

Meta Learning (Part 1)

Hung-yi Lee

Introduction

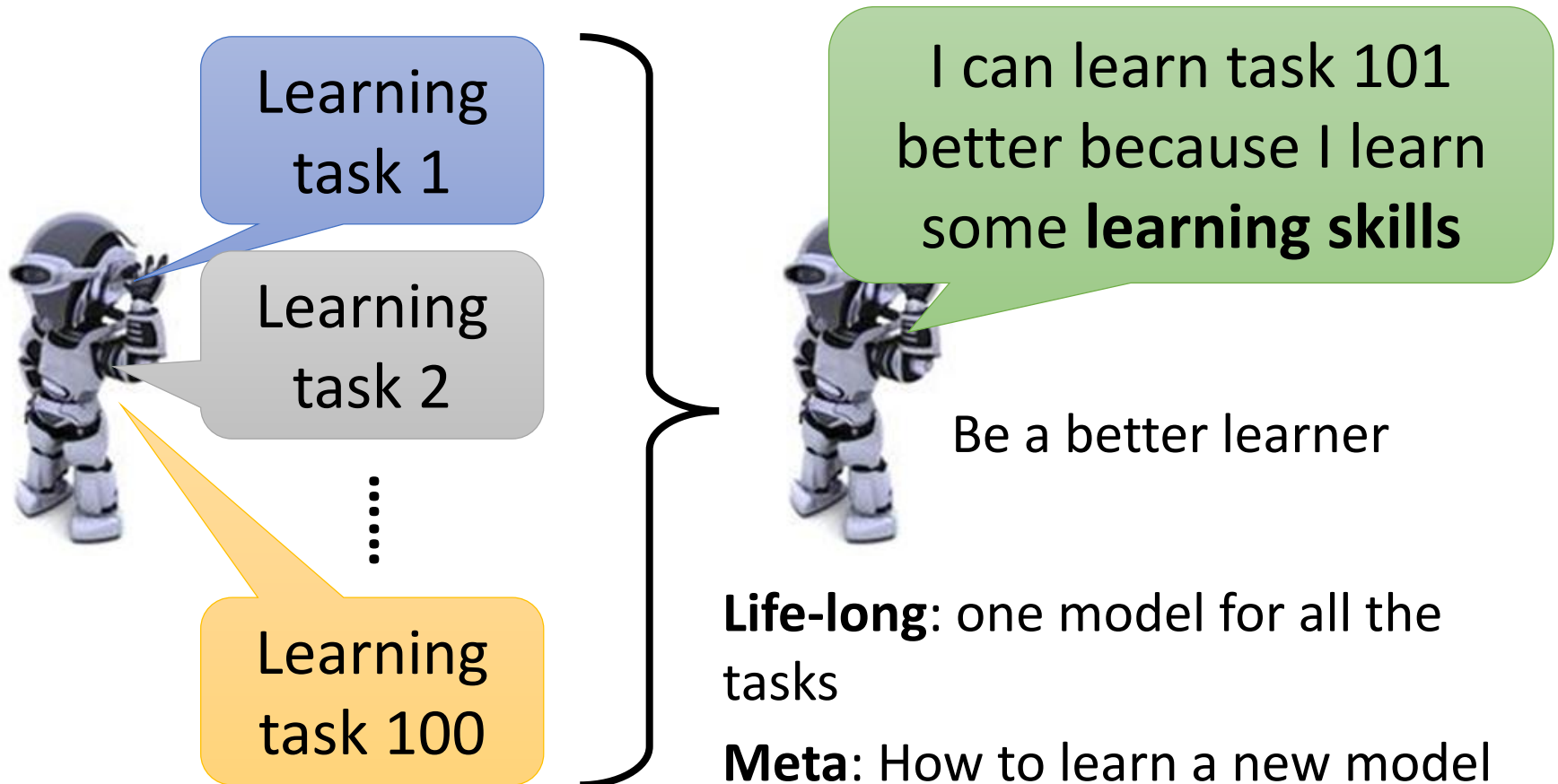
Task 1: speech recognition

Task 2: image recognition

⋮

Task 100: text classification

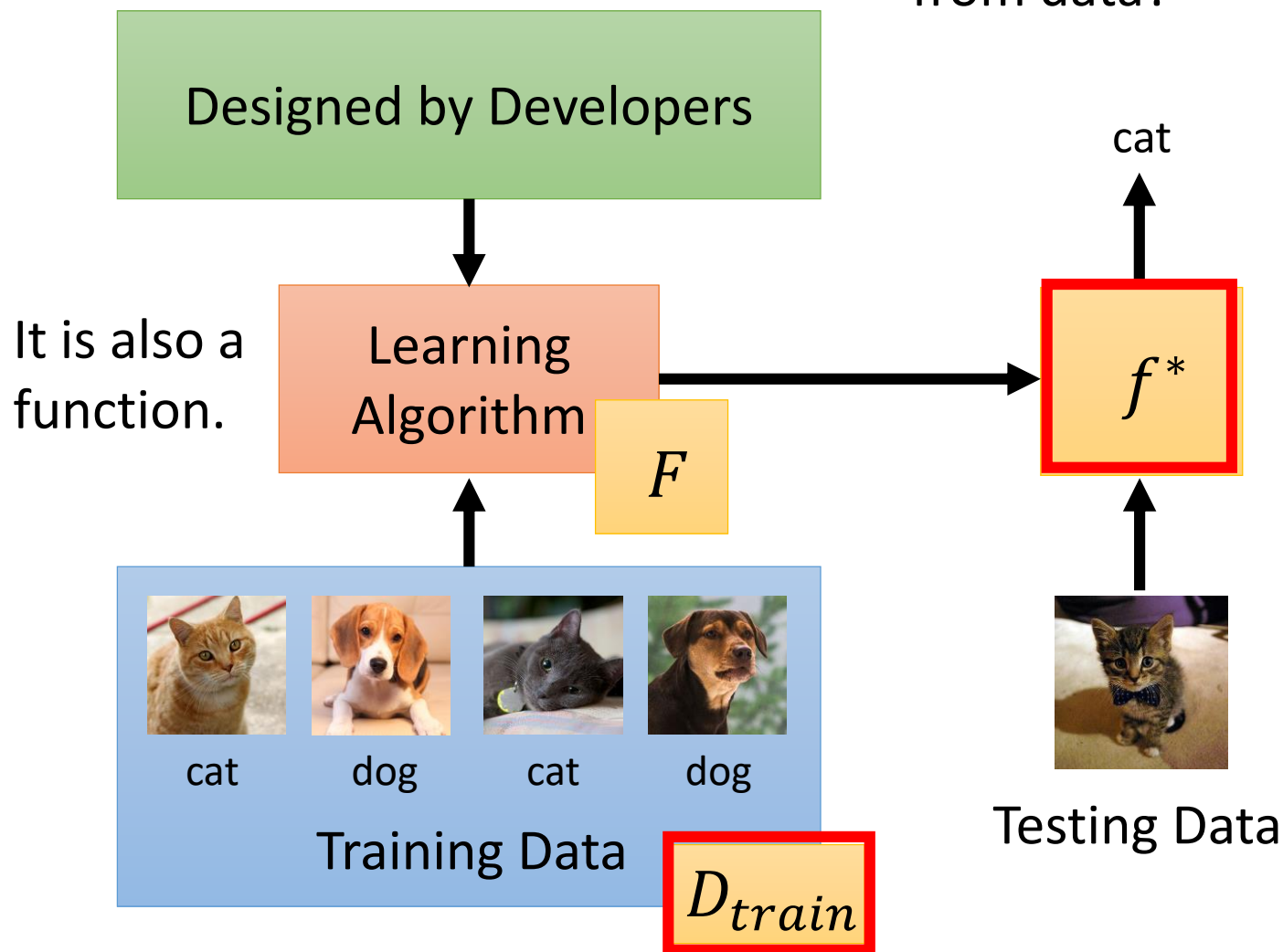
- Meta learning = Learn to learn



Meta Learning

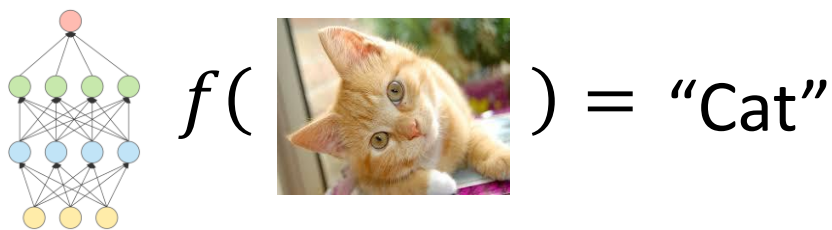
$$f^* = F(D_{train})$$

Can machine find F from data?



