

---

# MAML is a Noisy Contrastive Learner

---

NewInML Workshop@NeuroIPS'21

Chia-Hsiang Kao<sup>1</sup>, Wei-Chen Chiu<sup>1</sup>, Pin-Yu Chen<sup>2</sup>

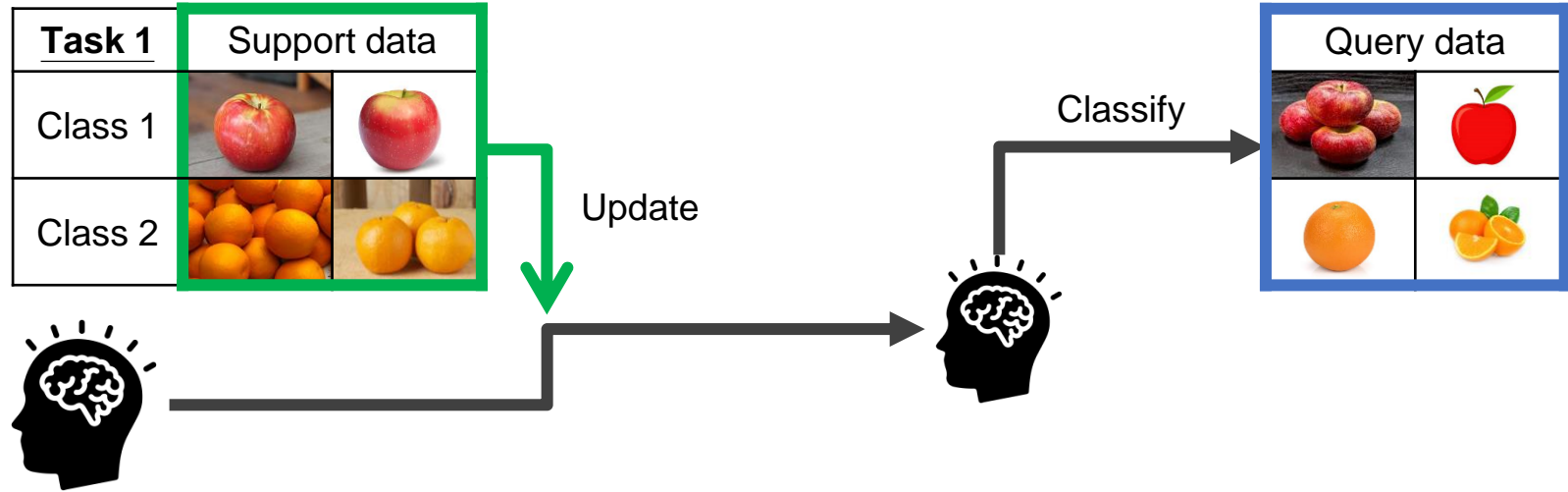
<sup>1</sup>National Yang Ming Chiao Tung University

<sup>2</sup>MIT-IBM Watson AI Lab



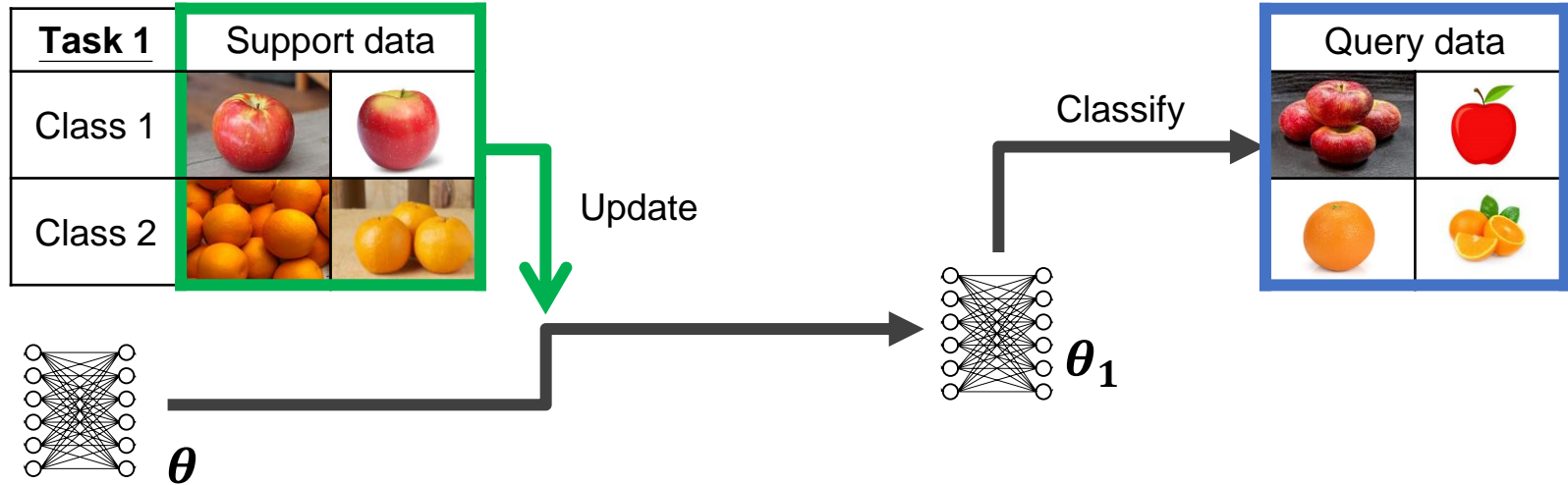
# Introduction

Humans learn to classify even with limited experience.



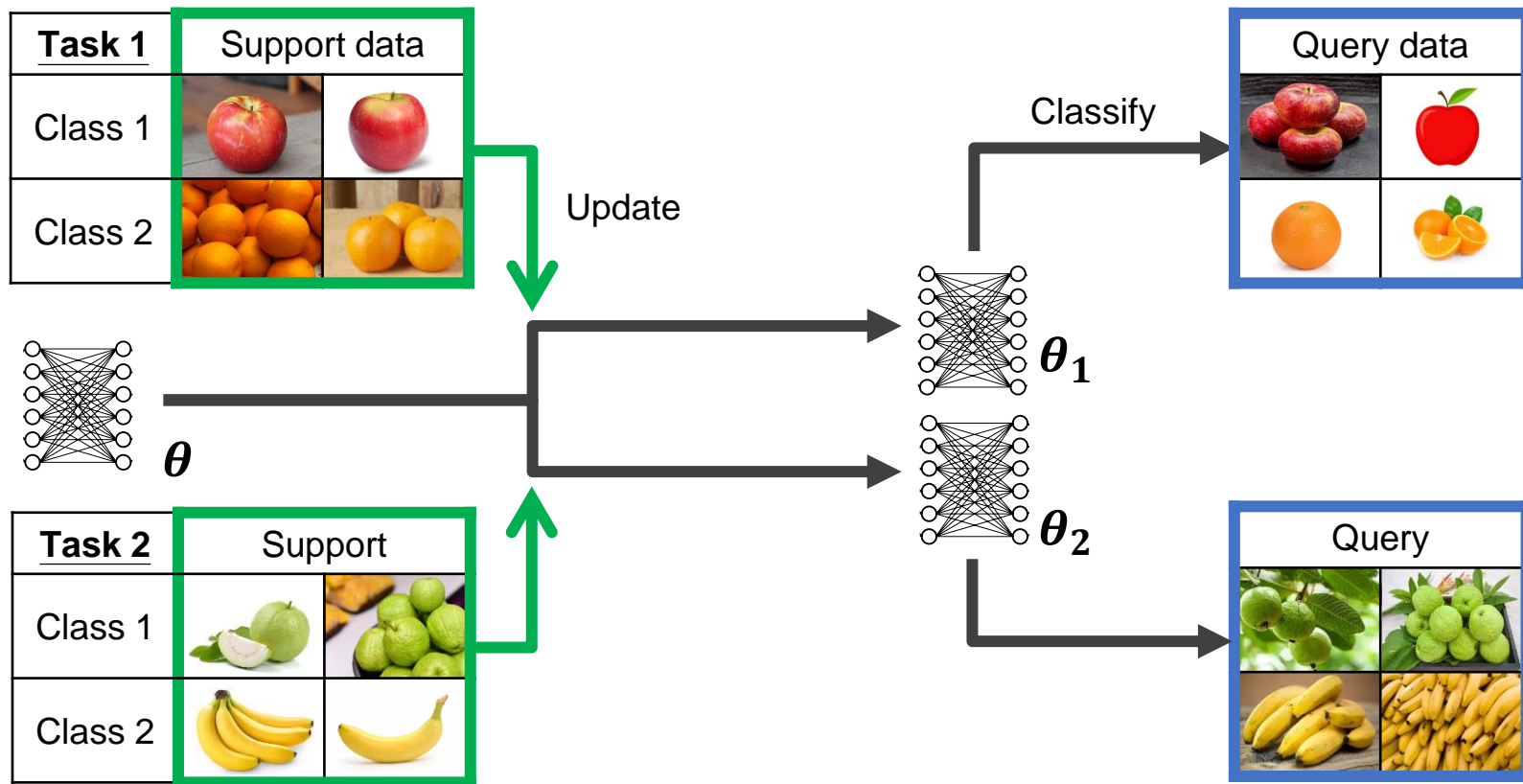
# Introduction

Humans learn to classify even with limited experience.



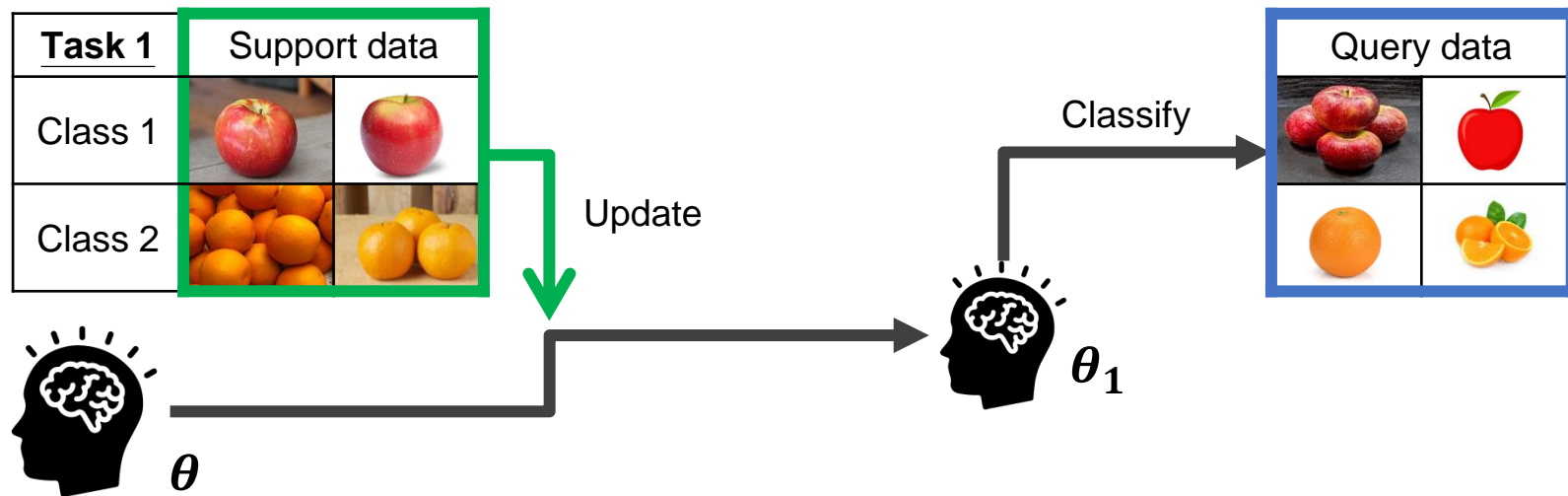
# Introduction

Humans learn to classify even with limited experience.



# Introduction

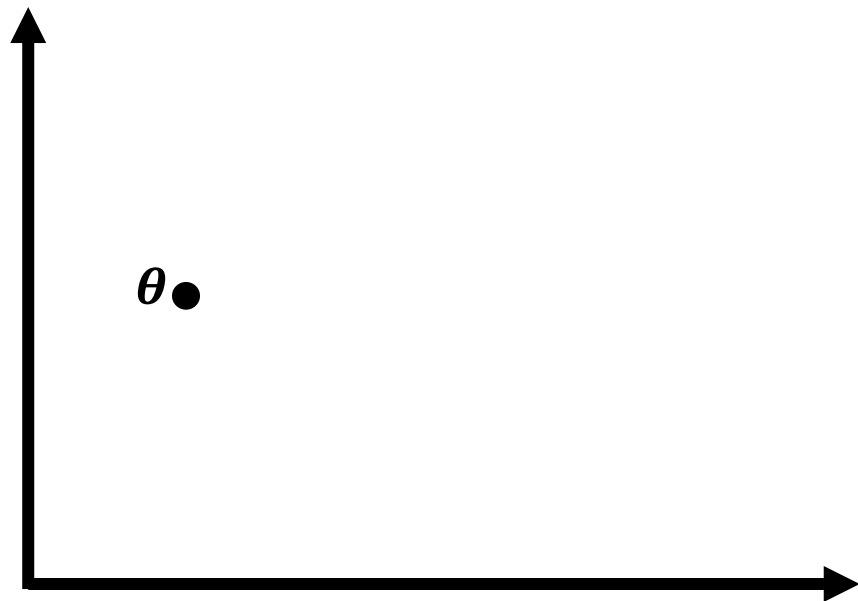
Humans learn to classify even with limited experience.











- MAML is a gradient-based meta-learning algorithm that finds a good  $\theta$ .









# Introduction

MAML makes model learn to classify after seeing little data.



