

# Revisit Multimodal Meta-Learning through the Lens of Multi-Task Learning

Milad Abdollahzadeh, Toubia Malekzadeh, Ngai-Man Cheung  
Singapore University of Technology and Design

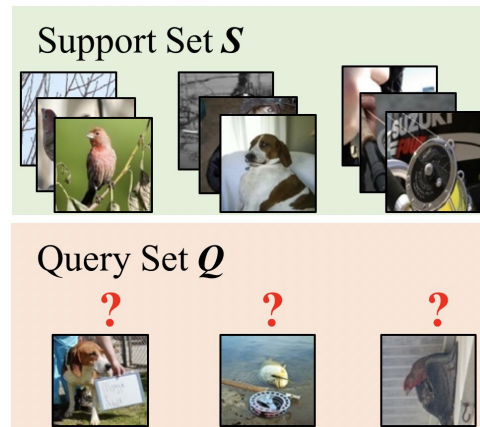
NeurIPS 2021

# Meta-Learning

**Meta-learning** aims to improve the learning algorithm

Providing an opportunity to tackle many challenges of deep learning: data-efficient AI

**Meta-Learning as a  
solution for few-shot  
learning**



# Meta-Learning

## Episodic Learning

Match the condition in which the model is trained (meta-train) and tested (meta-test)

Given 1 example of 5 classes:



Classify new examples



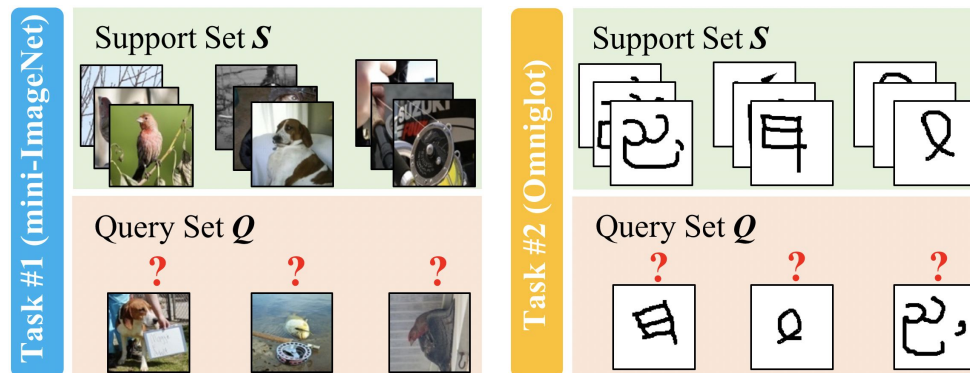
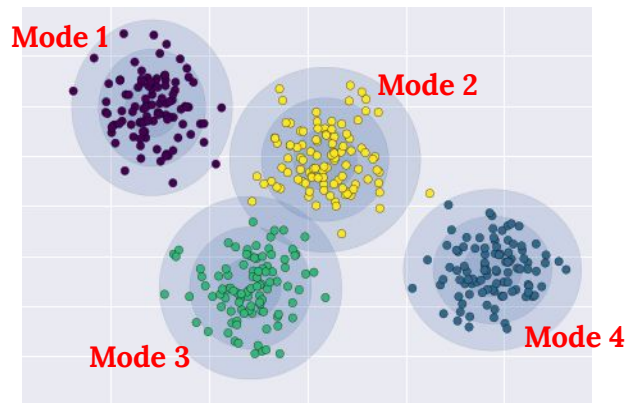
**5-way, 1-shot image classification problem**



# Multimodal Meta-Learning

Generalizing the conventional setup for more diverse task distributions [1]

Multimodality in task distribution



More challenging to current meta-learners

[1] Vuorio, Risto, et al. "Multimodal Model-Agnostic Meta-Learning via Task-Aware Modulation." *Advances in Neural Information Processing Systems* 32 (2019): 1-12.

# Research Gap

## **Main Claim** of previous work:

Improving generalization by transferring knowledge between different modes of task distributions

However, there are two main questions that need to be addressed:

- I) Not clear how task from one mode impacts the learning of task from another mode
- II) Can we make a better use of knowledge transfer and have a better generalization?