Meta Learning

Machine Learning \approx 根據資料找一個函數 f 的能力



Meta Learning

\approx 根據資料找一個找一個函數 f 的函數 F 的能力



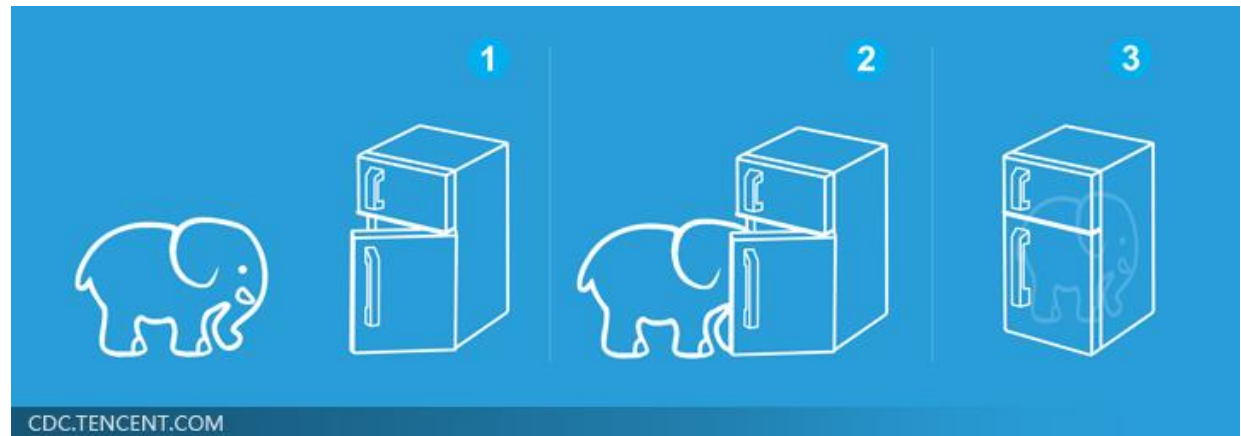
~~Machine~~ Learning is Simple

Meta



Function f \longrightarrow Learning algorithm F

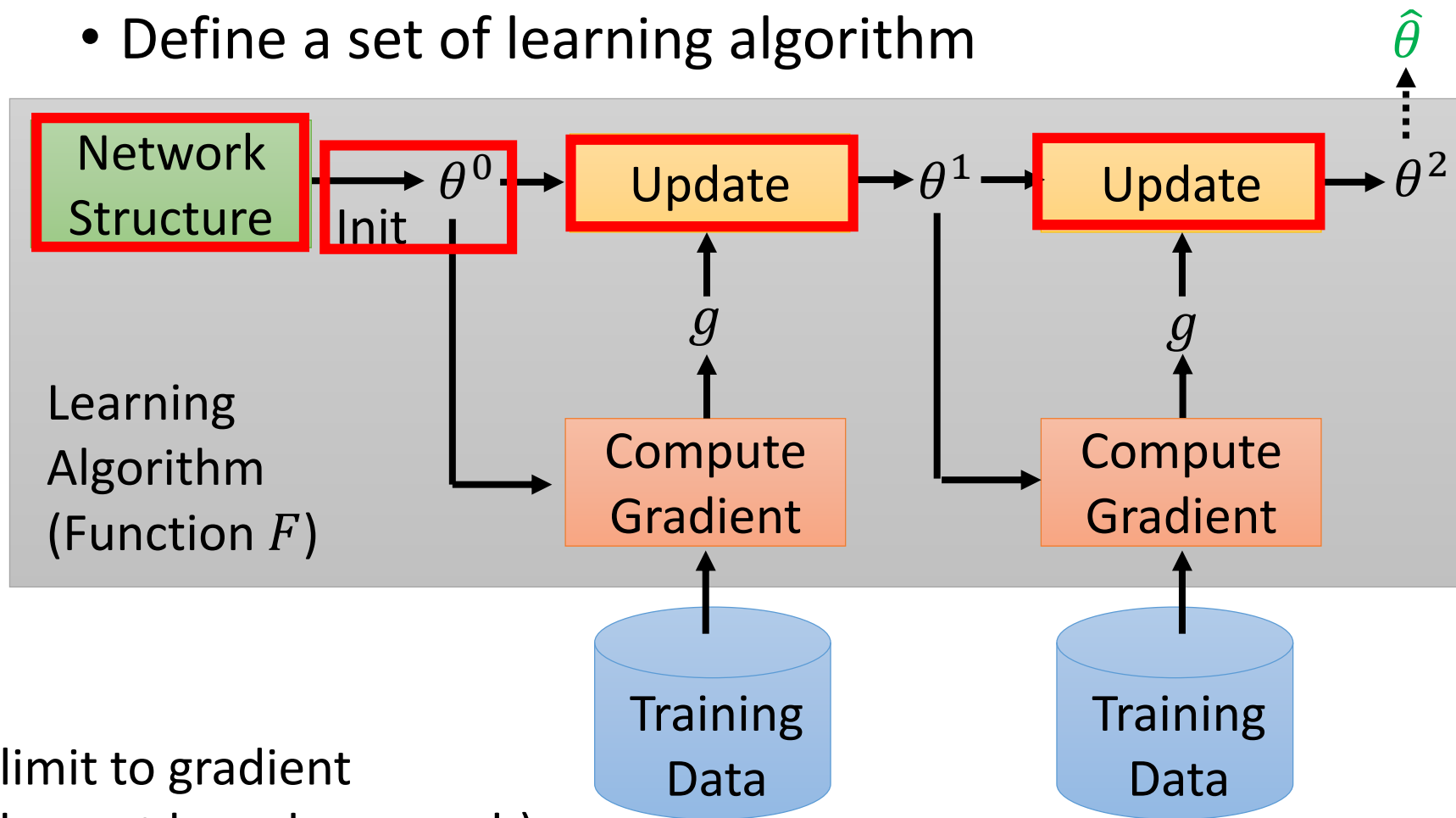
就好像把大象放進冰箱



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



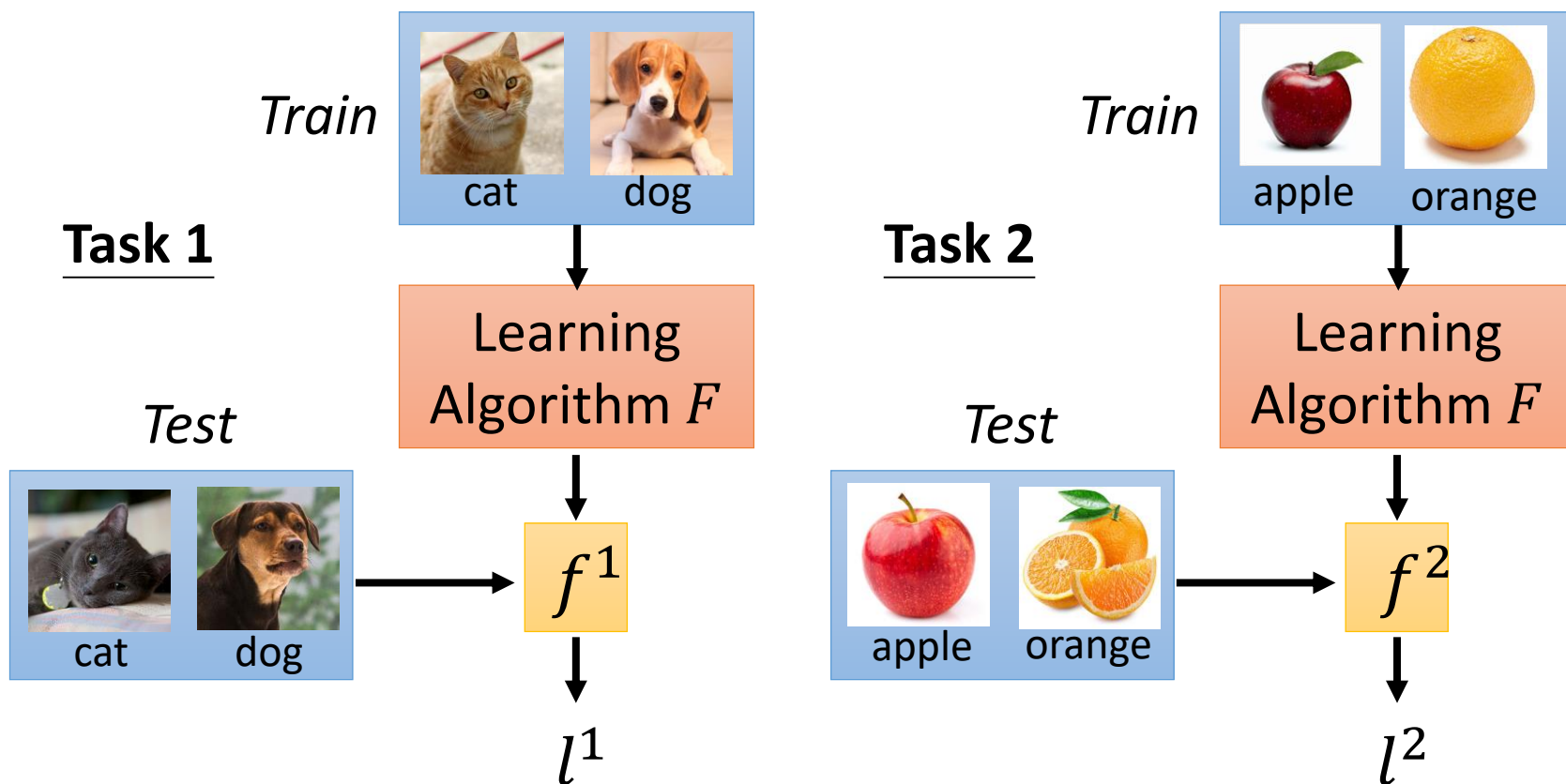
(limit to gradient descent based approach)

