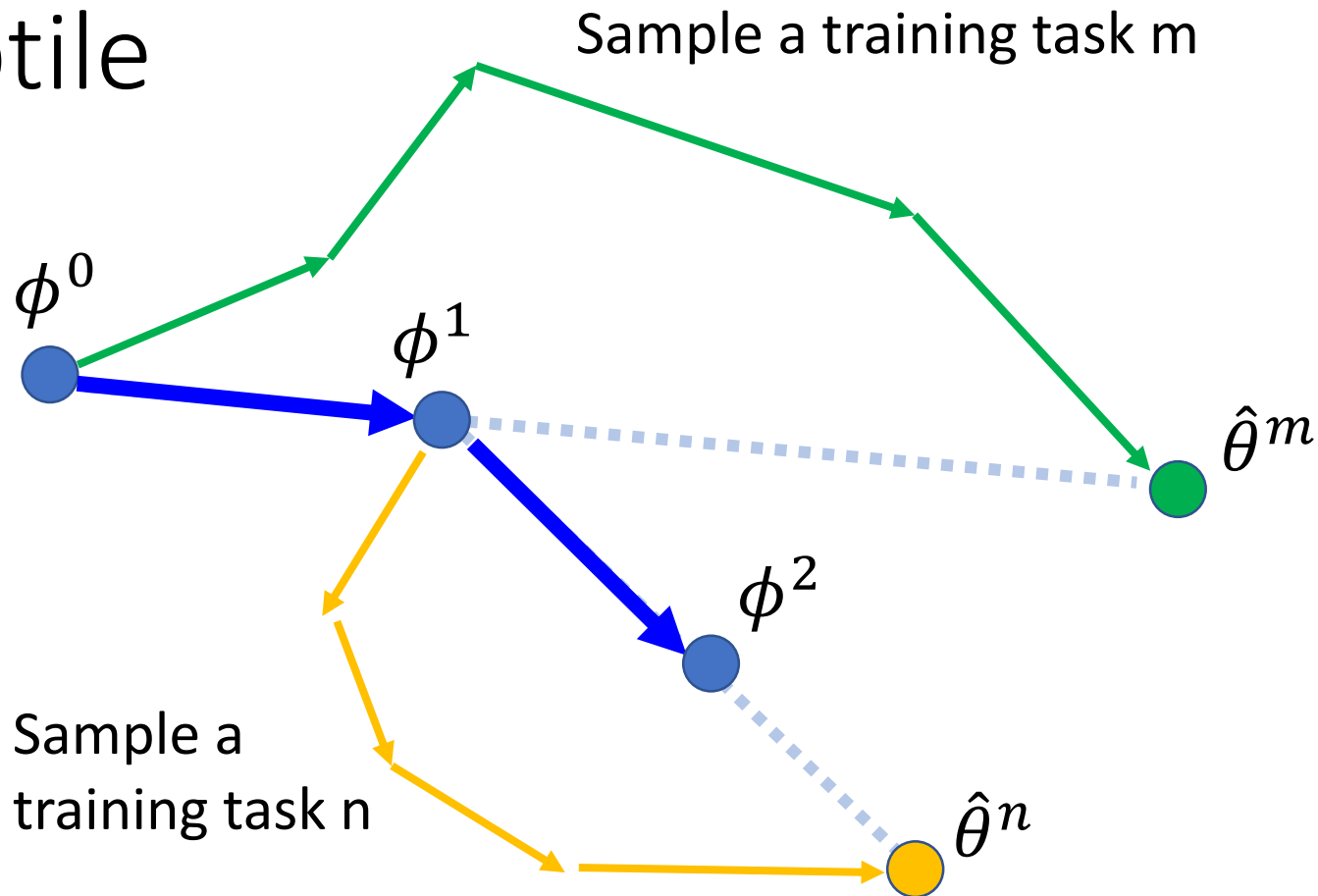# Techniques Today

- **MAML**
  - Chelsea Finn, Pieter Abbeel, and Sergey Levine, "Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks", ICML, 2017