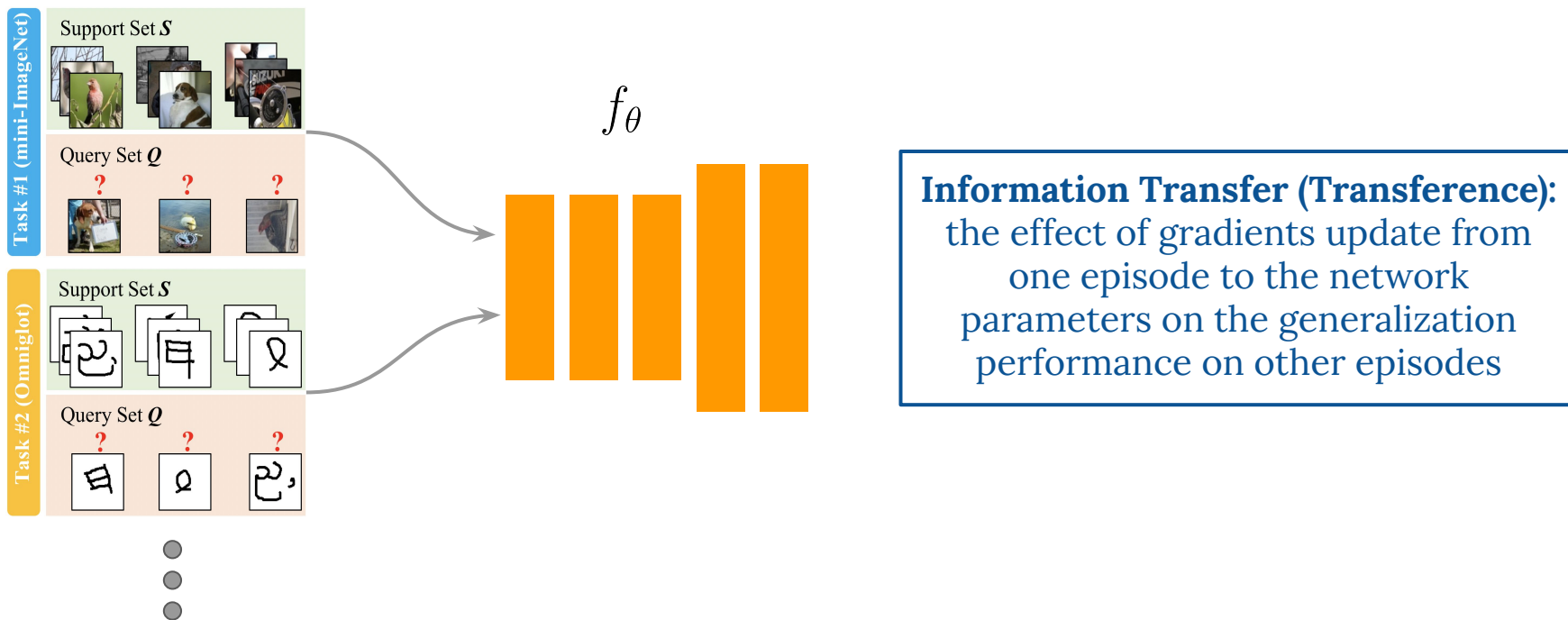
# Our Contributions

To **address** the discussed research gaps:

- 1) We extend the **transference idea** [1] from multi-task learning to episodic learning scenario of meta-learning to analyse information transfer between few-shot tasks
- 2) Propose **a new multimodal meta-learning** called **Kernel Modulation (KML)** which significantly advances state-of-the-art

# Analysis of Information Transfer (Transference)

## Episodic learning from the lens of Multi-Task learning



# Analysis of Information Transfer (Transference)

Measure transference from source (meta-train) task  $i$  to target (meta-test) task  $j$

Calculate the loss on target task before and after updating parameters wrt source task

The ratio between the loss of task  $j$  after and before parameter update w.r.t task  $i$

$$LR_{i \rightarrow j} = \frac{\mathcal{L}_{\mathcal{T}_j}(\mathcal{Q}_j; \theta_i^{t+1}, \mathcal{S}_j)}{\mathcal{L}_{\mathcal{T}_j}(\mathcal{Q}_j; \theta^t, \mathcal{S}_j)}$$