Meta Learning

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

N → N tasks
 l^n → Testing loss for task n after training

- Defining the goodness of a function F



Meta Learning

Widely considered in
few-shot learning

Training Tasks

Task 1

Train



cat



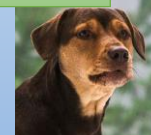
dog

Support set

Test



cat



dog

Query set

Task 2

Train

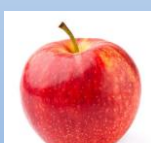


apple



orange

Test



apple



orange

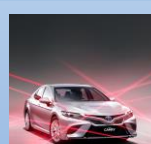
Sometimes you need
validation tasks

Testing Tasks

Train



bike

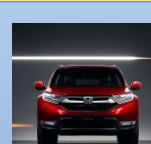


car

Test



bike



car

Machine Learning



cat



dog

Train



cat



dog

Test

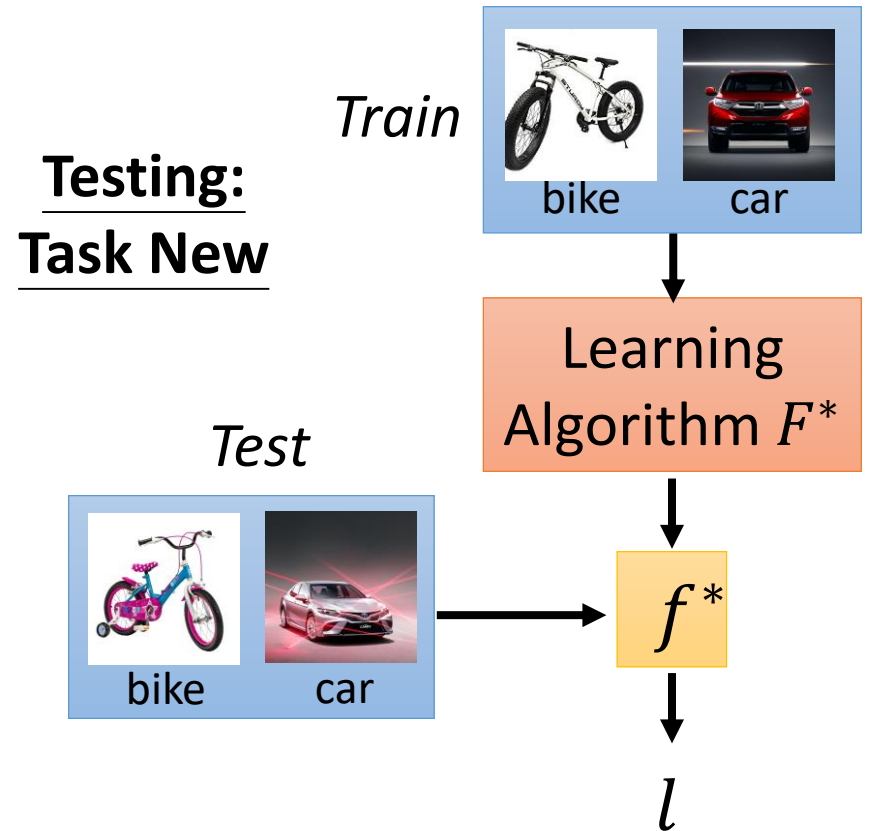
Meta Learning

- Defining the goodness of a function F

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

- Find the best function F^*

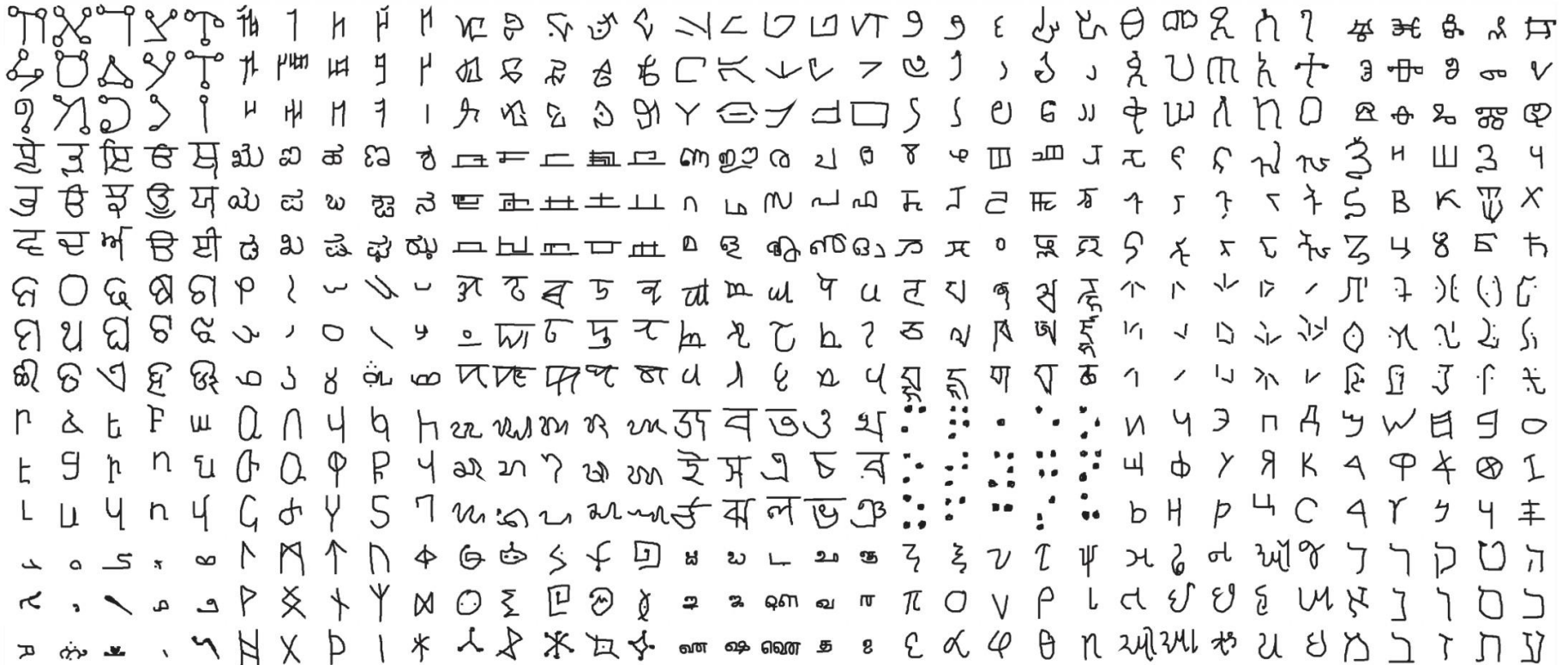
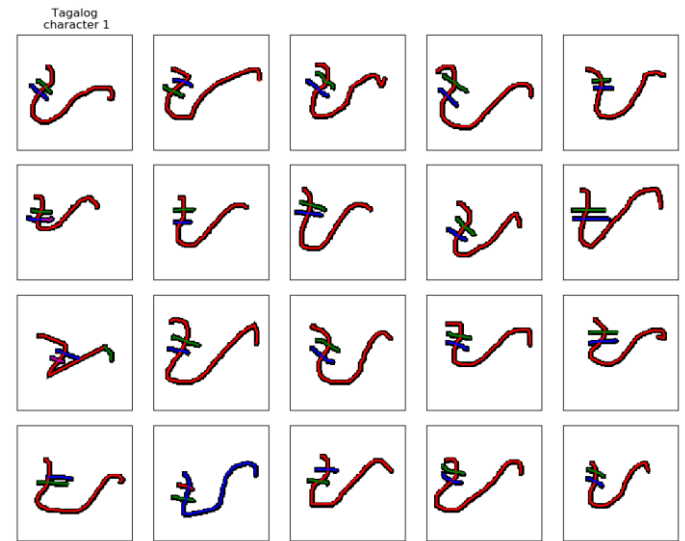
$$F^* = \arg \min_F L(F)$$



Omniglot

<https://github.com/brendenlake/omniglot>

- 1623 characters
- Each has 20 examples



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

𑀧	𑀩	𑀭	𑀮	𑀲
𑀵	𑀶	𑀷	𑀸	𑀹
𑀺	𑀻	𑀼	𑀽	𑀾
𑀿	𑁀	𑁁	𑁂	𑁃

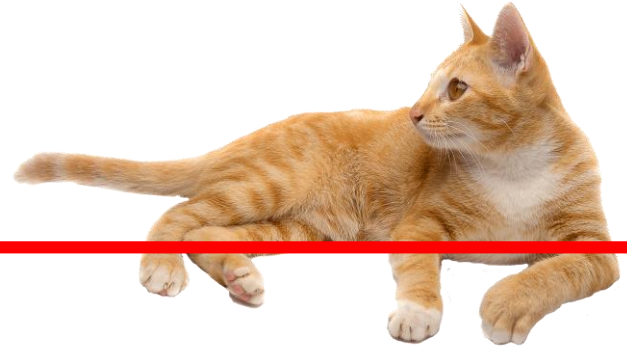


Testing set
(Query set)