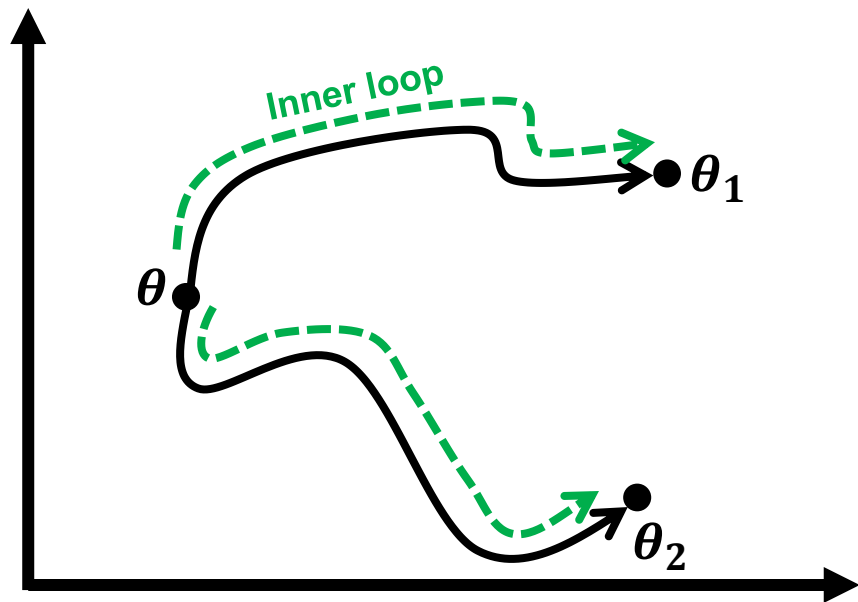
<u>Task 1</u>	Support $S_1$		Query $Q_1$	
Class 1				
Class 2				




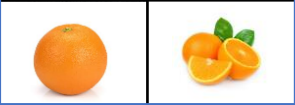












  

<u>Task 2</u>	Support $S_2$		Query $Q_2$	
Class 1				
Class 2				

# Introduction

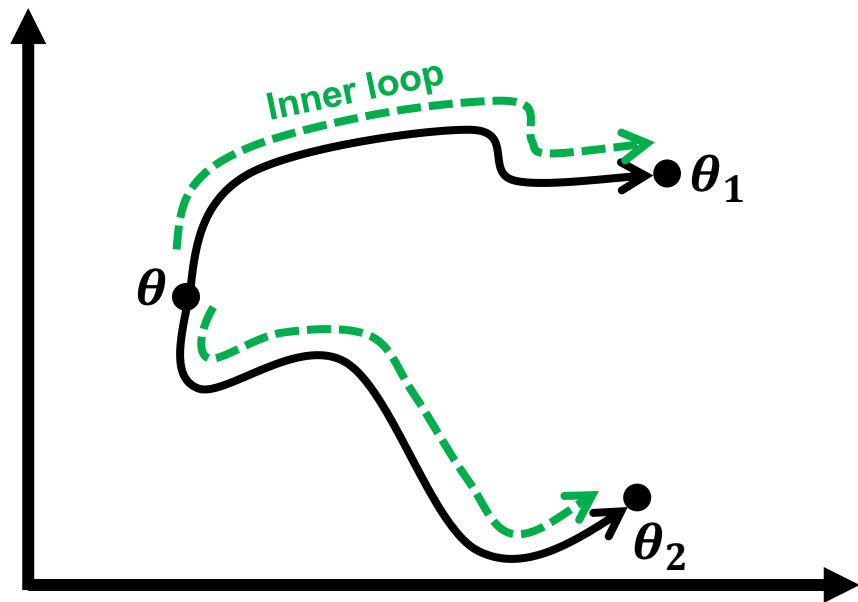
MAML makes model learn to classify after seeing little data.




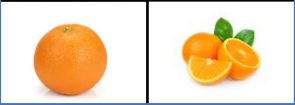



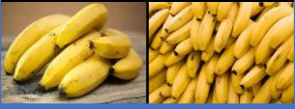


<u>Task 1</u>	Support $S_1$		Query $Q_1$	
Class 1				
Class 2				
<u>Task 2</u>	Support $S_2$		Query $Q_2$	
Class 1				
Class 2				

# Introduction

MAML makes model learn to classify after seeing little data.



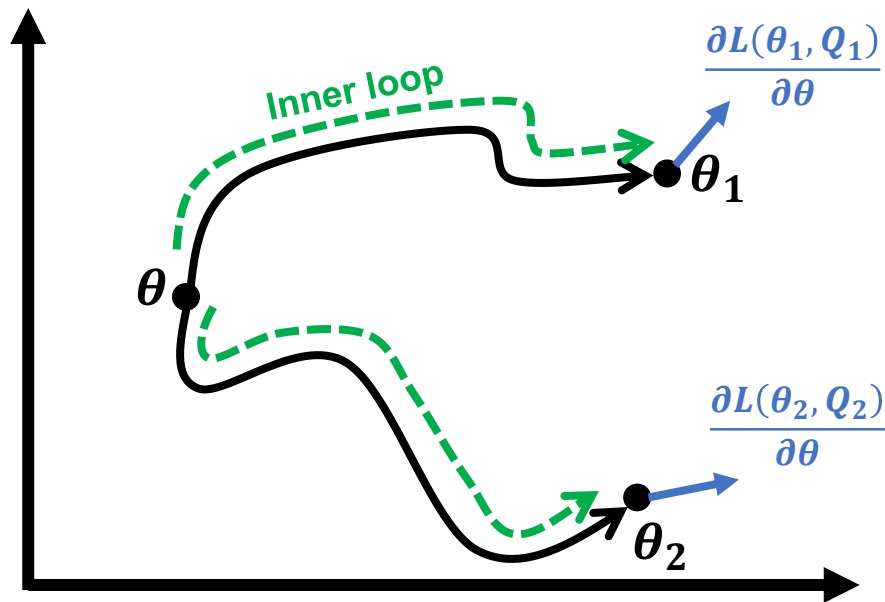
<u>Task 1</u>	Support $S_1$	Query $Q_1$
Class 1		
Class 2		
<u>Task 2</u>	Support $S_2$	Query $Q_2$
Class 1		
Class 2		









- Goal: Minimize  $L(\theta_1, Q_1)$  and  $L(\theta_2, Q_2)$  by finding best  $\theta$ .











# Introduction

MAML makes model learn to classify after seeing little data.



Task 1	Support $S_1$		Query $Q_1$	
Class 1				
Class 2				

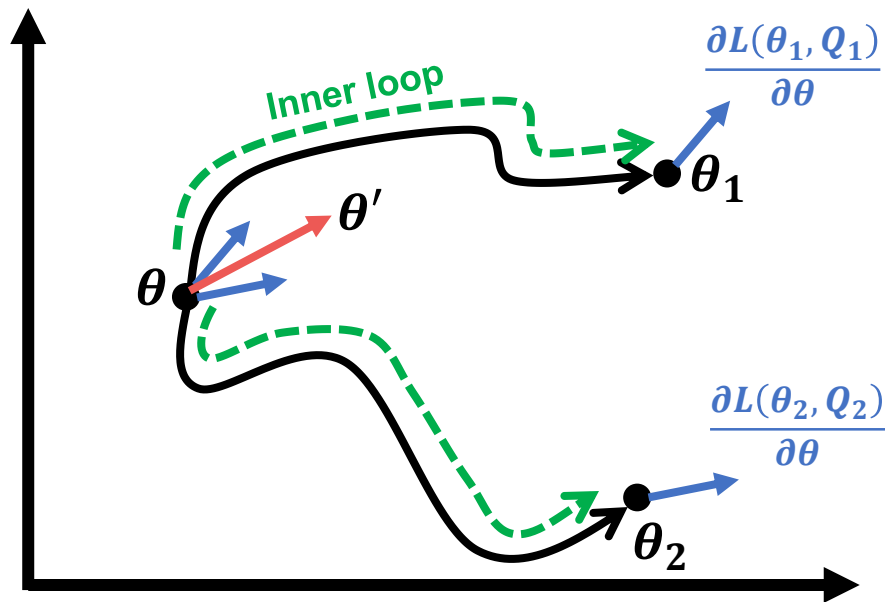
  









Task 2	Support $S_2$		Query $Q_2$	
Class 1				
Class 2				

- Goal: Minimize  $L(\theta_1, Q_1)$  and  $L(\theta_2, Q_2)$  by finding best  $\theta$ .









# Introduction

MAML makes model learn to classify after seeing little data.



Task 1	Support $S_1$		Query $Q_1$	
Class 1				
Class 2				

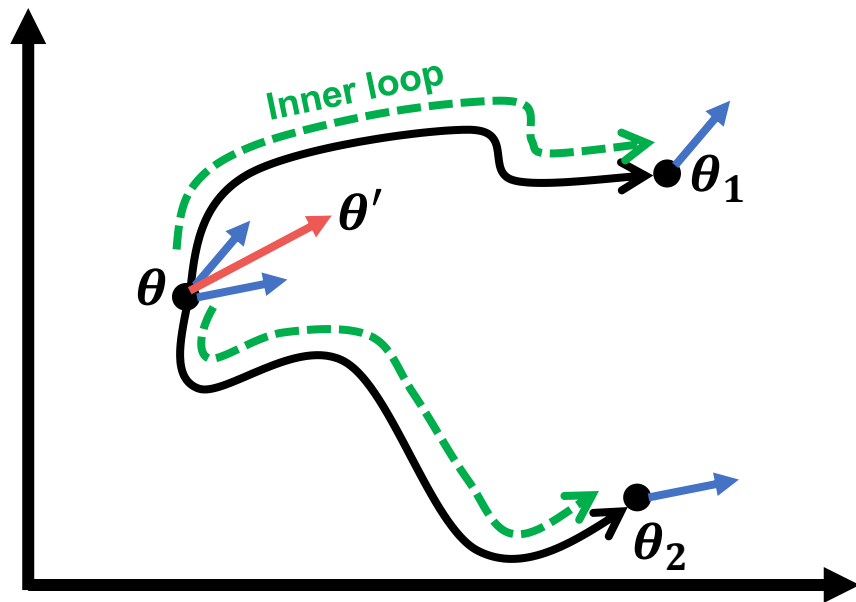
  

Task 2	Support $S_2$		Query $Q_2$	
Class 1				
Class 2				

- Goal: Minimize  $L(\theta_1, Q_1)$  and  $L(\theta_2, Q_2)$  by finding best  $\theta$ .

# Introduction

MAML makes model learn to classify after seeing little data.



## Algorithm 1 Second-order MAML

```
1: while not done do
2:   Sample tasks  $\{T_1, T_2\}$ 
3:   for  $n = 1, 2$  do
4:      $\{S_n, Q_n\} \leftarrow$  sample from  $T_n$ 
5:      $\theta_n = \theta$ 
6:     for  $i = 1, 2, \dots, N_{step}$  do
7:        $\theta_n \leftarrow \theta_n - \eta \nabla_{\theta_n} L(\theta_n, S_n)$ 
8:     end for
9:   end for
10:  Update  $\theta \leftarrow \theta - \rho \sum_{n=1}^{N_{batch}} \nabla_{\theta} L(\theta_n, Q_n)$ 
11: end while
```

- Goal: Minimize  $L(\theta_1, Q_1)$  and  $L(\theta_2, Q_2)$  by finding best  $\theta$ .
- Method: update  $\theta$  by  $\theta' = \theta - \sum \frac{\partial L(\theta_n, Q_n)}{\partial \theta}$

# Introduction

MAML makes model learn to classify after seeing little data.

Why is MAML successful?

- It is widely believed that MAML encourages models to learn a general-purpose representations which are applicable to novel tasks.

# Introduction

MAML makes model learn to classify after seeing little data.

Why is MAML successful?

- It is widely believed that MAML encourages models to learn a general-purpose representations which are applicable to novel tasks.

In this paper, we step further and ask:

- How does MAML encourage any model to learn general-purpose representations?
- What is the role of the support and query data and how do they interact with each other?

# Introduction

MAML makes model learn to classify after seeing little data.

Why is MAML successful?

- It is widely believed that MAML encourages models to learn a general-purpose representations which are applicable to novel tasks.

In this paper, we step further and ask:

- How does MAML encourage any model to learn general-purpose representations?
- What is the role of the support and query data and how do they interact with each other?

Our contribution:

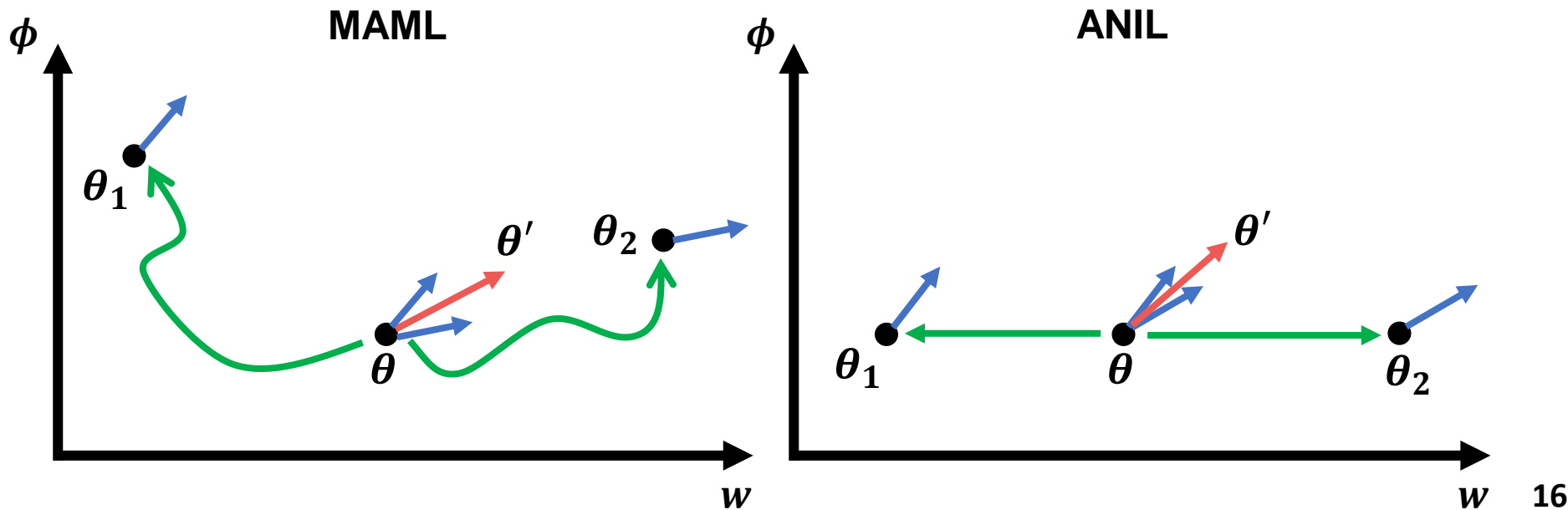
- We show that MAML is a noisy supervised contrastive learning algorithm.