$LR < 1$  : **Positive knowledge transfer** from source task  $i$  to target task  $j$

$LR > 1$  : **Negative knowledge transfer** from source task  $i$  to target task  $j$

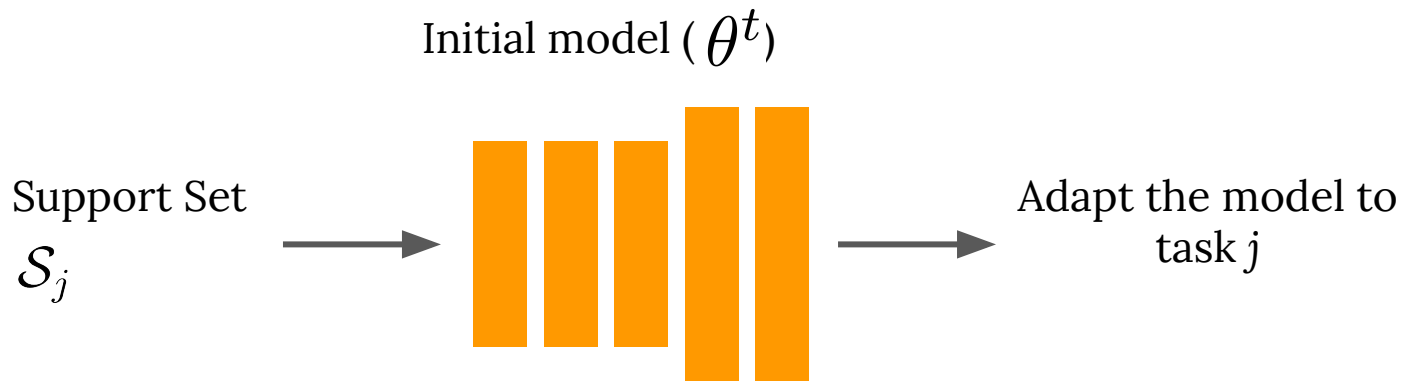


# Analysis of Information Transfer (Transference)

**Transference from source (meta-train) task  $i$  to target (meta-test) task  $j$**

**#1.** Calculate the loss of task  $j$  with initial model parameters

**a)** adapt model using the support set of target task  $j$

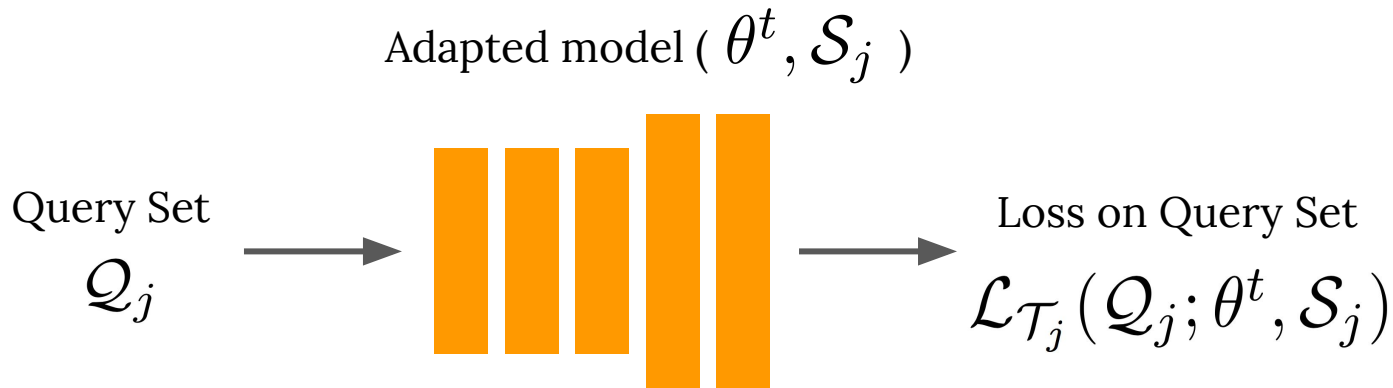


# Analysis of Information Transfer (Transference)

**Transference from source (meta-train) task  $i$  to target (meta-test) task  $j$**

**#1.** Calculate the loss of task  $j$  with initial model parameters

**b)** calculate the loss of adapted model using query set of target task  $j$

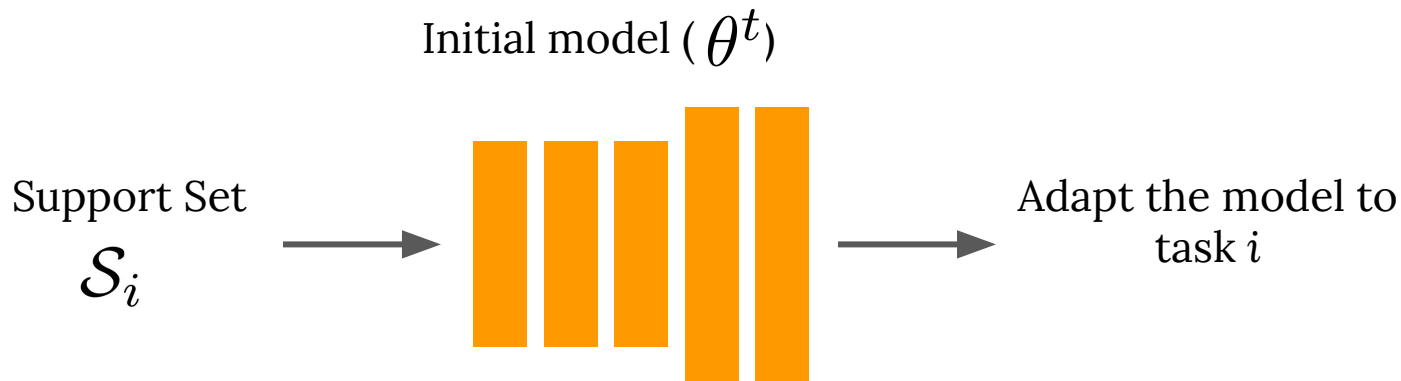