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

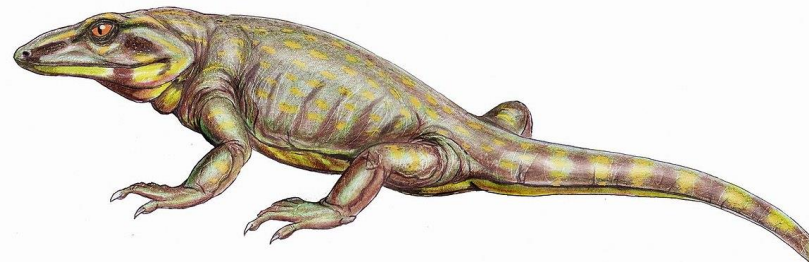


- **MAML**

- 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

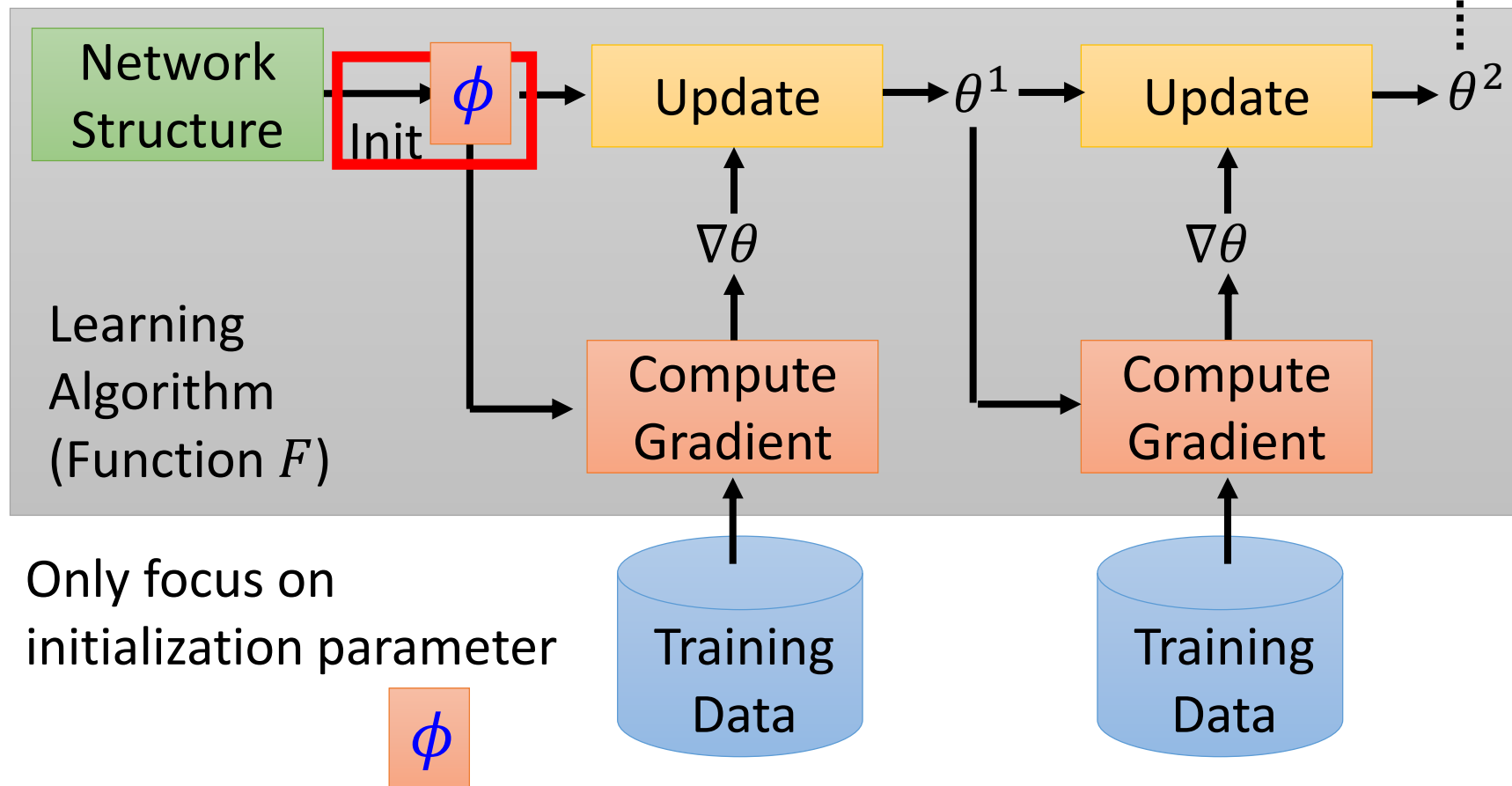
$\hat{\theta}^n$: model learned from task n

Loss Function:

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

$\hat{\theta}^n$ depends on ϕ

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



MAML

Loss Function:

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

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

$\hat{\theta}^n$ depends on ϕ

$l^n(\hat{\theta}^n)$: 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:

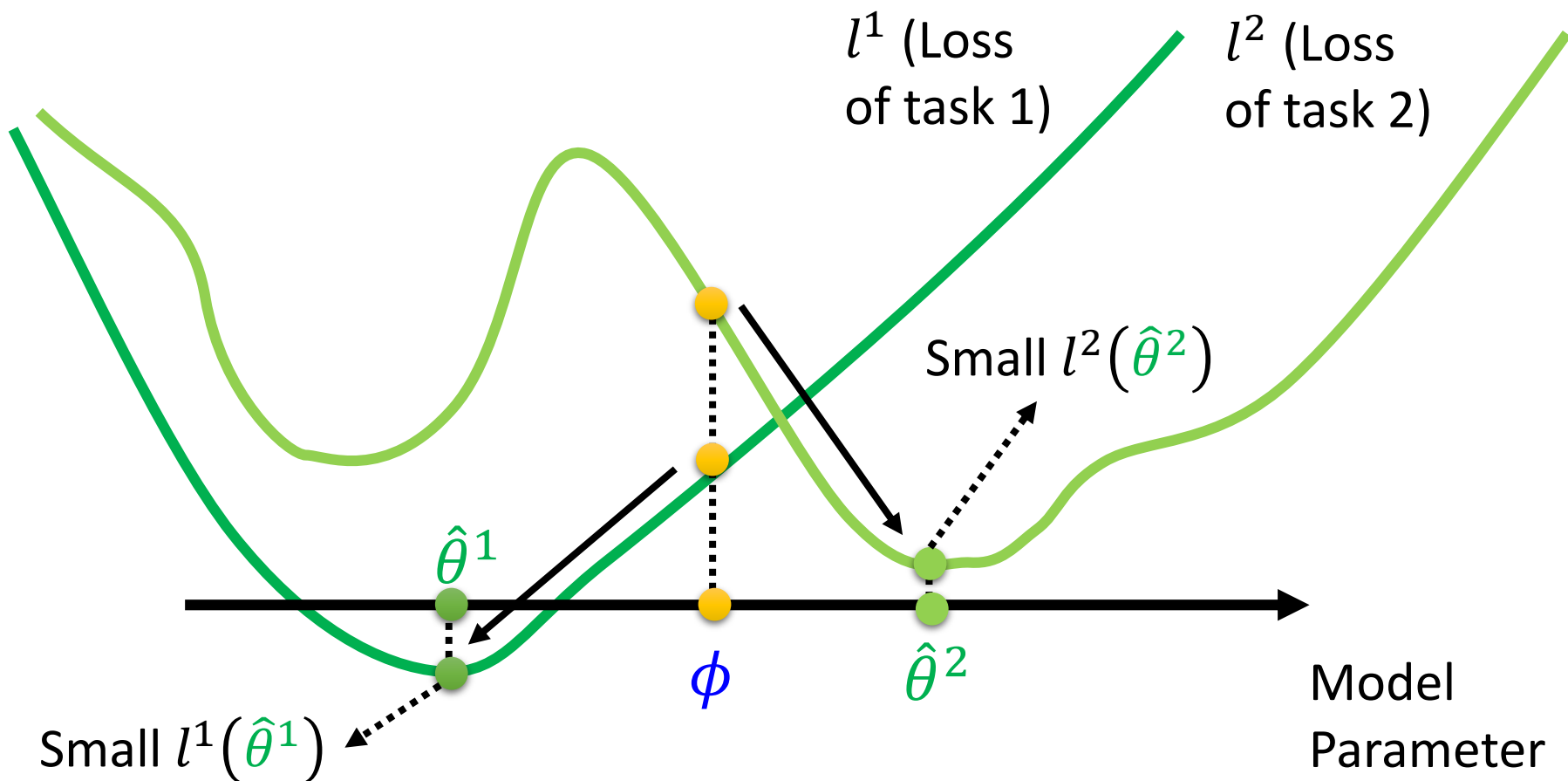
$$L(\phi) = \sum_{n=1}^N l^n(\phi)$$

MAML

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

我們不在意 ϕ 在 training task 上表現如何

我們在意用 ϕ 訓練出來的 $\hat{\theta}^n$ 表現如何

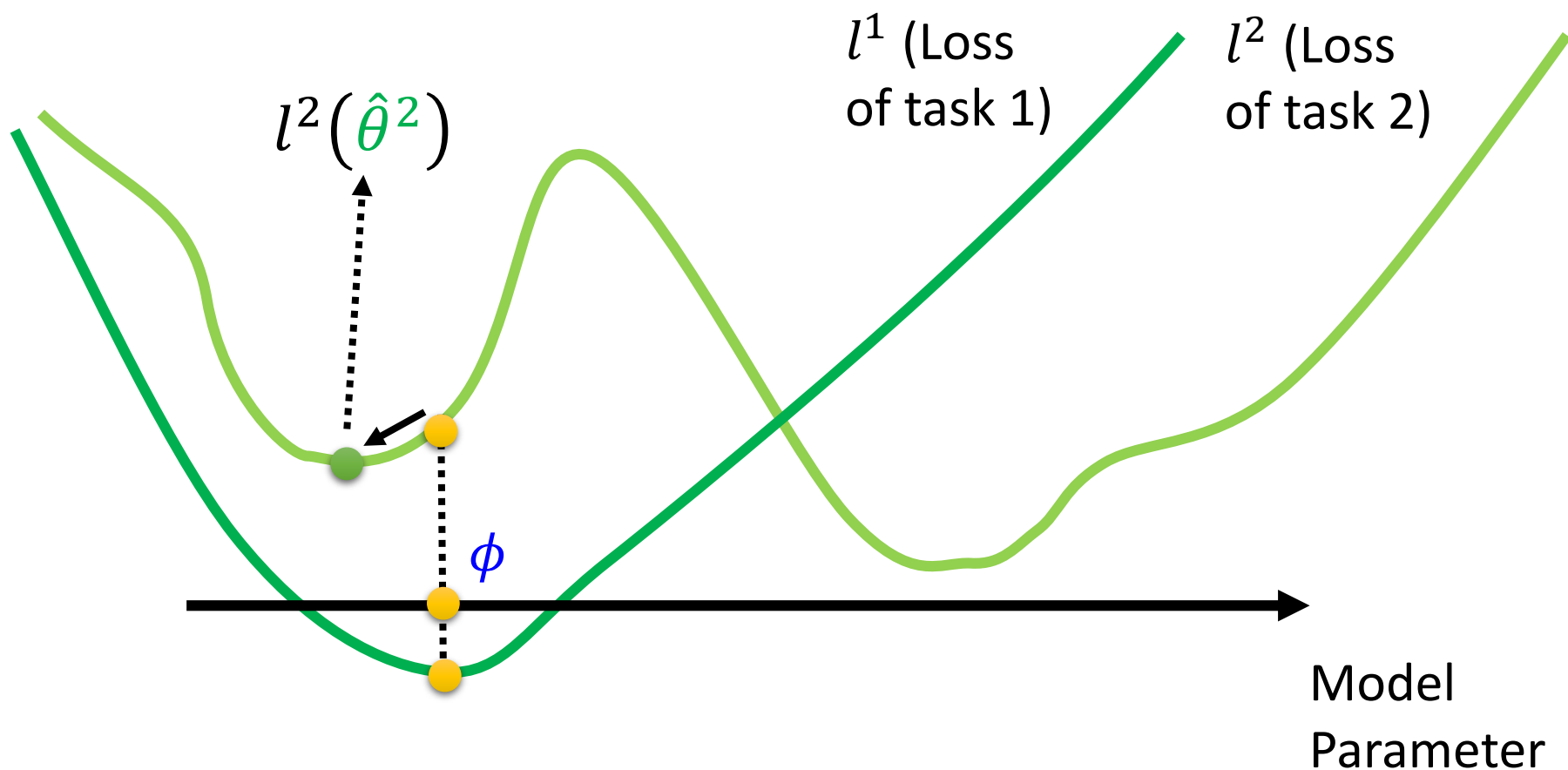


Model Pre-training

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

找尋在所有 task 都最好的 ϕ

並不保證拿 ϕ 去訓練以後會
得到好的 $\hat{\theta}^n$



MAML

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

Loss Function:

$\hat{\theta}^n$ depends on ϕ

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

$l^n(\hat{\theta}^n)$: 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 ϕ achieving good performance **after training**

潛力

Model Pre-training

Widely used in transfer learning

Loss Function:

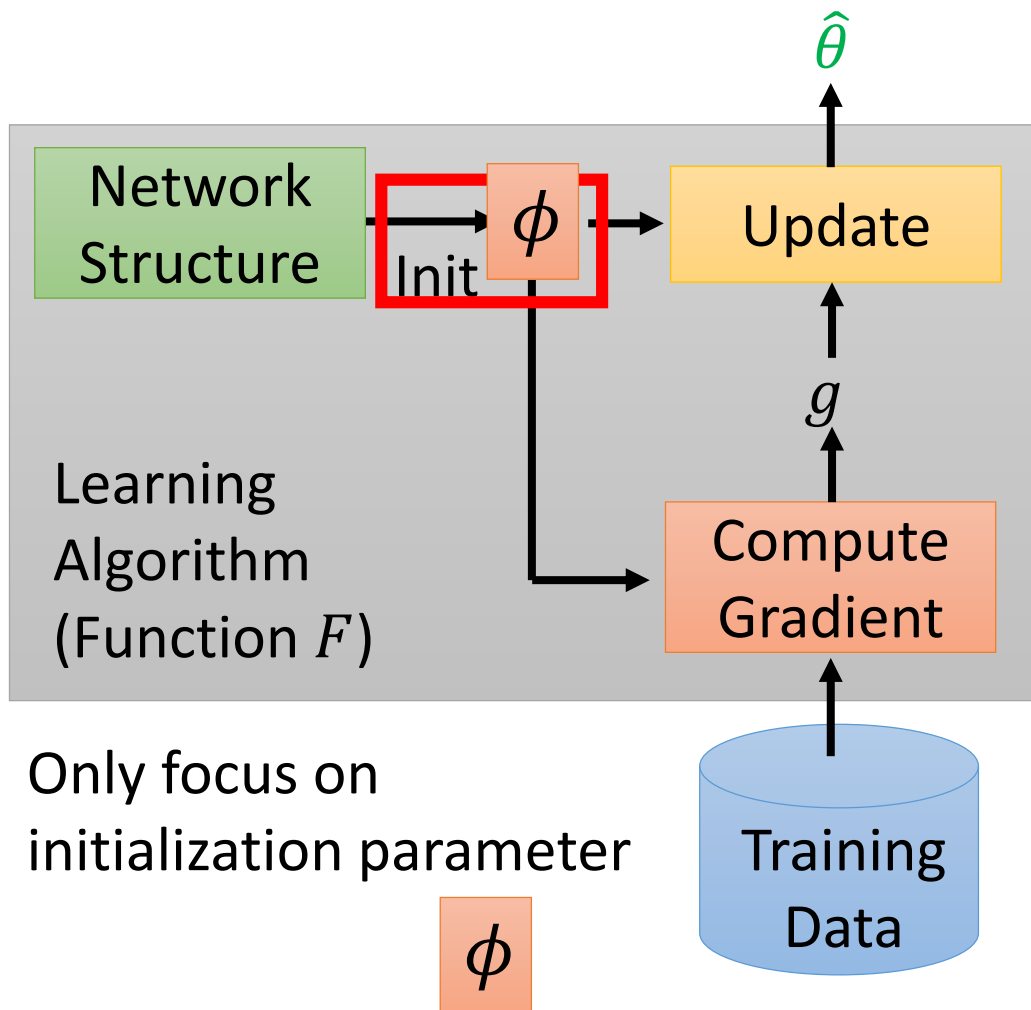
$$L(\phi) = \sum_{n=1}^N l^n(\phi)$$

Find ϕ achieving good performance

現在表現如何

MAML

- Fast ... Fast ... Fast ...
- Good to truly train a model with one step. 😊
- When using the algorithm, still update many times.
- Few-shot learning has limited data.



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

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

Considering one-step training:

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

Toy Example

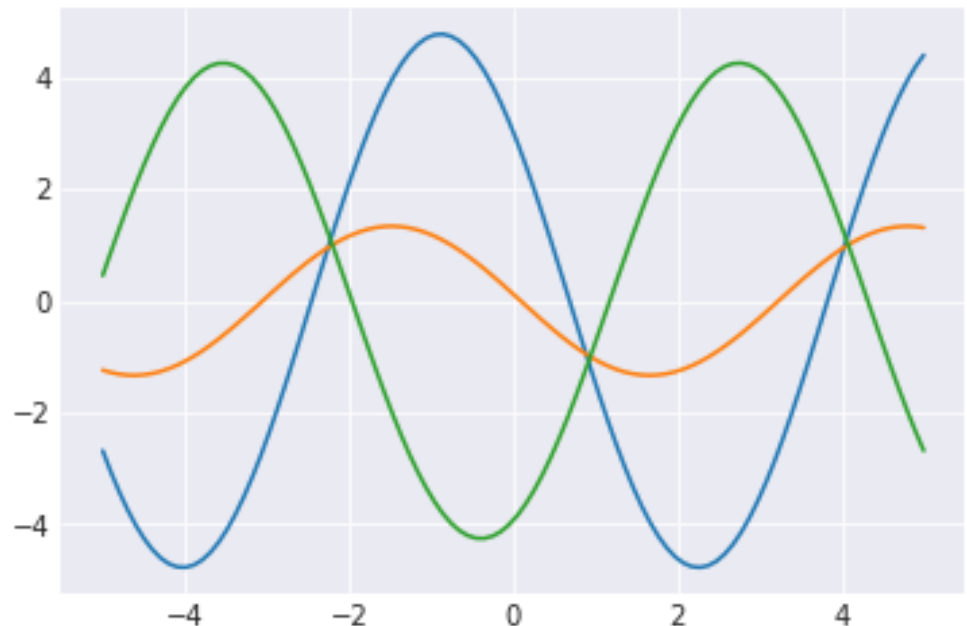
Source of images

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

Each task:

- Given a target sine function $y = a \sin(x + b)$
- Sample K points from the target function
- Use the samples to estimate the target function

Sample a and b to form a task



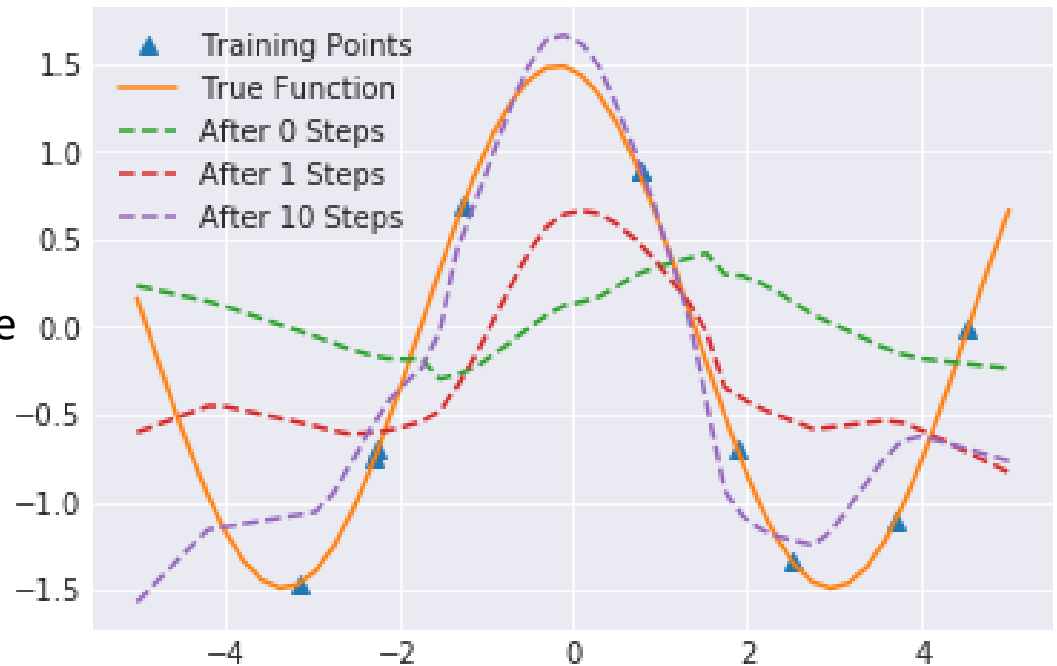
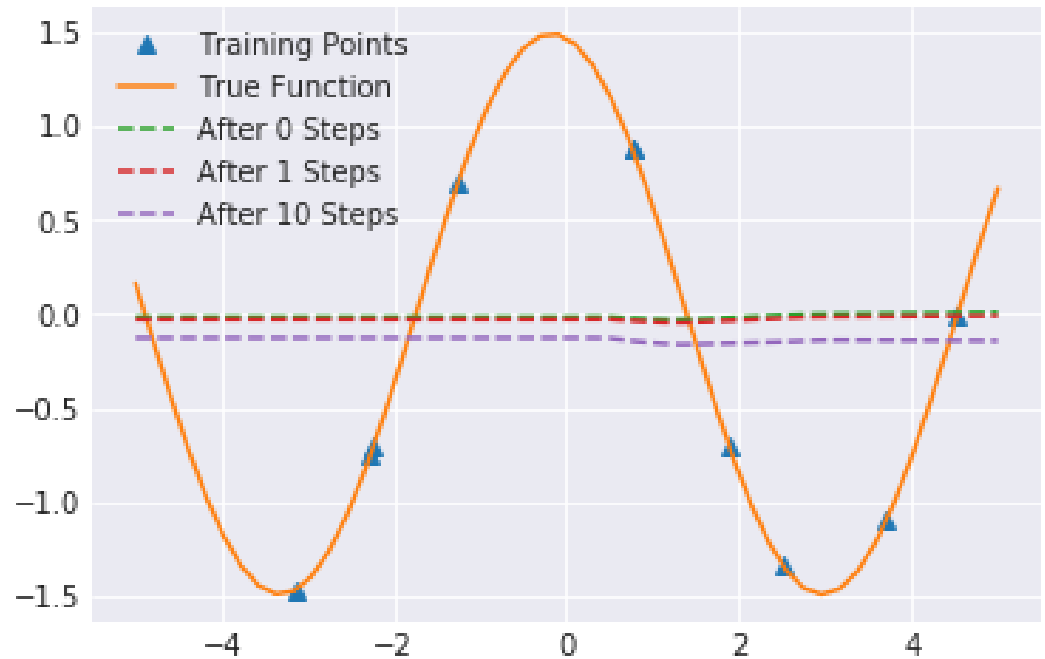
Toy Example