# Assumption

ANIL (Almost no inner loop)

Consider a model  $\theta = \{\phi, w\}$ , where  $\phi$  is an encoder and  $w$  is a linear classifier.

ANIL states that the encoder  $\phi$  is not updated during the inner loop.



# Assumption

## ANIL (Almost no inner loop)

Consider a model  $\theta = \{\phi, w\}$ , where  $\phi$  is an encoder and  $w$  is a linear classifier.

ANIL states that the encoder  $\phi$  is not updated during the inner loop.

The ANIL Assumption empirically sounds.

	Mini-ImageNet 5way-1shot	Mini-ImageNet 5way-5shot	Omniglot 20way-1shot	Omniglot 20way-5shot
MAML	<b>46.9<math>\pm</math>0.2</b>	<b>63.1<math>\pm</math>0.4</b>	93.7 $\pm$ 0.7	96.4 $\pm$ 0.1
ANIL	46.7 $\pm$ 0.4	61.5 $\pm$ 0.5	<b>96.2<math>\pm</math>0.5</b>	<b>98.0<math>\pm</math>0.3</b>



# Main Derivation

## Motivating example

Assumptions:

- ANIL.
- Linear classifier  $w$  is zeroed at the beginning.

Loss: Mean square error.

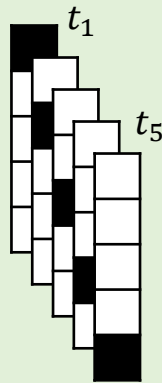
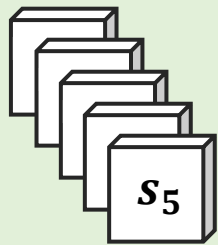
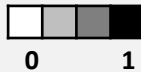
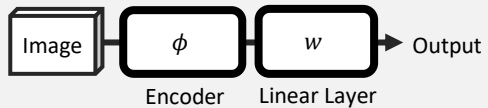
Condition: One inner loop update.

Setting

- 5-way: Each task contains 5 classes of images.
- 1-shot: Only one image per class in the support data.

**Setting:**  
5-way 1-shot using MAML with one  
inner-loop update under MSE loss.

**Model:**

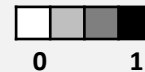
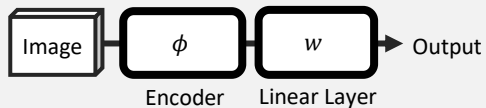


**Inner Loop (1 step)**

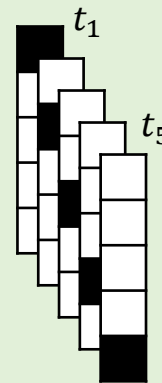
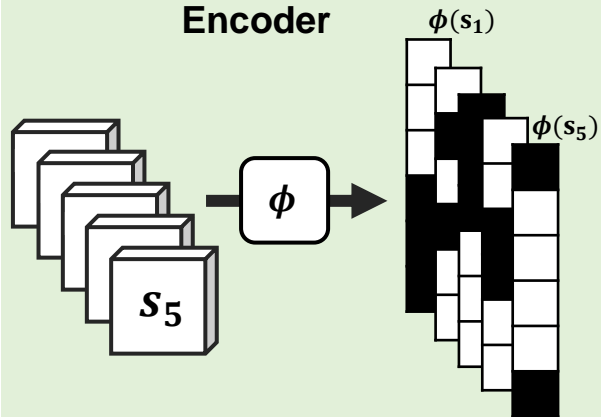
Setting:

**5-way 1-shot** using MAML with one  
inner-loop update under MSE loss.

Model:



**Encoder**

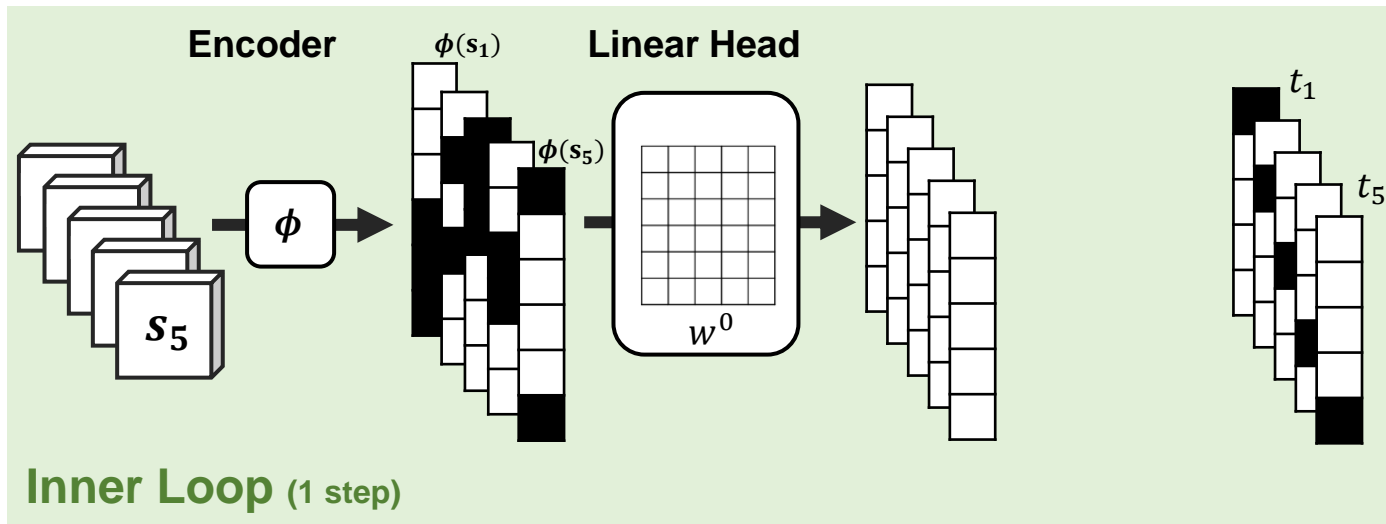
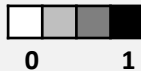
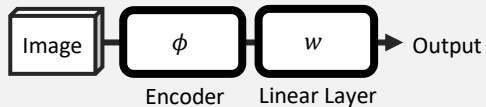


**Inner Loop (1 step)**

Setting:

**5-way 1-shot** using MAML with one  
inner-loop update under MSE loss.

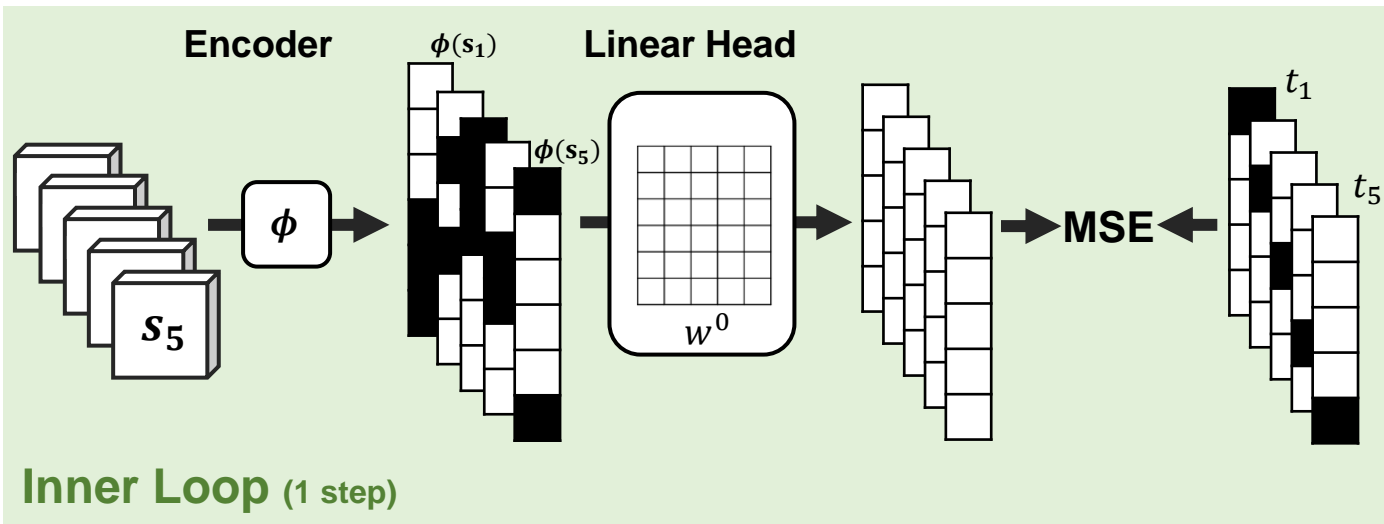
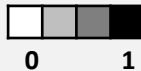
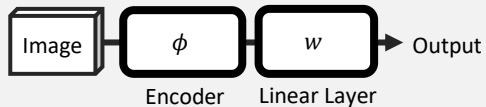
Model:



Setting:

**5-way 1-shot** using MAML with one  
inner-loop update under MSE loss.

Model:

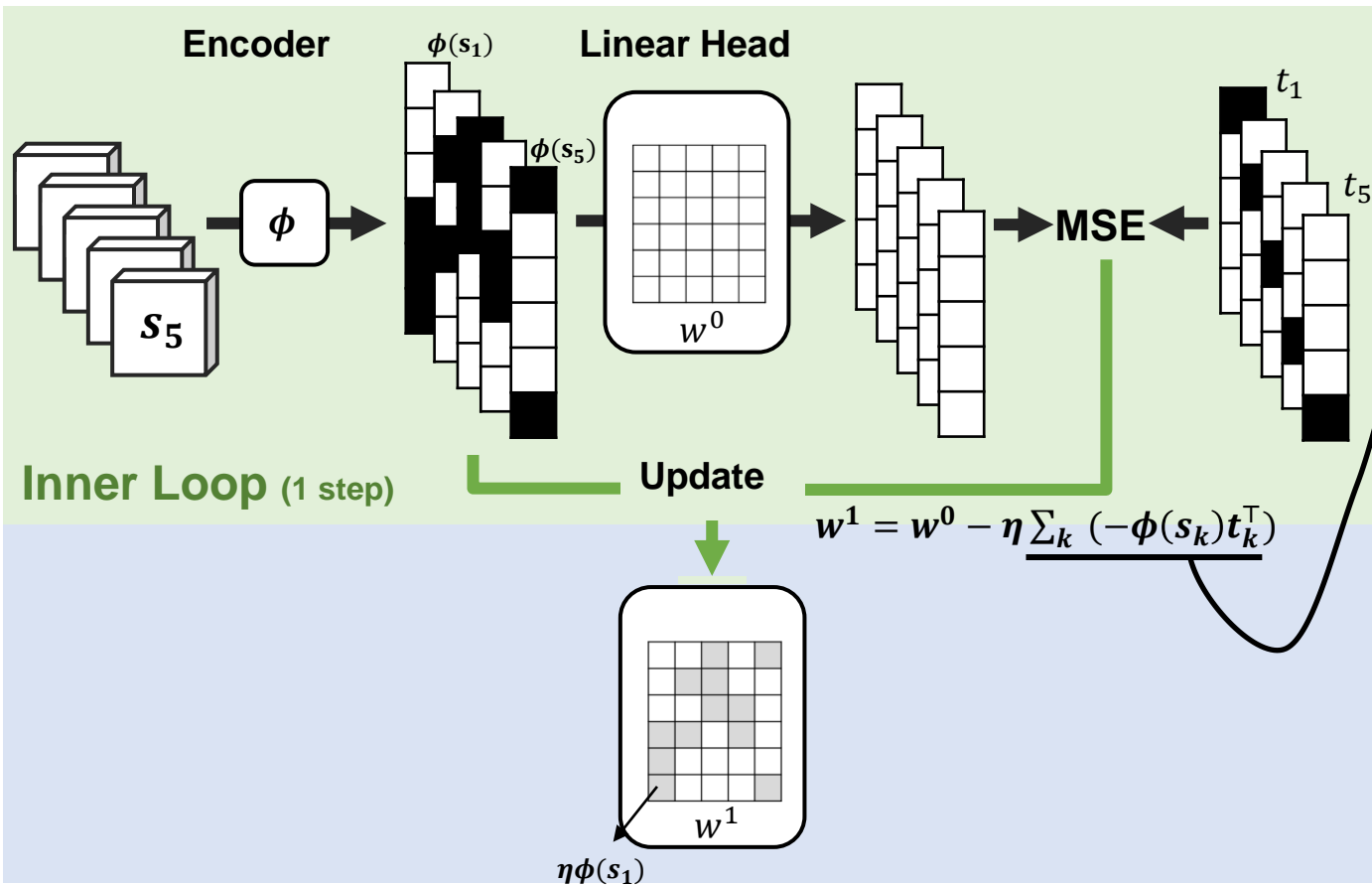


**5-way 1-shot** using MAML with **one** inner-loop update under **MSE loss**.

```

graph LR
    Image[Image] --> Encoder["\phi"]
    Encoder --> LinearLayer["w"]
    LinearLayer --> Output[Output]
    subgraph Labels
        Encoder --- EncoderLabel[Encoder]
        LinearLayer --- LinearLabel[Linear Layer]
    end