# Analysis of Information Transfer (Transference)

**Transference from source (meta-train) task  $i$  to target (meta-test) task  $j$**

**#2.** Update the initial model with respect to task  $i$

**a)** adapt the model using support set of source task  $i$

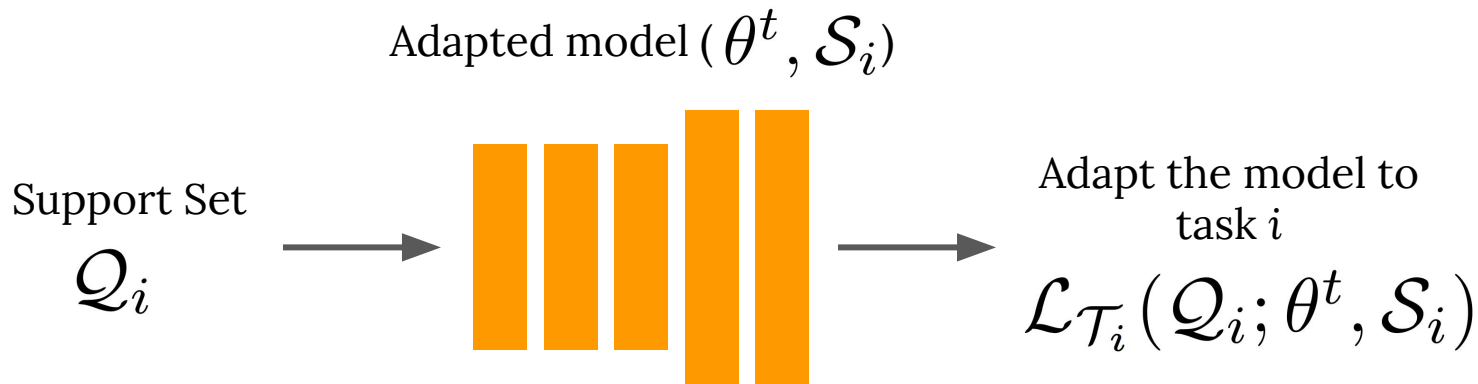


# Analysis of Information Transfer (Transference)

**Transference from source (meta-train) task  $i$  to target (meta-test) task  $j$**

**#2.** Update the initial model with respect to task  $i$

**b)** calculate the loss of adapted model using query set of source task  $i$



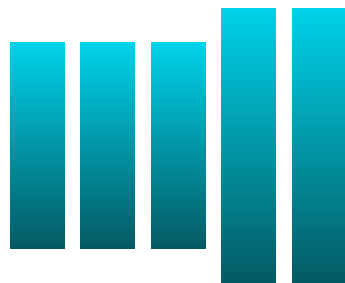
# Analysis of Information Transfer (Transference)