Model Pre-training

Source of images

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

MAML



Omniglot & Mini-ImageNet

	5-way Accuracy		20-way Accuracy	
	1-shot	5-shot	1-shot	5-shot
Omniglot (Lake et al., 2011)				
MANN, no conv (Santoro et al., 2016)	82.8%	94.9%	–	–
MAML, no conv (ours)	89.7 ± 1.1%	97.5 ± 0.6%	–	–
Siamese nets (Koch, 2015)	97.3%	98.4%	88.2%	97.0%
matching nets (Vinyals et al., 2016)	98.1%	98.9%	93.8%	98.5%
neural statistician (Edwards & Storkey, 2017)	98.1%	99.5%	93.2%	98.1%
memory mod. (Kaiser et al., 2017)	98.4%	99.6%	95.0%	98.6%
MAML (ours)	98.7 ± 0.4%	99.9 ± 0.1%	95.8 ± 0.3%	98.9 ± 0.2%

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

Warning of Math

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

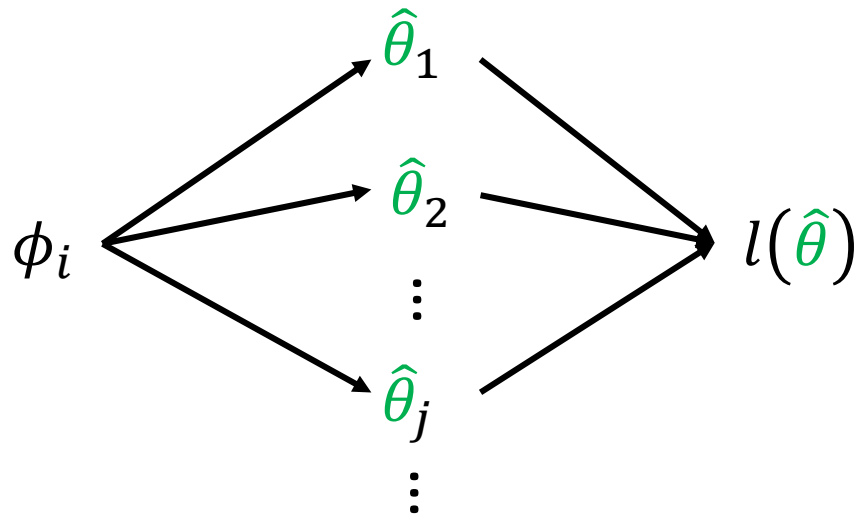
$$L(\phi) = \sum_{n=1}^N l^n(\hat{\theta}^n)$$

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

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

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

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



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

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

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

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

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

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

$i \neq j$:

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

$i = j$:

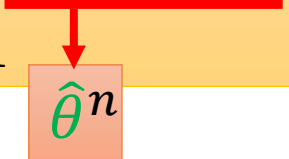
$$\frac{\partial \hat{\theta}_j}{\partial \phi_i} = 1 - \varepsilon \frac{\partial l(\phi)}{\partial \phi_i \partial \phi_j} \approx 1$$

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

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

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

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

$$\nabla_{\phi} L(\phi) = \nabla_{\phi} \sum_{n=1}^N l^n(\hat{\theta}^n) = \sum_{n=1}^N \nabla_{\phi} l^n(\hat{\theta}^n)$$
A red arrow points from the term $\nabla_{\phi} l^n(\hat{\theta}^n)$ in the sum to a red box containing $\hat{\theta}^n$.

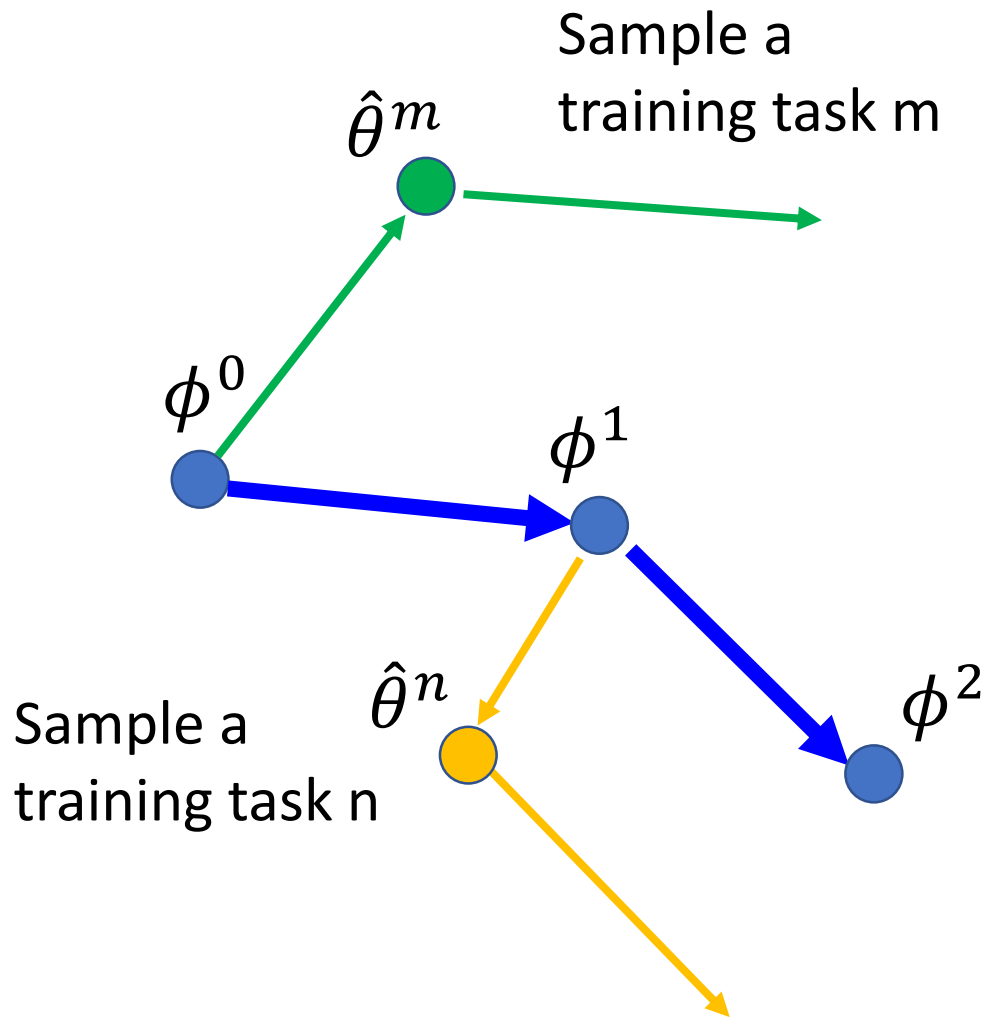
$$\frac{\partial l(\hat{\theta})}{\partial \phi_i} = \sum_j \frac{\partial l(\hat{\theta})}{\partial \hat{\theta}_j} \frac{\partial \hat{\theta}_j}{\partial \phi_i} \approx \frac{\partial l(\hat{\theta})}{\partial \hat{\theta}_i}$$

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.