```

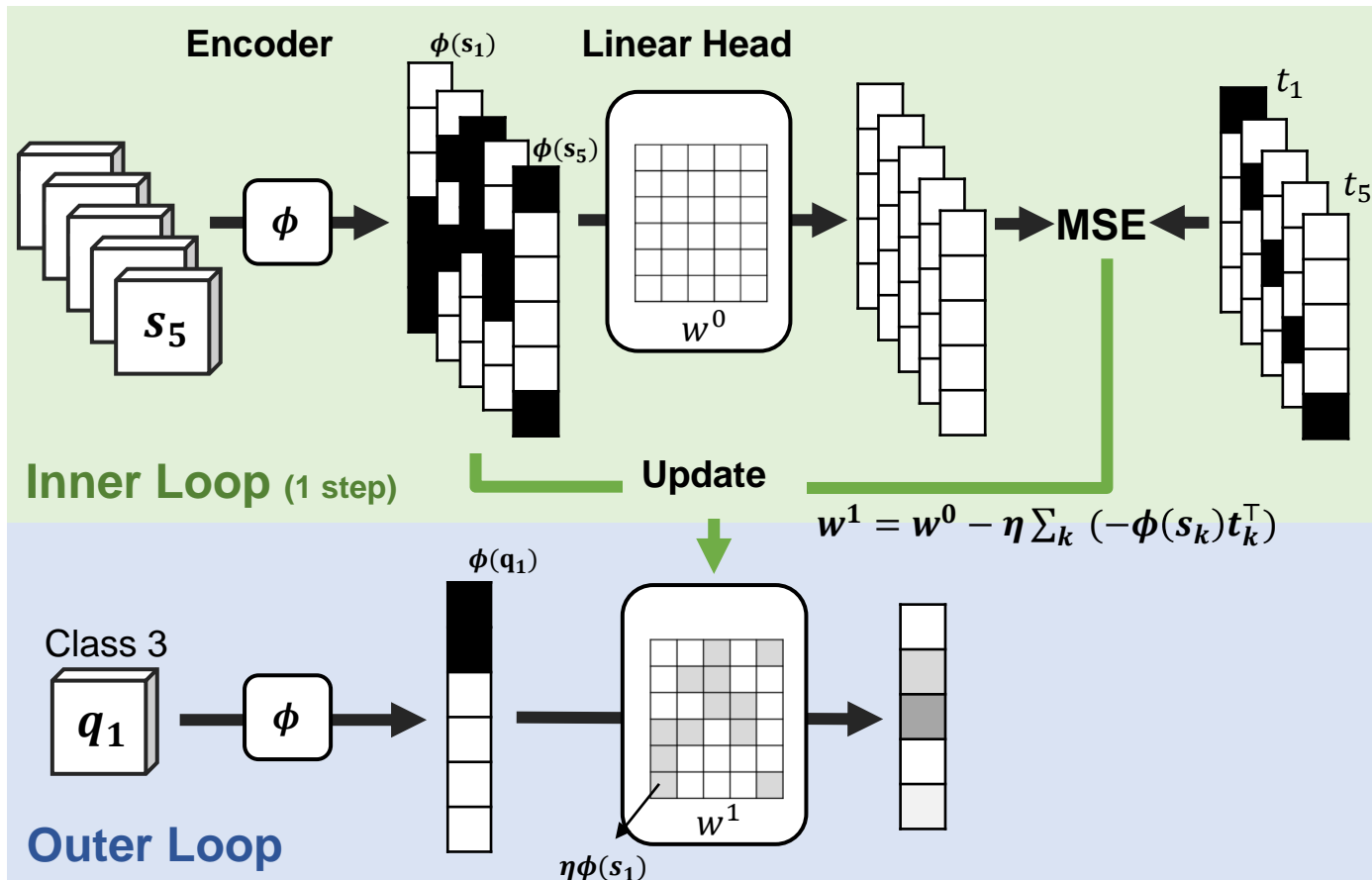
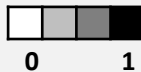
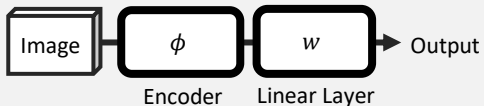


The diagram shows a 5x5 grid being converted into a column vector  $\phi(s_1)$  and a row vector  $t_1^\top$ . The grid is divided into four quadrants: top-left (white), top-right (black), bottom-left (gray), and bottom-right (white). The column vector  $\phi(s_1)$  is a 5x1 vector where the top two cells are white, the middle two are black, and the bottom one is gray. The row vector  $t_1^\top$  is a 1x5 vector where the first two cells are black, the middle two are white, and the last one is gray. This process is repeated for other states  $s_2, \dots$ .

Setting:

**5-way 1-shot** using MAML with one  
inner-loop update under MSE loss.

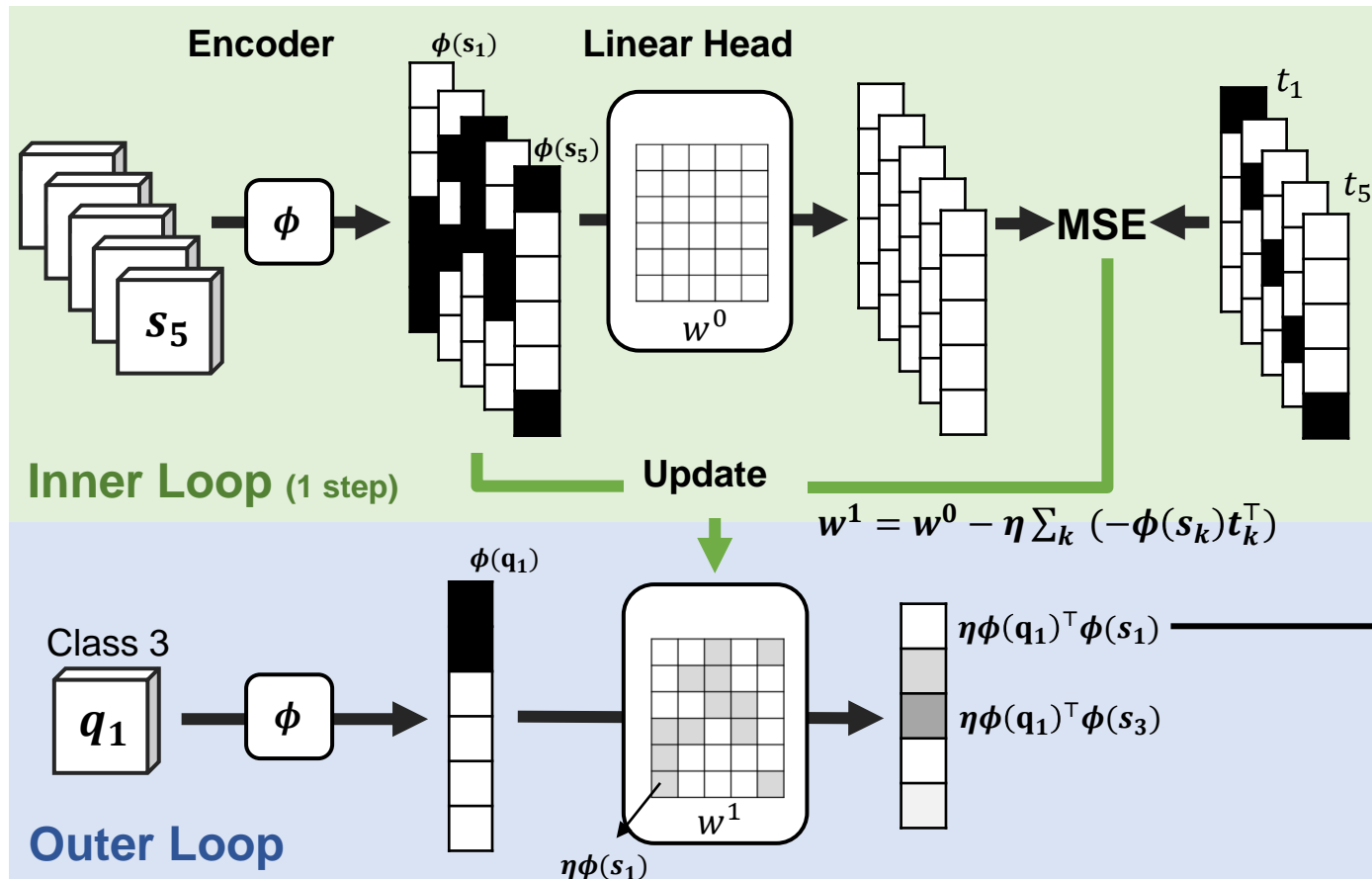
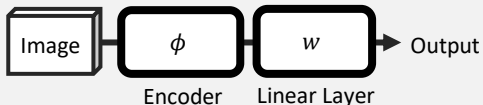
Model:



Setting:

**5-way 1-shot** using MAML with one  
inner-loop update under MSE loss.

Model:



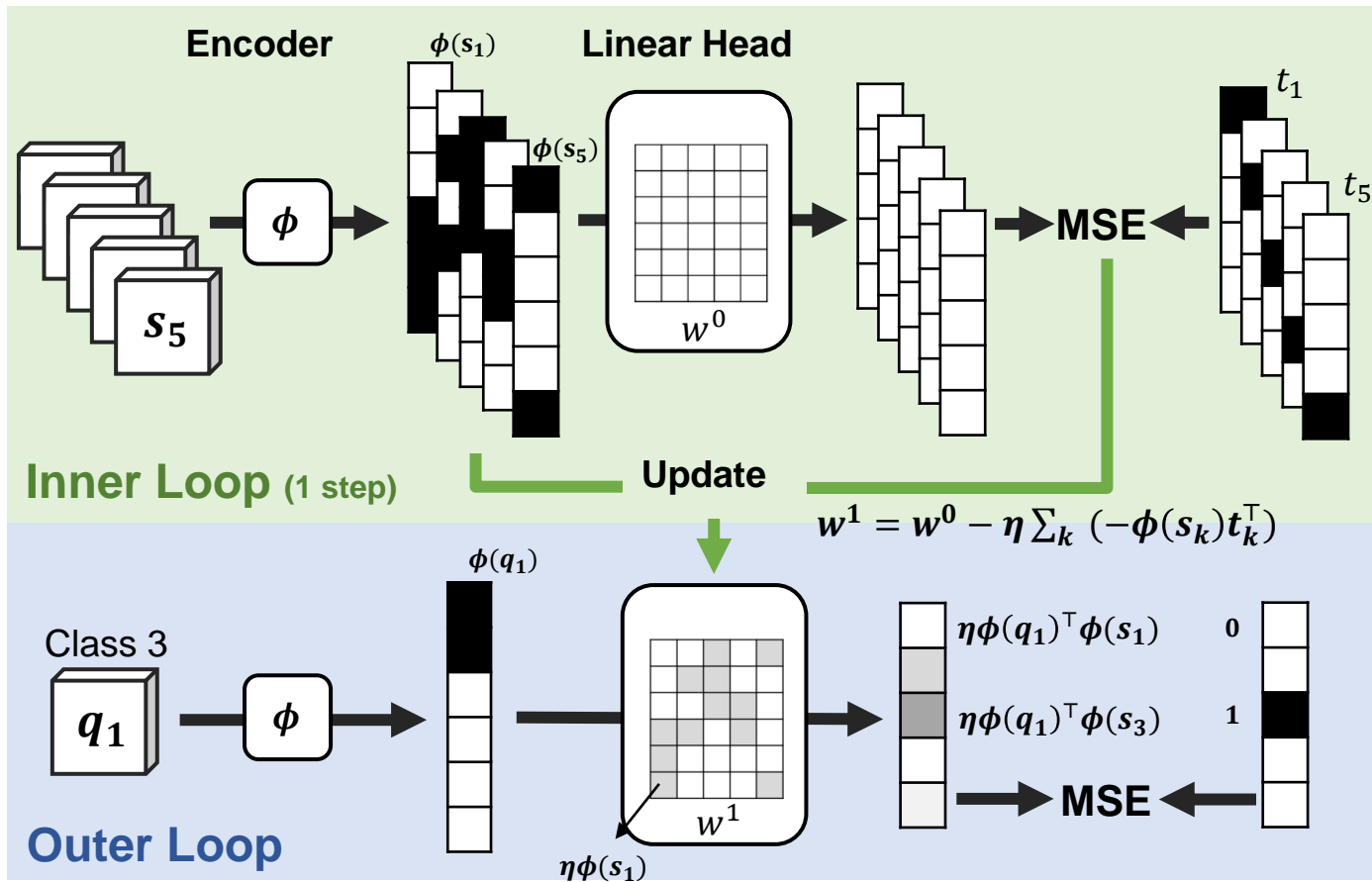
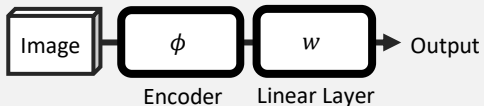
The inner product  
between support feature  
 $s_1$  and query feature  $q_1$ .