**Transference from source (meta-train) task  $i$  to target (meta-test) task  $j$**

**#2.** Update the initial model with respect to task  $i$

**c)** use the gradient of the loss to update model parameters

Updated model (  $\theta_i^{t+1}$  )



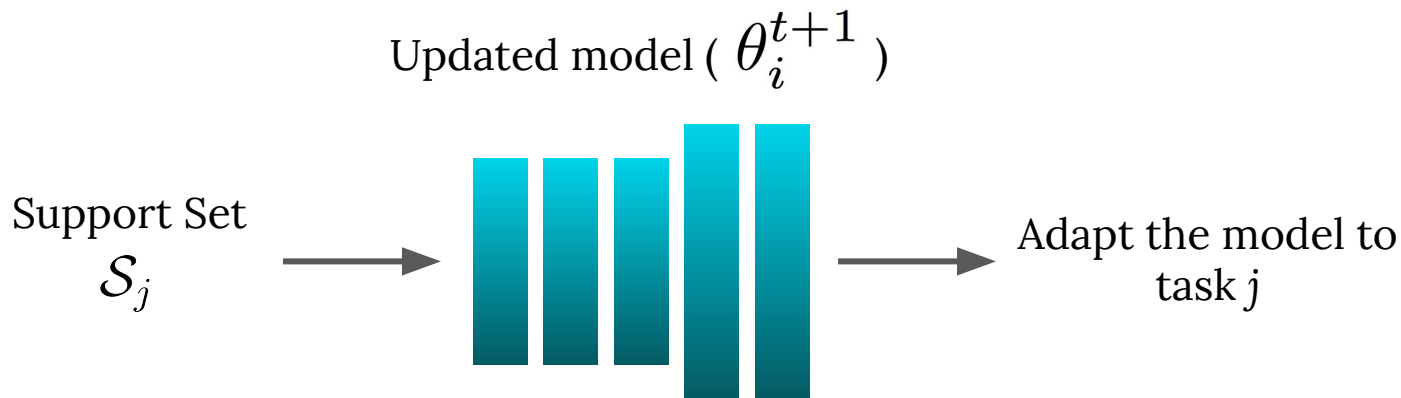
$$\theta_i^{t+1} = \theta^t - \alpha \nabla_{\theta^t} \mathcal{L}_{\mathcal{T}_i}(\mathcal{Q}_i; \theta^t, \mathcal{S}_i)$$