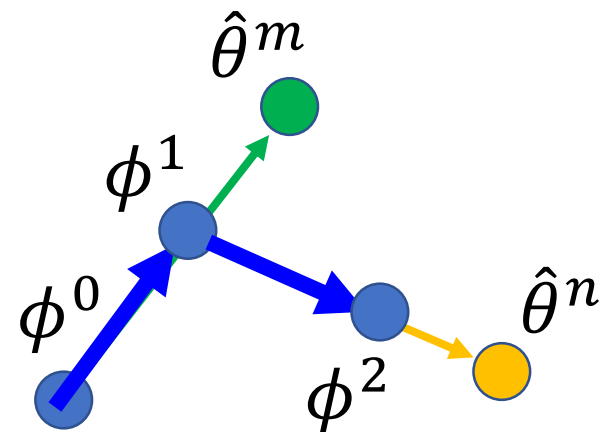
$$\nabla_{\phi} l(\hat{\theta}) = \begin{bmatrix} \partial l(\hat{\theta}) / \partial \phi_1 \\ \partial l(\hat{\theta}) / \partial \phi_2 \\ \vdots \\ \partial l(\hat{\theta}) / \partial \phi_i \\ \vdots \end{bmatrix} = \begin{bmatrix} \partial l(\hat{\theta}) / \partial \hat{\theta}_1 \\ \partial l(\hat{\theta}) / \partial \hat{\theta}_2 \\ \vdots \\ \partial l(\hat{\theta}) / \partial \hat{\theta}_i \\ \vdots \end{bmatrix} = \nabla_{\hat{\theta}} l(\hat{\theta})$$

End of Warning

MAML – Real Implementation

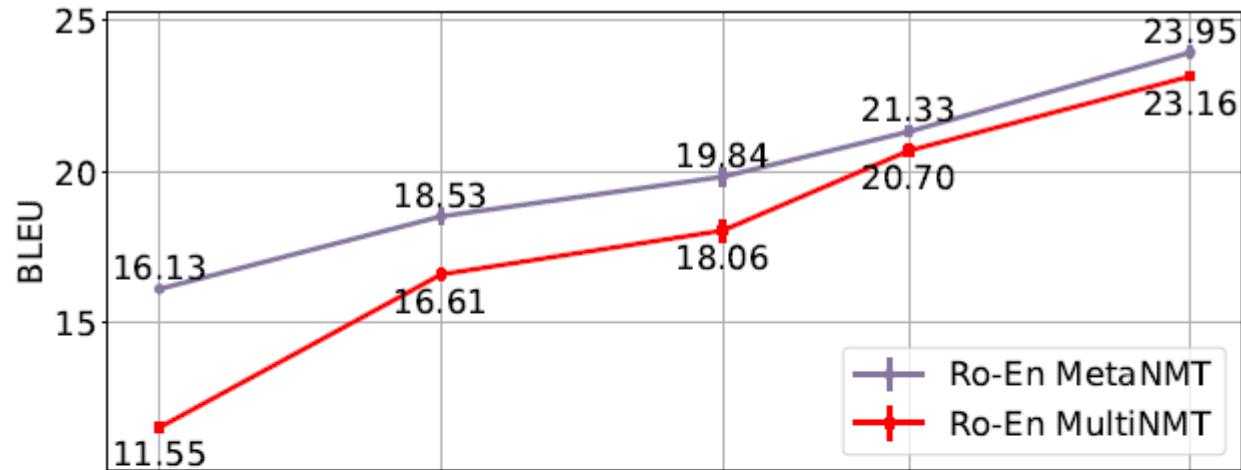


Model Pre-training

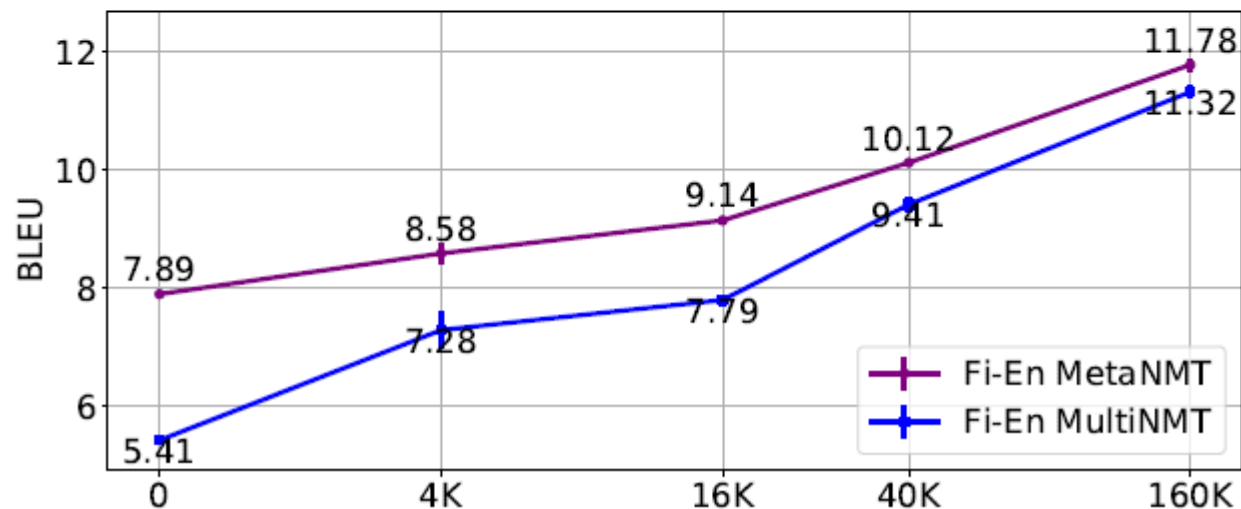


Translation

18 training tasks: 18 different languages translating to English
2 validation tasks: 2 different languages translating to English