Setting:

**5-way 1-shot** using MAML with one  
inner-loop update under MSE loss.

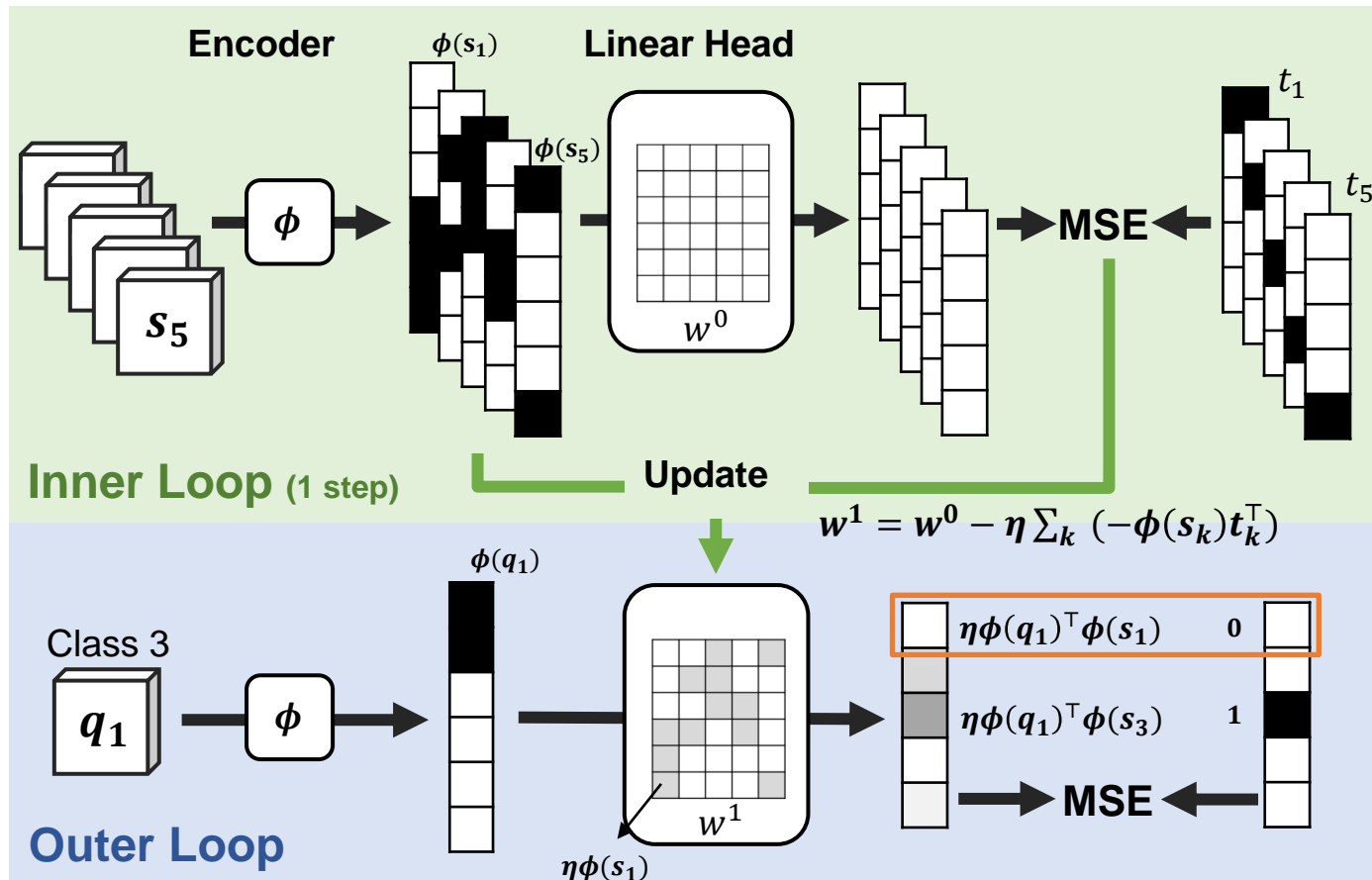
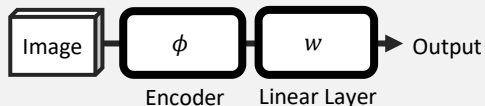
Model:



Setting:

**5-way 1-shot** using MAML with one inner-loop update under MSE loss.

Model:



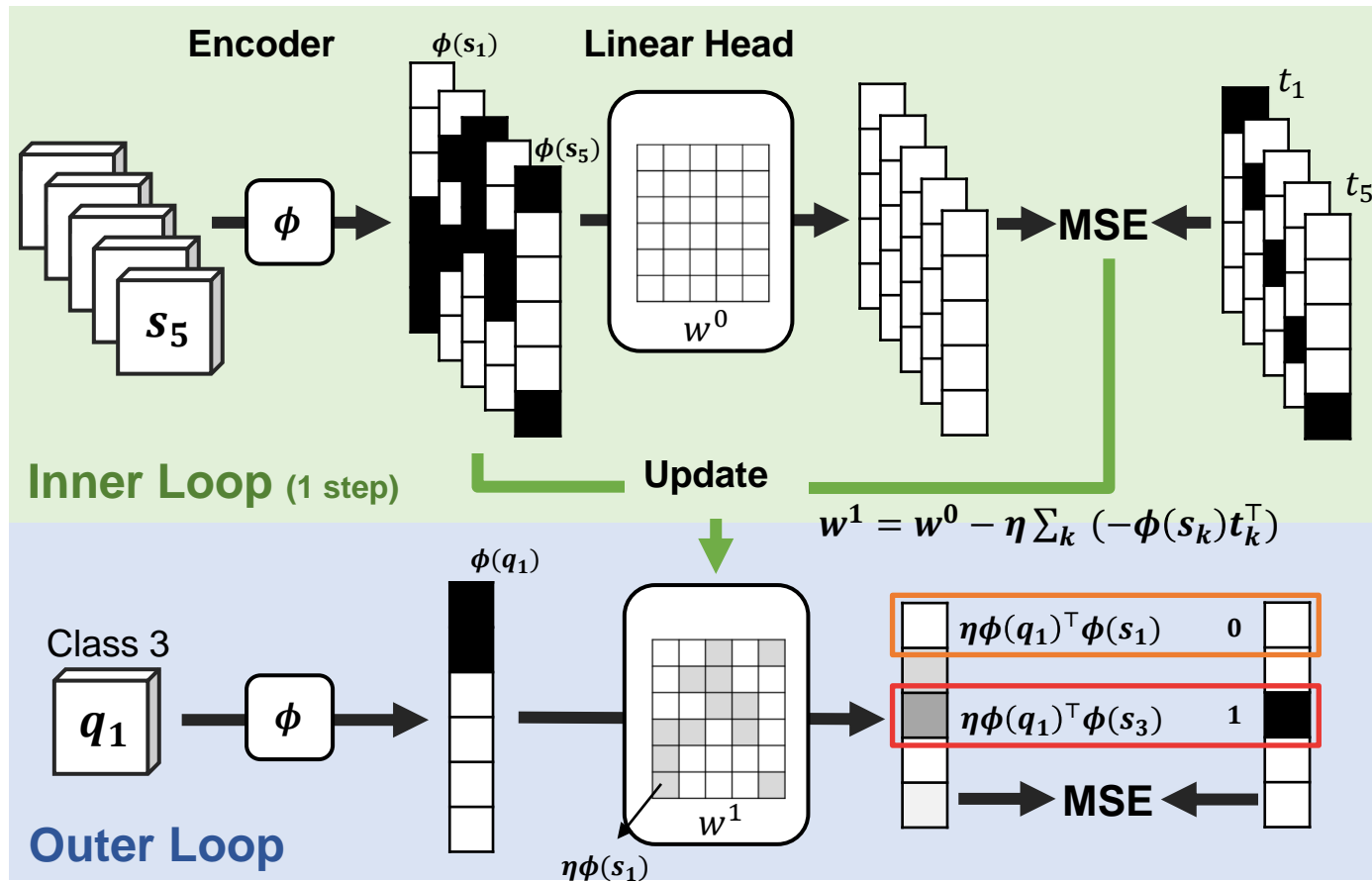
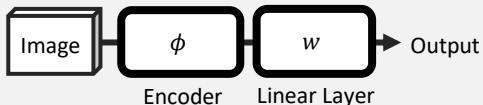
Negative sample

- $q_1$  and  $s_1$  have different labels
- Their inner product of their features should be zero.

Setting:

**5-way 1-shot** using MAML with one inner-loop update under MSE loss.

Model:



**Negative sample**

- $q_1$  and  $s_1$  have different labels
- Their inner product of their features should be zero.

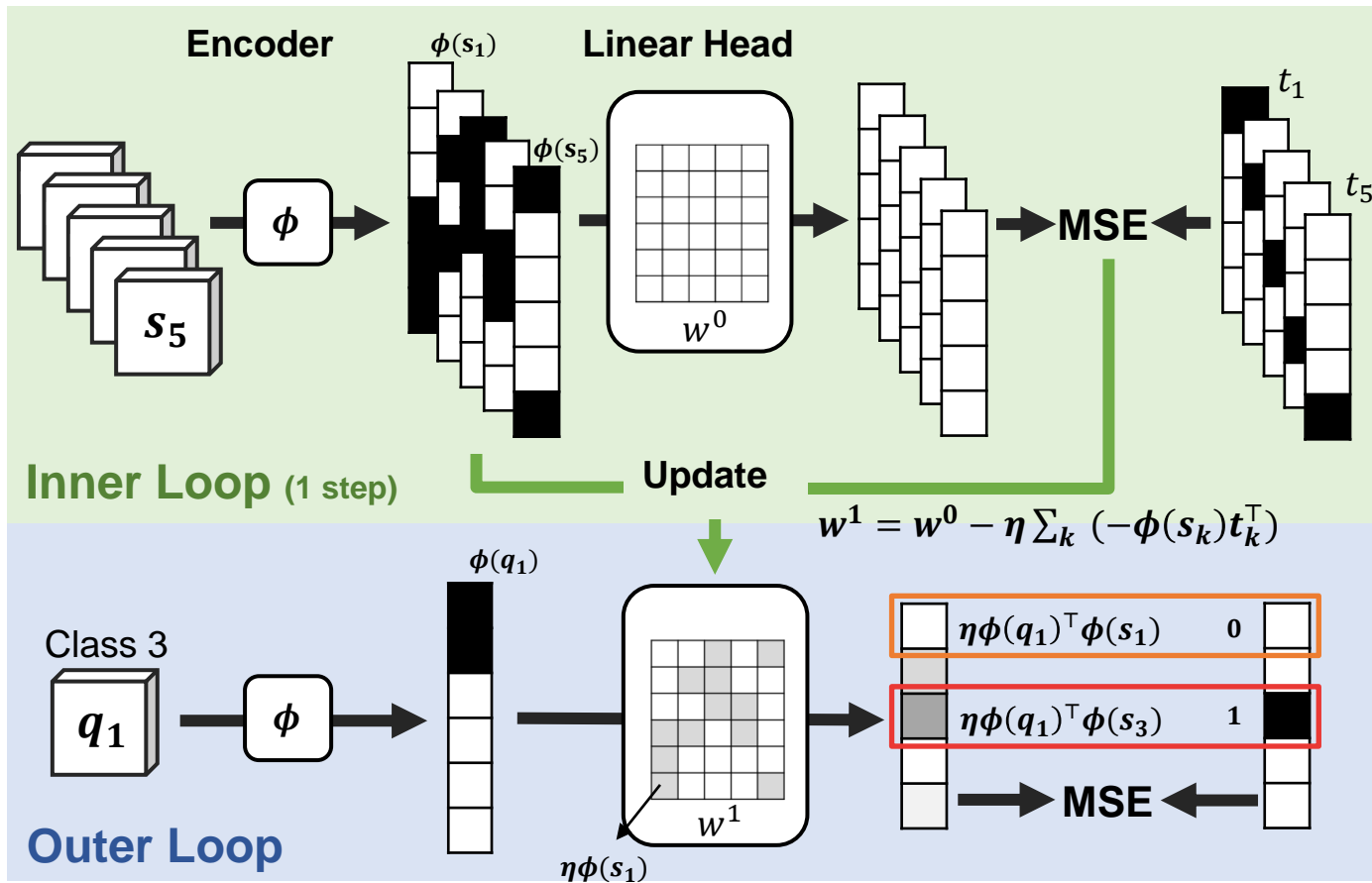
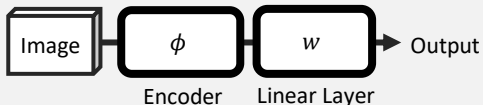
**Positive sample**

- $q_1$  and  $s_3$  have same labels,
- Their inner product of their features should be one.

Setting:

**5-way 1-shot** using MAML with one inner-loop update under MSE loss.

Model:



**Supervised  
contrastive  
learning**



**Negative sample**

- $q_1$  and  $s_1$  have different labels
- Their inner product of their features should be zero.

**Positive sample**

- $q_1$  and  $s_3$  have same labels,
- Their inner product of their features should be one.

# Main Derivation

## Main result

Consider support data  $S = \{(s, t)\}$  and one query data  $(q, u)$ .

Under **ANIL assumption**, the **loss for the encoder** is :

- First-order MAML:

$$L = \sum_{i=1}^{N_{class}} \underbrace{(q_i - 1_{i=u}) \mathbf{w}_i^\top}_{\text{stop gradient}} \phi(q) + \eta \mathbf{E}_{(s,t) \sim S} \underbrace{\left[ - \sum_{i=1}^{N_{class}} q_i s_i + s_u + q_t - 1_{t=u} \right] \phi(s)^\top \phi(q)}_{\text{stop gradient}}$$

Contrastive coefficientInner product

# Main Derivation

## Main result

Consider support data  $S = \{(s, t)\}$  and one query data  $(q, u)$ .