- **Reptile**
  - Alex Nichol, Joshua Achiam, John Schulman, On First-Order Meta-Learning Algorithms, arXiv, 2018

https://openai.com/blog/reptile/

# Reptile



Sample a training task m

$\phi^0$

$\phi^1$

$\hat\theta^m$

$\phi^2$

Sample a
training task n

$\hat\theta^n$

You might be thinking "isn't this the same as training on the expected loss $\mathbb{E}_\tau[L_\tau]$?" and then checking if the date is April 1st. Indeed, if the partial minimization consists of a single gradient step, then this algorithm corresponds to minimizing the expected loss:
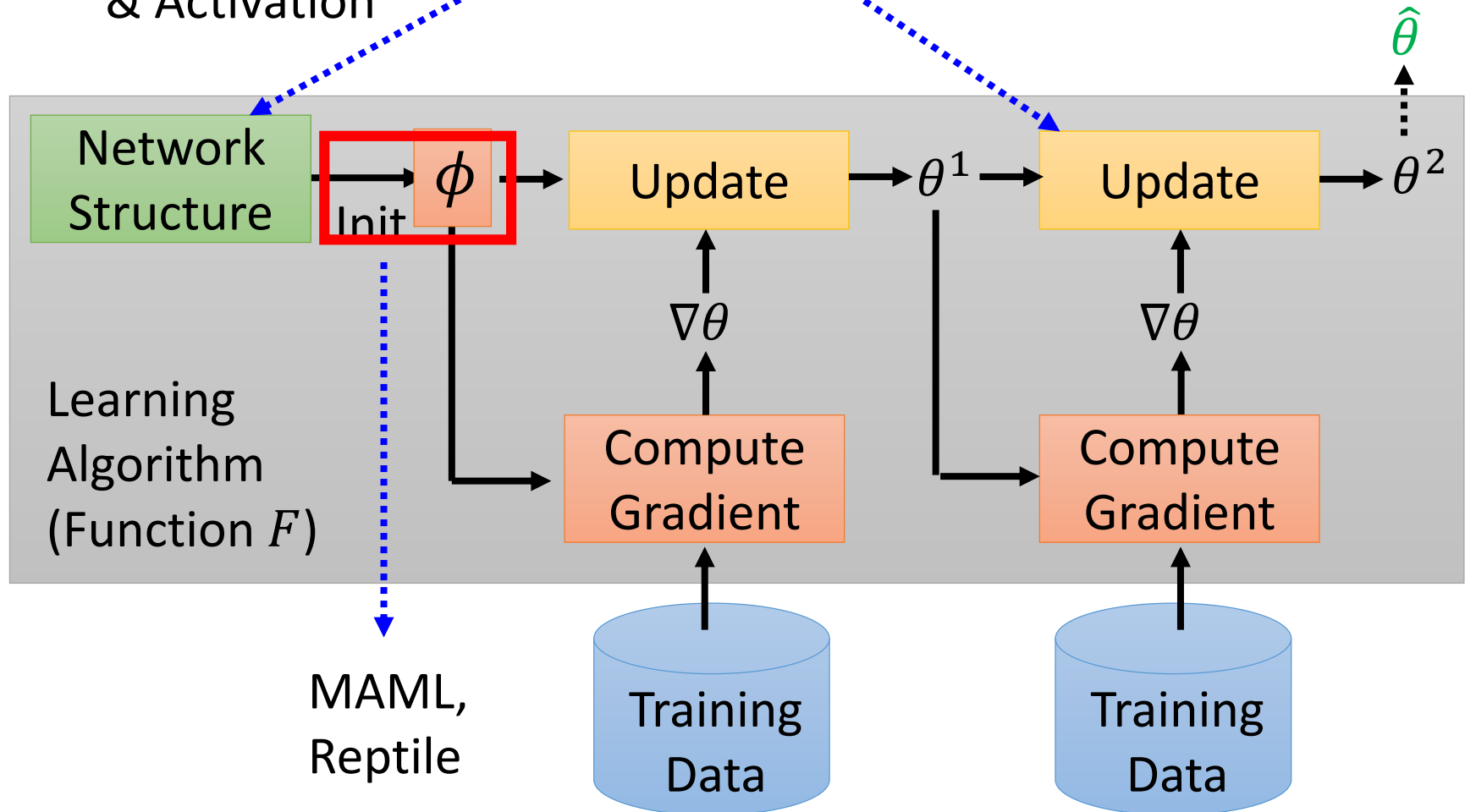
(this sentence is removed in the updated version)

# *More ...*  Video: https://www.youtube.com/watch?v=c10nxBcSH14

Training a network (by RL) to determine ...

Architecture
& Activation

How to update

$\hat{\theta}$

Network
Structure

Init  $\phi$  → Update  →  $\theta^1$  → Update  →  $\theta^2$

Learning
Algorithm
(Function $F$)

$\nabla\theta$

$\nabla\theta$

MAML,
Reptile

Compute
Gradient

Compute
Gradient

Training
Data

Training
Data

# Turtles all the way down ...... ?



- We learn the initialization parameter $\phi$ by gradient descent

- What is the initialization parameter $\phi^0$ for initialization parameter $\phi$?

Learn

Learn to learn

Learn to learn to learn

# Crazy Idea?

- How about learning algorithm beyond gradient descent?



Just a network

Design network architecture = Design training algorithm?

Learn an even bigger function