# Analysis of Information Transfer (Transference)

**Transference from source (meta-train) task  $i$  to target (meta-test) task  $j$**

**#3.** Calculate the loss of task  $j$  with updated model parameters

**a)** adapt model using the support set of target task  $j$

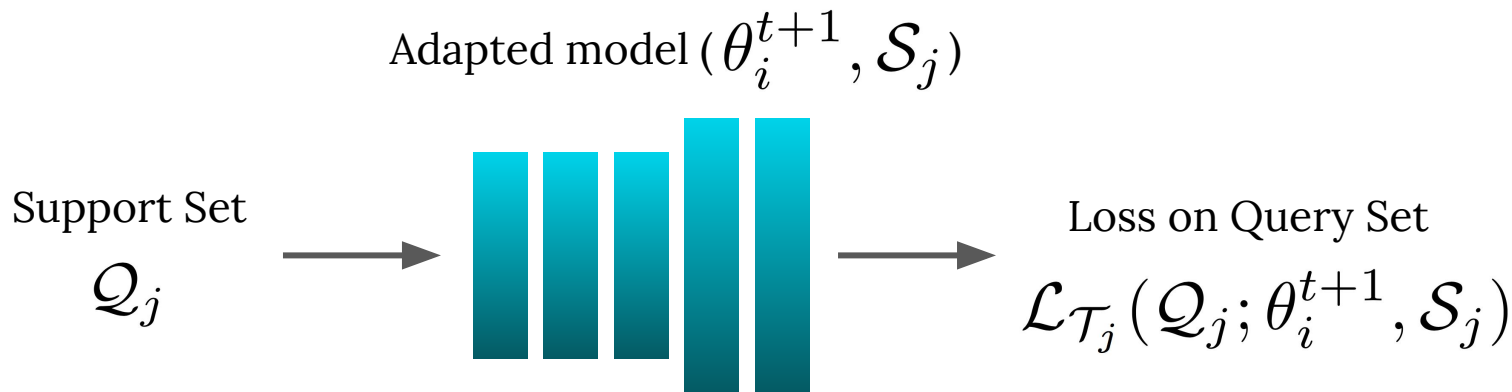


# Analysis of Information Transfer (Transference)

**Transference from source (meta-train) task  $i$  to target (meta-test) task  $j$**

**#3.** Calculate the loss of task  $j$  with updated model parameters

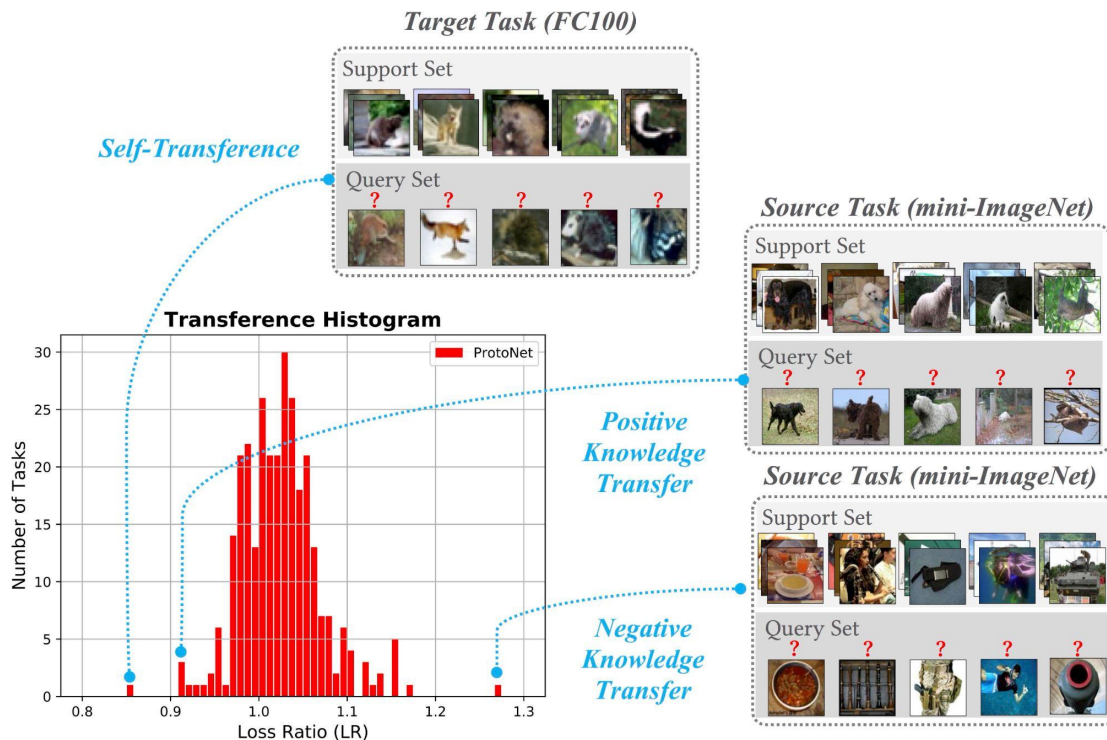
**b)** calculate the loss of adapted model using query set of target task  $j$



# Analysis of Information Transfer (Transference)

Transference histogram from  
300 meta-train  
mini-ImageNet tasks to a  
meta-test FC100 target task