Under **ANIL assumption**, the **loss for the encoder** is :

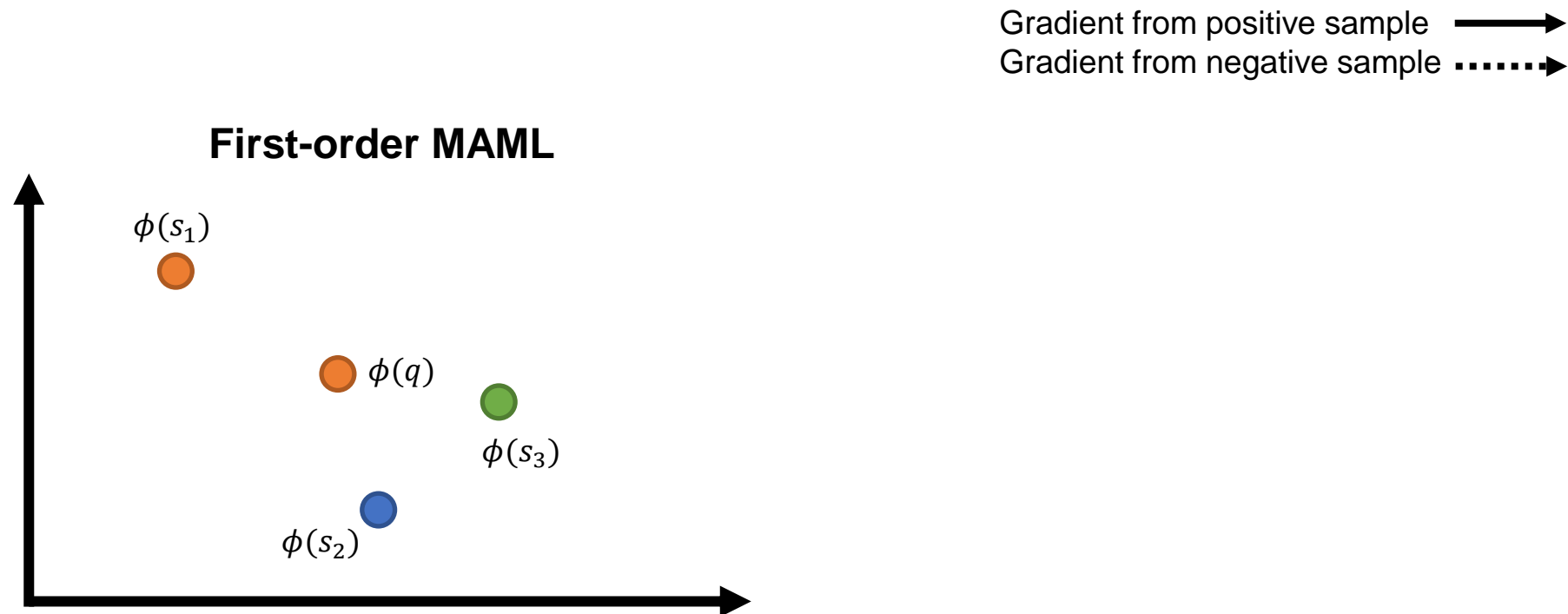
- First-order MAML:

$$L = \sum_{i=1}^{N_{class}} \underbrace{(q_i - 1_{i=u}) \mathbf{w}_i^\top}_{\text{stop gradient}} \phi(q) + \eta \underbrace{\mathbf{E}_{(s,t) \sim S} \left[ - \sum_{i=1}^{N_{class}} q_i s_i + s_u + q_t - 1_{t=u} \right] \phi(s)^\top \phi(q)}_{\substack{\text{Contrastive coefficient (CC)} \quad \text{Inner product}}}$$

- Expected effect of the second term:
  - If  $q$  and  $s$  have sample label, then the coefficients is negative, meaning that we are to update  $\phi$  s.t.  $\phi(q)$  is closer to  $\phi(s)$ .
  - If  $q$  and  $s$  have different labels, then the coefficients should be positive, meaning that we are to update  $\phi$  s.t.  $\phi(q)$  is further to  $\phi(s)$ .

# Main Derivation

Main result. Feature space illustration.



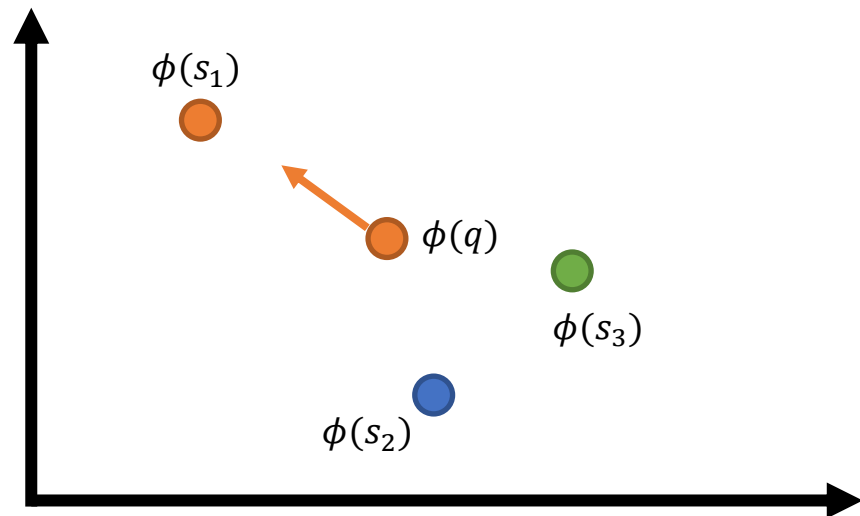
update  $\phi$  such that  $\phi(q)$  is closer/further to  $\phi(s)$

# Main Derivation

Main result. Feature space illustration.

Gradient from positive sample  $\longrightarrow$   
Gradient from negative sample  $\cdots\cdots\longrightarrow$

## First-order MAML

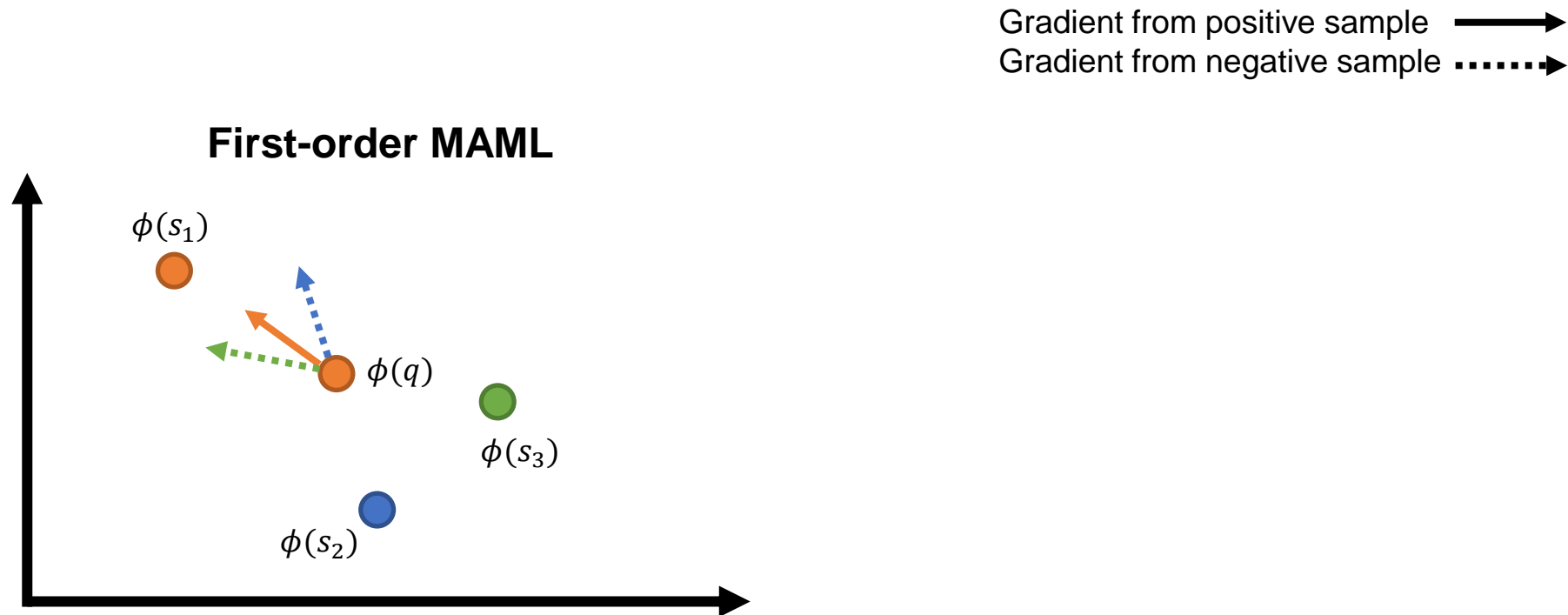


update  $\phi$  such that  $\phi(q)$  is closer/further to  $\phi(s)$



# Main Derivation

Main result. Feature space illustration.



update  $\phi$  such that  $\phi(q)$  is closer/further to  $\phi(s)$

# Main Derivation

## Main result

Consider support data  $S = \{(s, t)\}$  and one query data  $(q, u)$ .

Under **ANIL assumption**, the **loss for the encoder** is :

- First-order MAML:

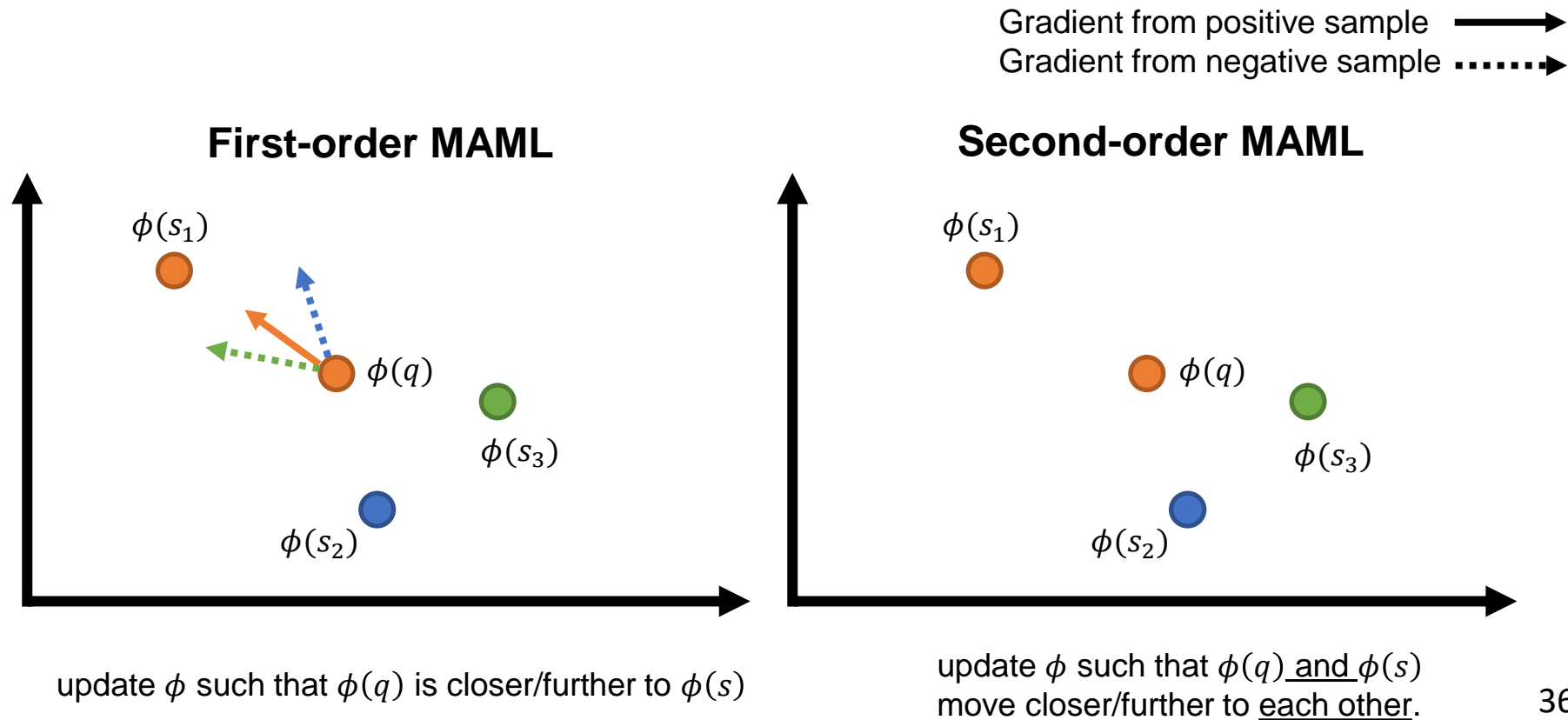
$$L = \sum_{i=1}^{N_{class}} \underbrace{(q_i - 1_{i=u}) \mathbf{w}_i^\top}_{\text{stop gradient}} \phi(q) + \eta \underbrace{\mathbf{E}_{(s,t) \sim S} \left[ - \sum_{i=1}^{N_{class}} q_i s_i + s_u + q_t - 1_{t=u} \right] \phi(s)^\top \phi(q)}_{\text{Contrastive coefficient (CC) Inner product}}$$

- Second-order MAML:

$$L = \sum_{i=1}^{N_{class}} \underbrace{(q_i - 1_{i=u}) \mathbf{w}_i^\top}_{\text{stop gradient}} \phi(q) + \eta \underbrace{\mathbf{E}_{(s,t) \sim S} \left[ - \sum_{i=1}^{N_{class}} q_i s_i + s_u + q_t - 1_{t=u} \right] \phi(s)^\top \phi(q)}_{\text{Contrastive coefficient Inner product}}$$

# Main Derivation

Main result. Feature space illustration.

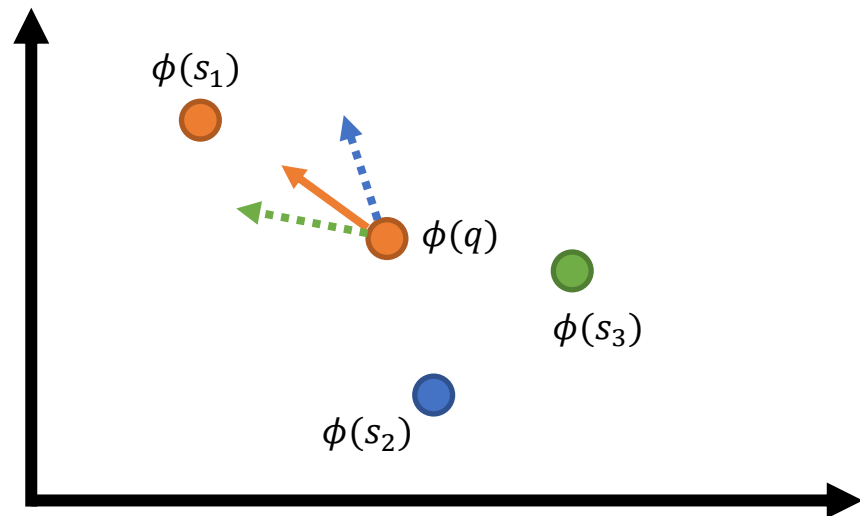


# Main Derivation

Main result. Feature space illustration.