Ro = Romanian



Fi = Finnish

<https://arxiv.org/abs/1808.08437>

Techniques Today



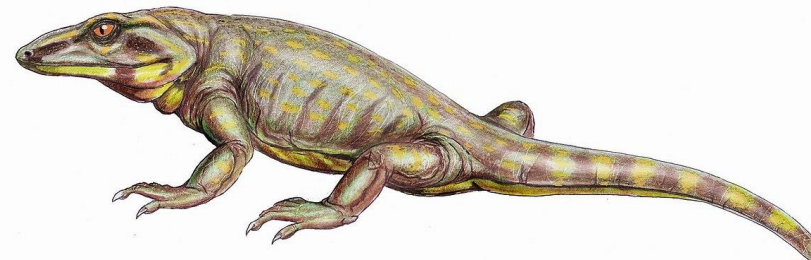
- **MAML**

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

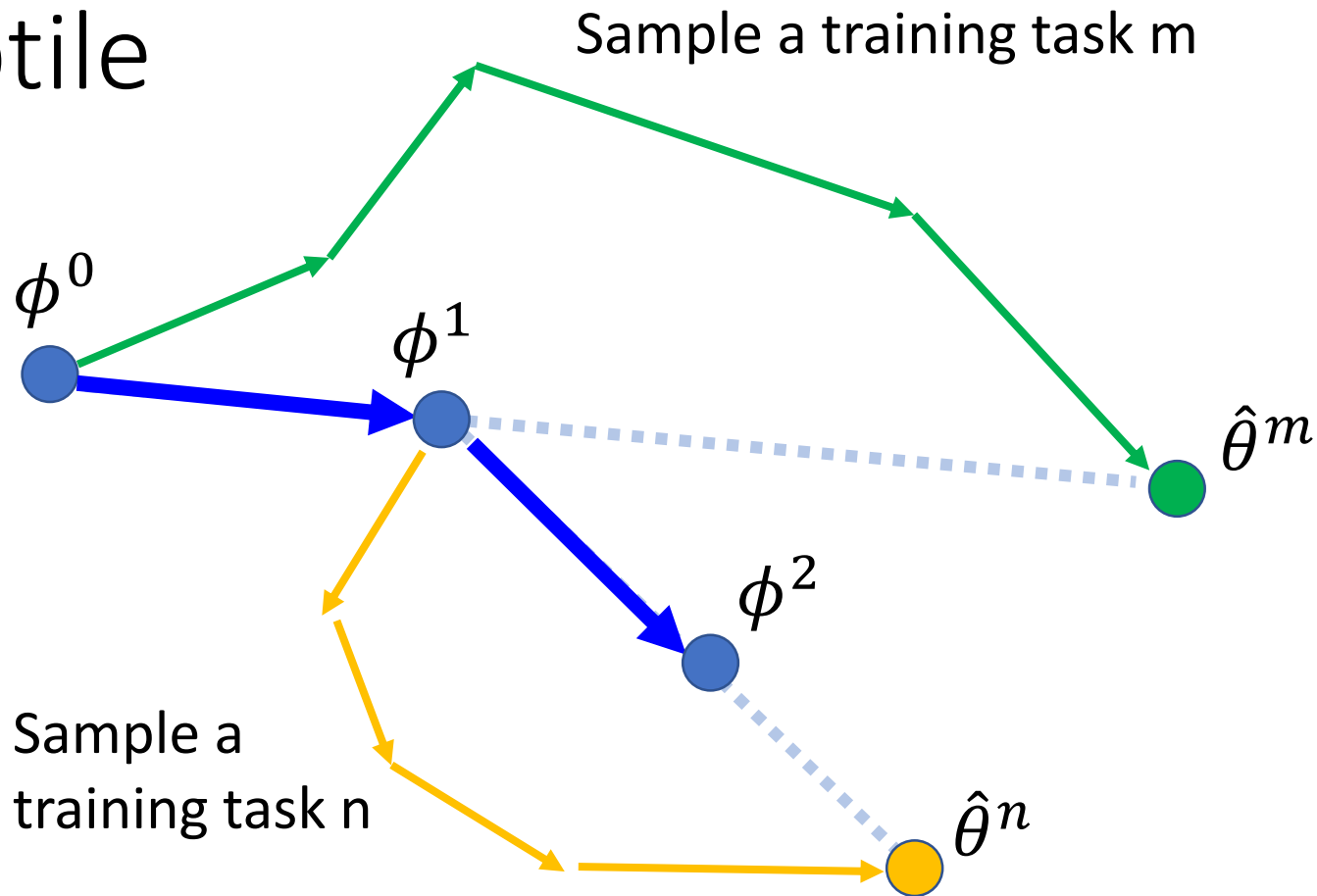
- **Reptile**

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

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



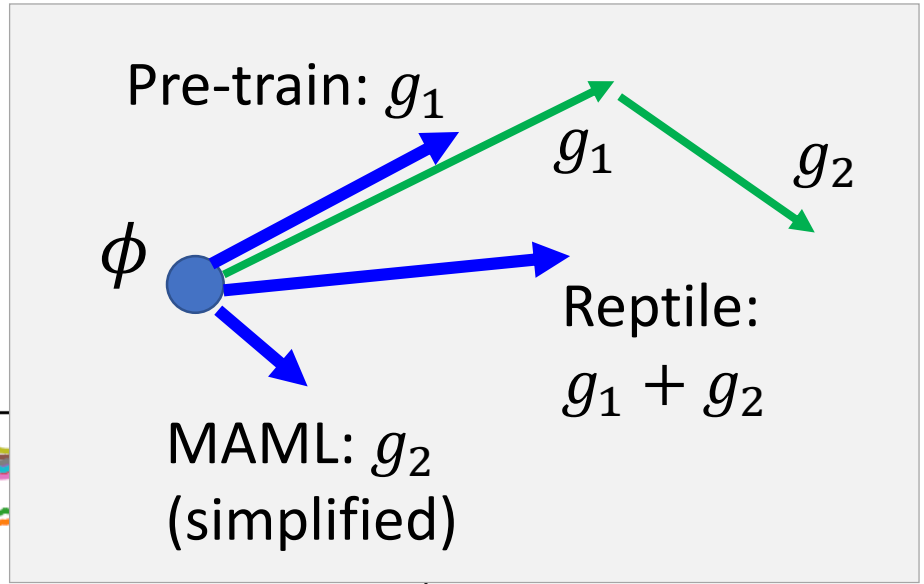
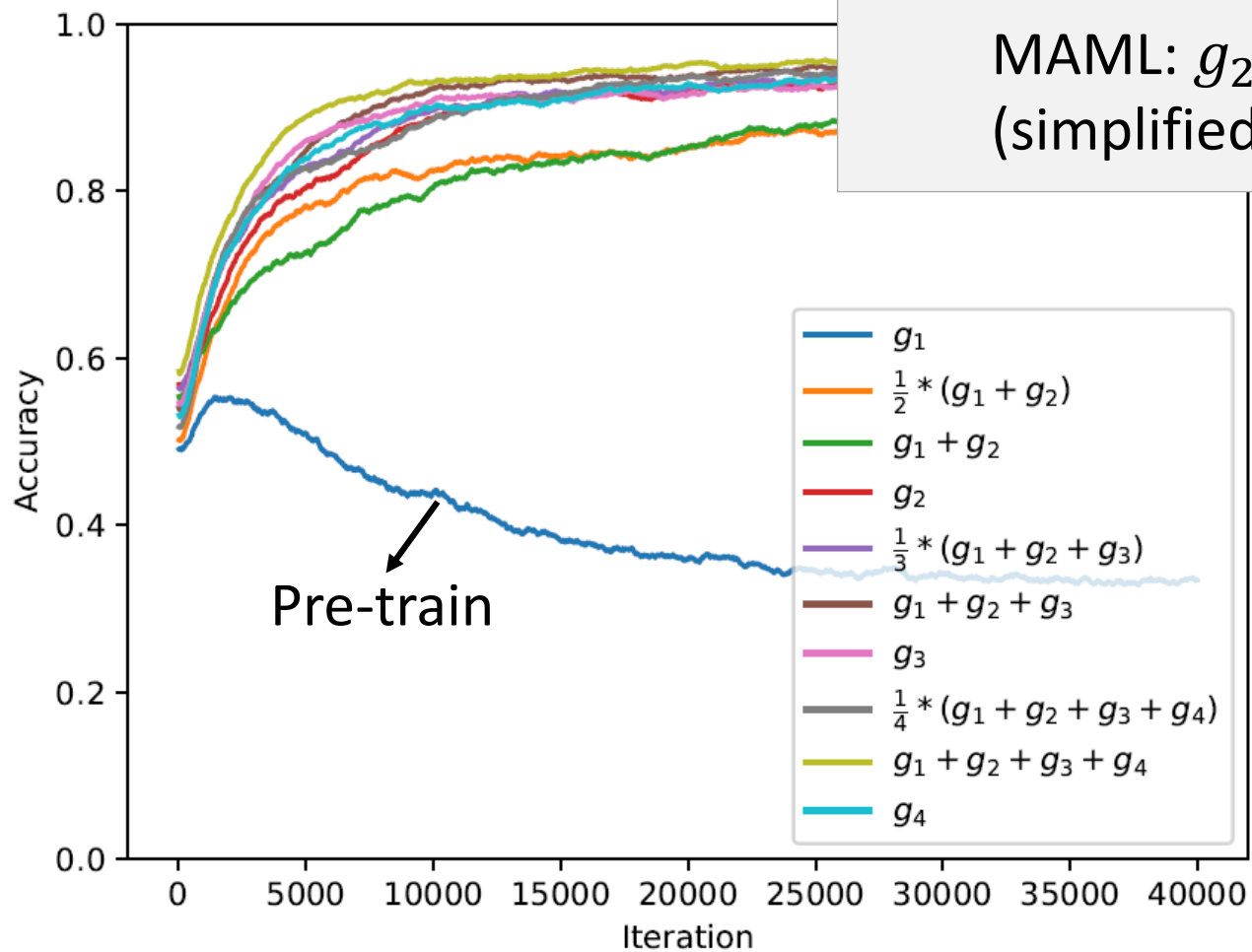
Reptile



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:

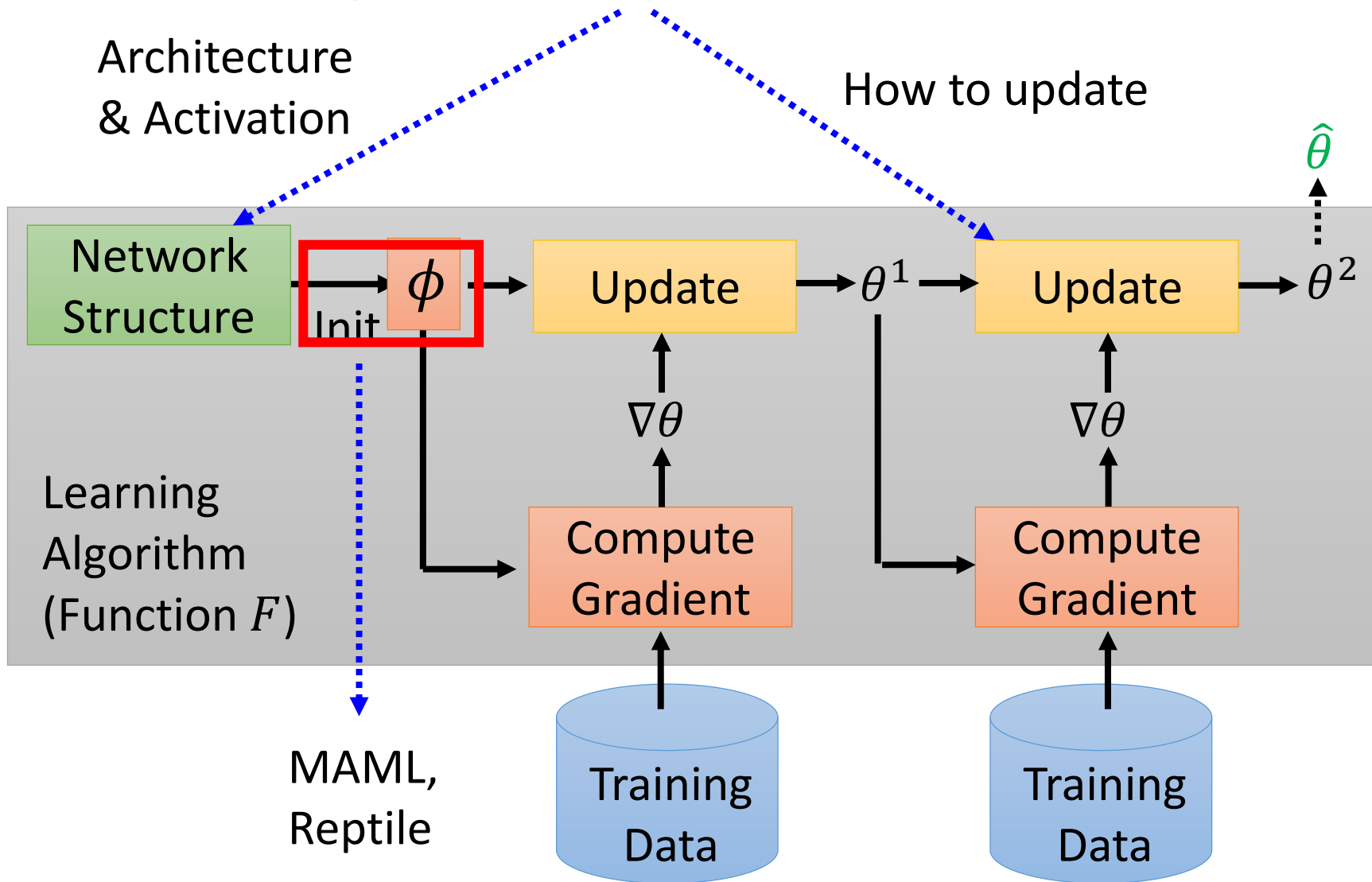
(this sentence is removed in the updated version)

Reptile



More ... Video: <https://www.youtube.com/watch?v=c10nxBcSH14>

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



Turtles all the way down ?



- We learn the initialization parameter ϕ by gradient descent
 - What is the initialization parameter ϕ^0 for initialization parameter ϕ ?
- Learn
- Learn to learn
- Learn to learn to learn

Crazy Idea?

下回分解 😊

- How about learning algorithm beyond gradient descent?