LR<1: Positive Knowledge Transfer  
LR>1: Negative Knowledge Transfer





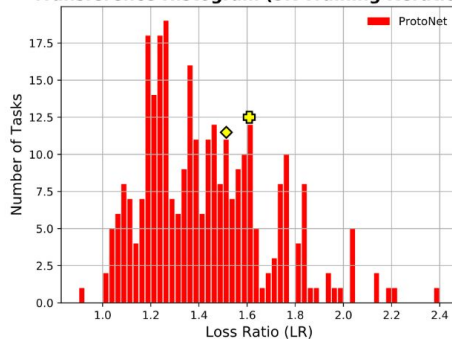
# Analysis of Information Transfer (Transference)

Transference  
Histogram  
During Training

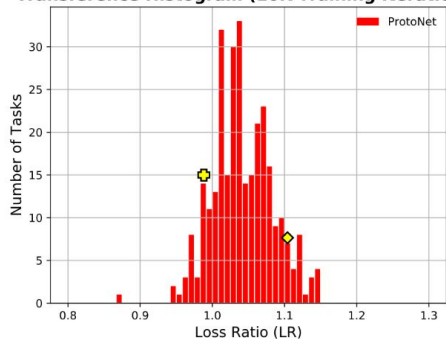


More positive  
transfer as  
training proceeds

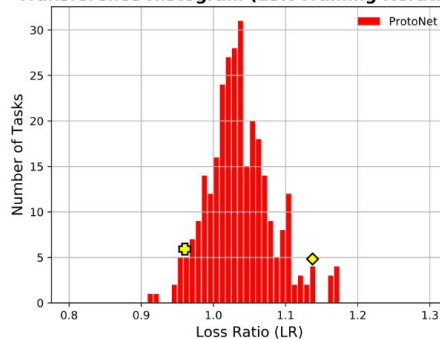
Transference Histogram (5K Training Iterations)



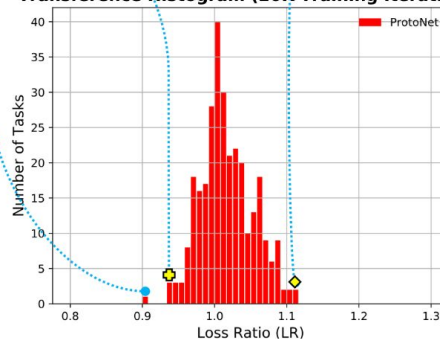
Transference Histogram (10K Training Iterations)



Transference Histogram (15K Training Iterations)



Transference Histogram (20K Training Iterations)



# Towards Reducing Negative Knowledge Transfer

## **Analysis results**

both positive and negative transference

## **How does MTL reduce negative transfer?**

hard parameter sharing and grouping tasks during training

## **Ideal learning episodes in our episodic training scenario**

most compatible with meta-test task

## **Grouping and hard parameter sharing is not possible in episodic training**

1. meta-test tasks are unseen and unknown during meta-training
2. episodic training involves tens of thousands of tasks