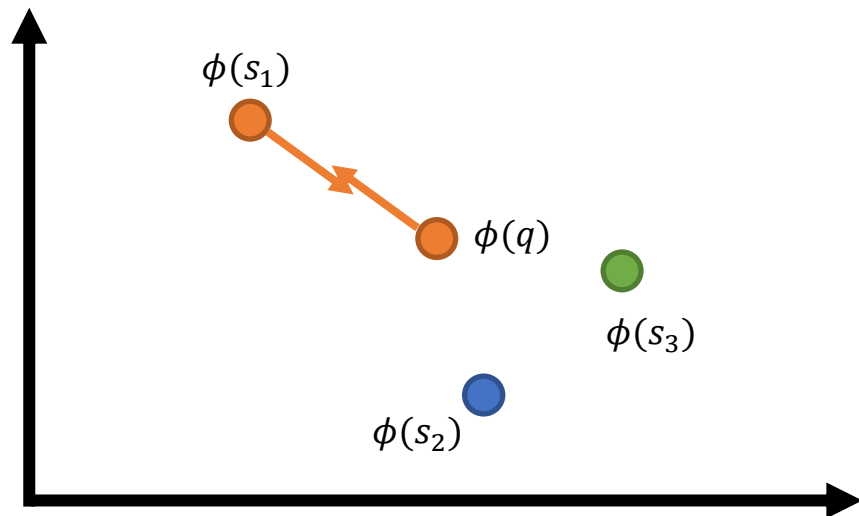
Gradient from positive sample  $\longrightarrow$   
Gradient from negative sample  $\cdots\cdots\longrightarrow$

## First-order MAML



update  $\phi$  such that  $\phi(q)$  is closer/further to  $\phi(s)$

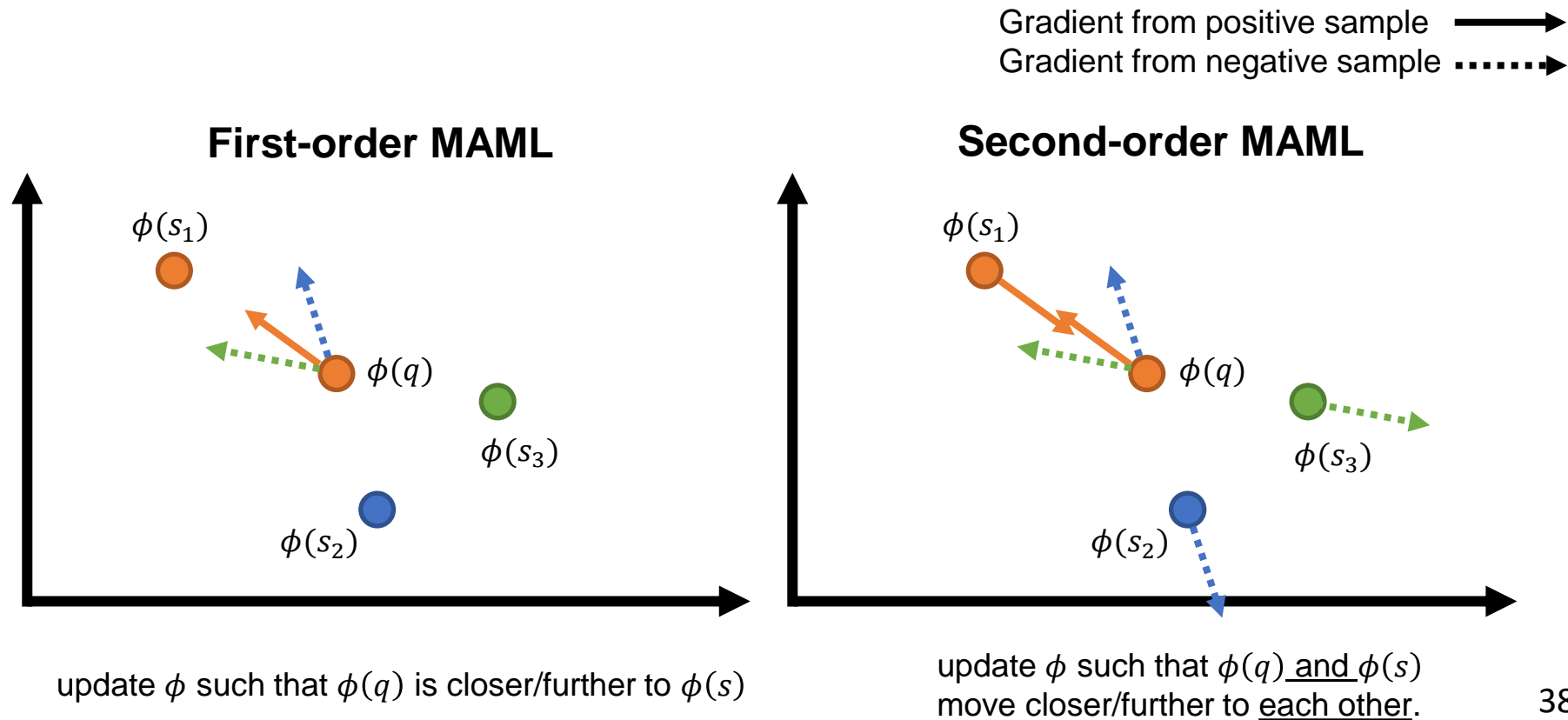
## Second-order MAML



update  $\phi$  such that  $\phi(q)$  and  $\phi(s)$   
move closer/further to each other.

# Main Derivation

Main result. Feature space illustration.



# Main Derivation

## Main result

Consider support data  $S = \{(s, t)\}$  and one query data  $(q, u)$ .

Under **ANIL assumption**, the **loss for the encoder** is :

- First-order MAML:

$$L = \underbrace{\sum_{i=1}^{N_{class}} (\mathbf{q}_i - 1_{i=u}) \mathbf{w}_i^\top \phi(q)}_{\text{stop gradient}} + \eta \underbrace{\mathbf{E}_{(s,t) \sim S} \left[ - \sum_{i=1}^{N_{class}} \mathbf{q}_i \mathbf{s}_i + \mathbf{s}_u + \mathbf{q}_t - 1_{t=u} \right] \phi(s)^\top \phi(q)}_{\text{stop gradient}}$$

- Random initialization
- Cross task interference
- The coefficient is not always positive when  $s$  and  $q$  come from different classes

**Theorem 1** *With the assumption of (a) no inner loop update of the encoder, FOMAML is a noisy SCL algorithm. With assumptions of (a) no inner loop update of the encoder and (b) a single inner-loop update, SOMAML is a noisy SCL algorithm.*

# Main Derivation

## Main result. Introducing the zeroing trick.

Consider support data  $S = \{(s, t)\}$  and one query data  $(q, u)$ .

Under ANIL assumption and the zeroing trick, the loss for the encoder is :

- First-order MAML:

$$L = \eta \mathbf{E}_{(s,t) \sim S} \underbrace{(\mathbf{q}_t - \mathbf{t}_{t=u})}_{\text{stop gradient}} \phi(s)^\top \phi(q)$$

- Second-order MAML:

$$L = \eta \mathbf{E}_{(s,t) \sim S} \underbrace{(\mathbf{q}_t - \mathbf{t}_{t=u})}_{\text{stop gradient}} \phi(s)^\top \phi(q)$$

**Corollary 1** *With mild assumptions of (a) no inner loop update of the encoder, (b) a single inner-loop update and (c) training with the zeroing trick (i.e., the linear layer is zeroed at the end of each outer loop), both FOMAML and SOMAML are SCL algorithms.*

# Main Derivation

Main result. Introducing the zeroing trick.

---

**Algorithm 1** Second-order MAML

---

```
1: while not done do
2:   Sample tasks  $\{T_1, T_2\}$ 
3:
4:   for  $n = 1, 2$  do
5:      $\{S_n, Q_n\} \leftarrow$  sample from  $T_n$ 
6:      $\theta_n = \theta$ 
7:     for  $i = 1, 2, \dots, N_{step}$  do
8:        $\theta_n \leftarrow \theta_n - \eta \nabla_{\theta_n} L(\theta_n, S_n)$ 
9:     end for
10:   end for
11:   Update  $\theta \leftarrow \theta - \rho \sum_{n=1}^{N_{batch}} \nabla_{\theta} L(\theta_n, Q_n)$ 
12:
13: end while
```

---

---

**Algorithm 2** Second-order MAML with zeroing trick

---

```
1: while not done do
2:   Sample tasks  $\{T_1, T_2\}$ 
3:    $w = 0$ 
4:   for  $n = 1, 2$  do
5:      $\{S_n, Q_n\} \leftarrow$  sample from  $T_n$ 
6:      $\theta_n = \theta$ 
7:     for  $i = 1, 2, \dots, N_{step}$  do
8:        $\theta_n \leftarrow \theta_n - \eta \nabla_{\theta_n} L(\theta_n, S_n)$ 
9:     end for
10:   end for
11:   Update  $\theta \leftarrow \theta - \rho \sum_{n=1}^{N_{batch}} \nabla_{\theta} L(\theta_n, Q_n)$ 
12:    $w = 0$ 
13: end while
```

---



# Results

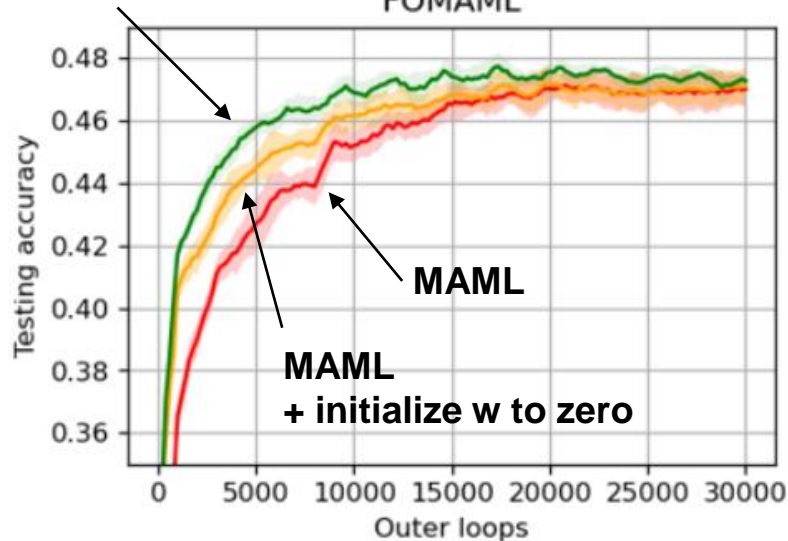
Using the zeroing trick improves performance.

Setting: MinImageNet.

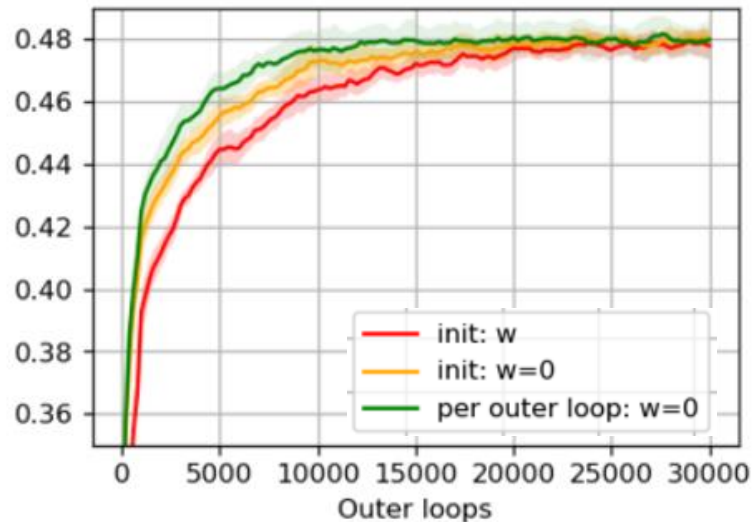
5-way 1-shot setting

MAML + zeroing trick

FOMAML



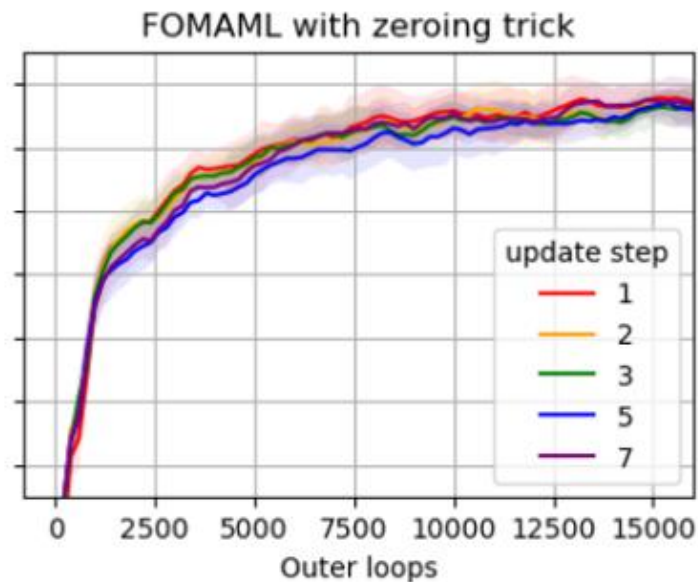
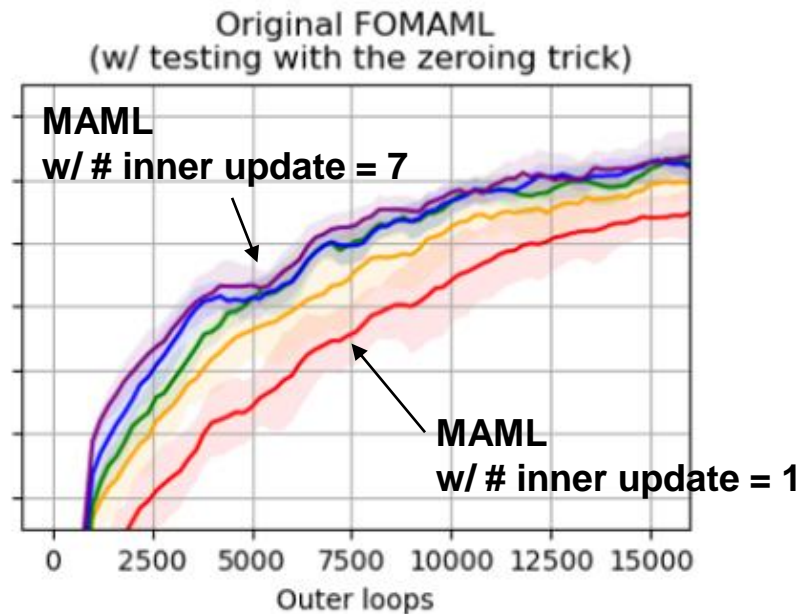
SOMAML



# Results

With zeroing trick, # of inner loop steps no longer matters

Setting: MinImageNet 5-way 1-shot.





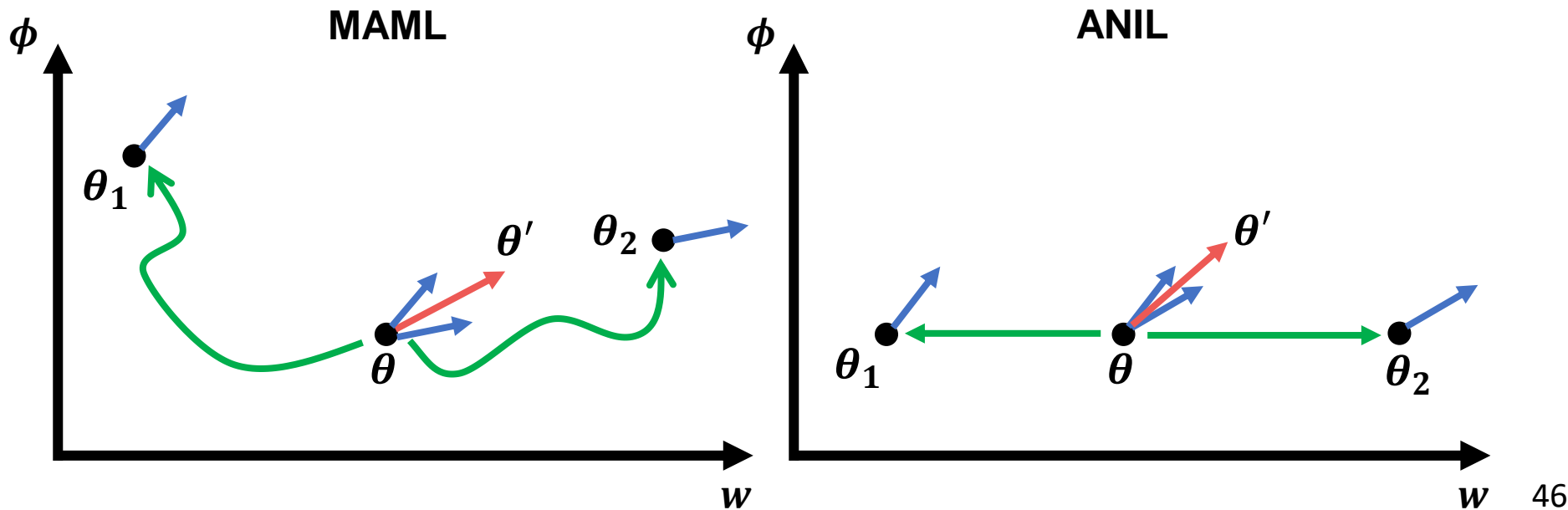
# Wrap up

# Wrap up

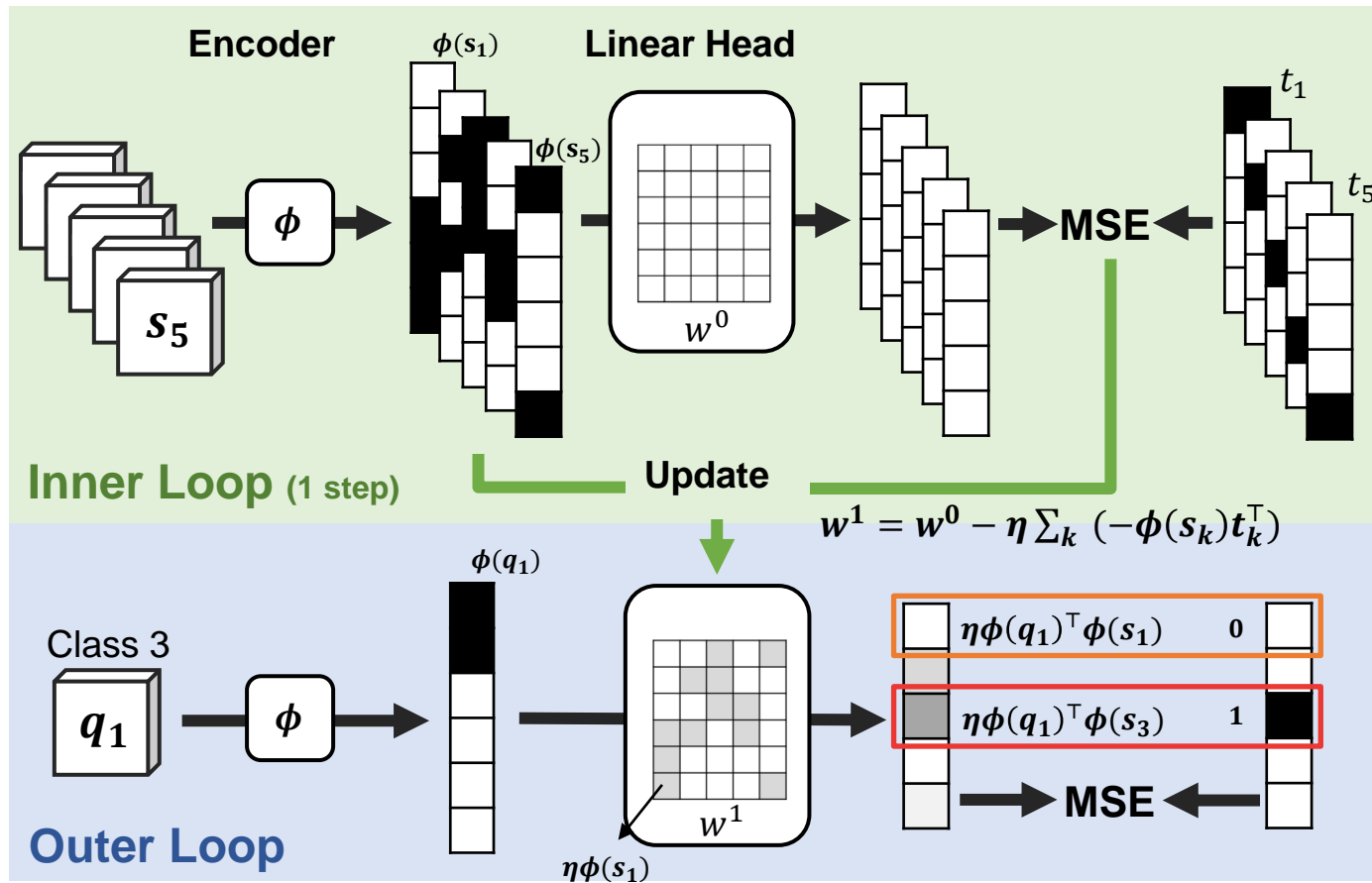
We use the ANIL assumption for derivation.

Consider a model  $\theta = \{\phi, w\}$ , where  $\phi$  is an encoder and  $w$  is a linear classifier.

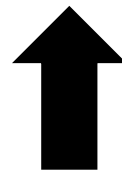
ANIL states that the encoder  $\phi$  is not updated during the inner loop.



# Wrap up



**Supervised  
contrastive  
learning**



**Negative sample**

- $q_1$  and  $s_1$  have different labels
- Their inner product of their features should be zero.

**Positive sample**

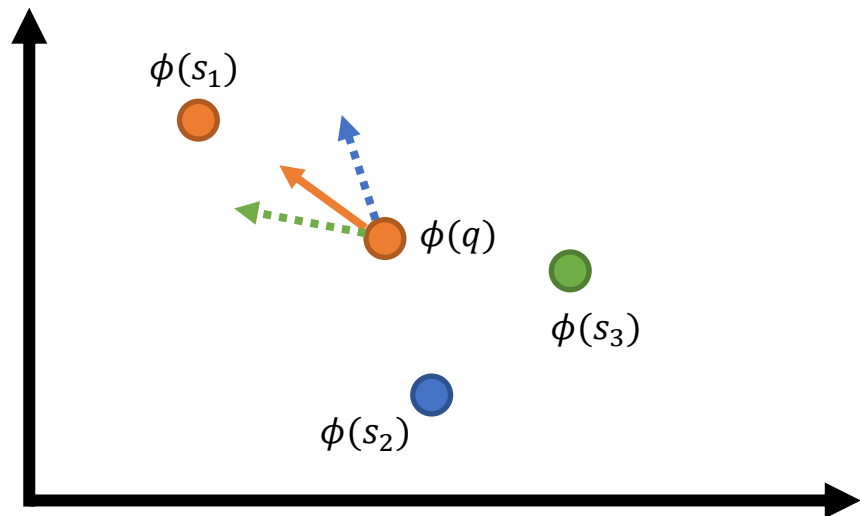
- $q_1$  and  $s_3$  have same labels,
- Their inner product of their features should be one.

# Wrap up

We show how FOMAML different from SOMAML.

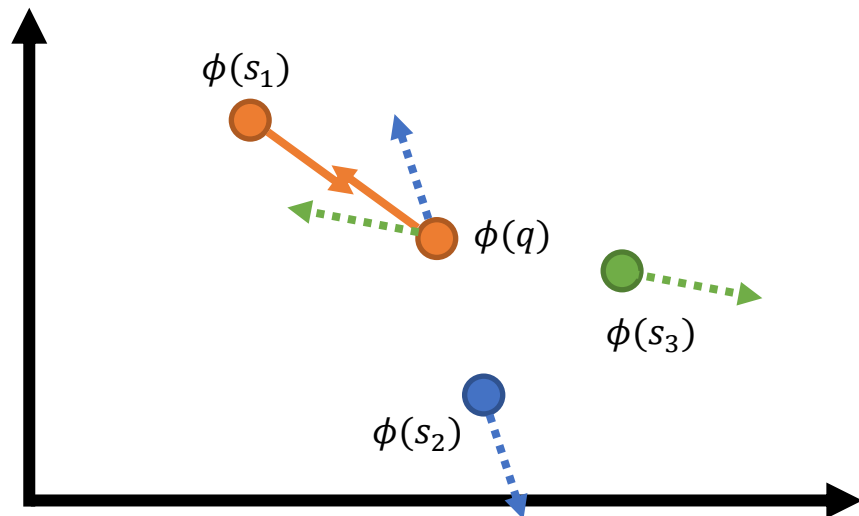
Gradient from positive sample  $\longrightarrow$   
Gradient from negative sample  $\cdots\cdots\longrightarrow$

**First-order MAML**



update  $\phi$  such that  $\phi(q)$  is closer/further to  $\phi(s)$

**Second-order MAML**



update  $\phi$  such that  $\phi(q)$  and  $\phi(s)$   
move closer/further to each other.

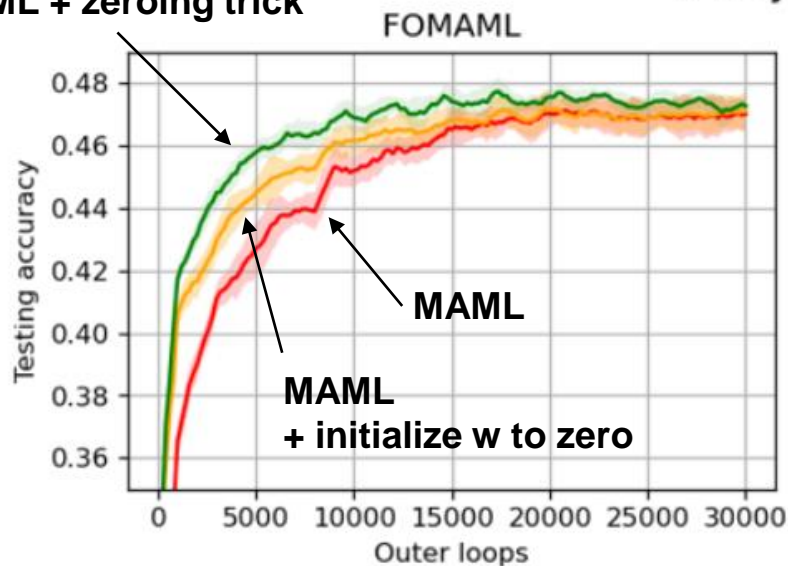
# Wrap up

We show that the zeroing trick improves MAML.

Setting: MinImageNet.

5-way 1-shot setting

MAML + zeroing trick



SOMAML