**Solution: Task-Aware Layers**

# Proposed Multimodal Meta-Learner

## MMAML General Framework [1]

Task Encoder

$$\mathbf{v}_{\mathcal{T}} = h_{\phi}(\mathcal{S}_{\mathcal{T}})$$

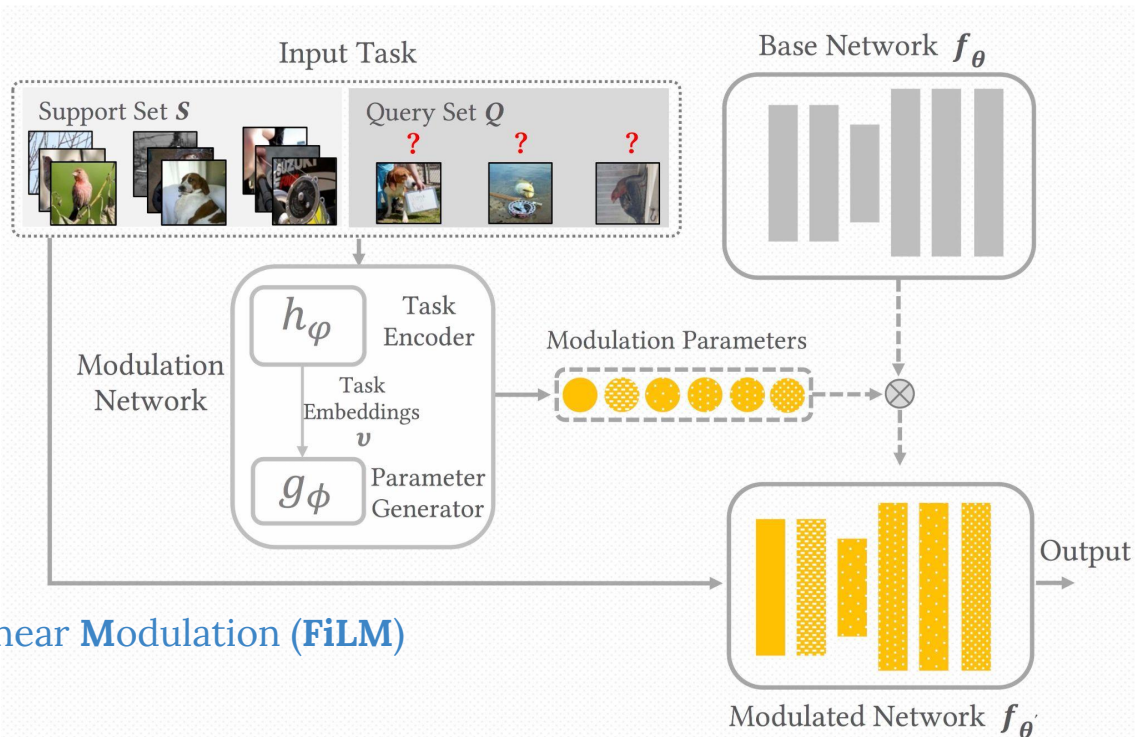
Parameter Generator

$$\omega_{\mathcal{T}} = g_{\phi}(\mathbf{v}_{\mathcal{T}})$$

Modulation Scheme

$$\omega_{\mathcal{T}} = \{\eta_{\mathcal{T}}, \gamma_{\mathcal{T}}\}$$

$$\hat{\mathbf{Y}}_i = \eta_i \mathbf{Y}_i + \gamma_i \quad \text{Feature-wise Linear Modulation (FiLM)}$$



# Proposed Multimodal Meta-Learner

## Limitation of FiLM for Multimodal Meta-Learning

Convolution operator for each channel

$$\mathbf{Y}_i = \mathbf{W}_i * \mathbf{X} + b_i$$

Re-writing FiLM modulation

$$\hat{\mathbf{Y}}_i = (\eta_i \mathbf{W}_i) * \mathbf{X} + (\eta_i b_i + \gamma_i)$$

New interpretation of FiLM modulation:

**Convolution with modulated parameters**

$$\hat{\mathbf{W}}_i = \eta_i \mathbf{W}_i$$

**modulated kernel**

$$\hat{b}_i = \eta_i b_i + \gamma_i$$

**modulated bias**

# Proposed Multimodal Meta-Learner

## Proposed **Kernel Modulation (KML)**

Modulate every parameter within network

$$\hat{\mathbf{W}}_{\mathcal{T}}^l = \mathbf{W}^l \odot (\mathbf{J} + \mathbf{M}^l(\mathbf{v}_{\mathcal{T}}, \phi))$$

$$\hat{\mathbf{b}}_{\mathcal{T}}^l = \mathbf{b}^l + \Delta \mathbf{b}^l(\mathbf{v}_{\mathcal{T}}, \phi)$$

Modulated parameters for whole network

$$\hat{\theta}_{\mathcal{T}} = \{\hat{\mathbf{W}}_{\mathcal{T}}^1, \dots, \hat{\mathbf{W}}_{\mathcal{T}}^L, \hat{\mathbf{b}}_{\mathcal{T}}^1, \dots, \hat{\mathbf{b}}_{\mathcal{T}}^L\}$$

Generator Design for KML

use **MLP** to be computationally efficient

Limitation: large number of required parameters in MLP

# Proposed Multimodal Meta-Learner

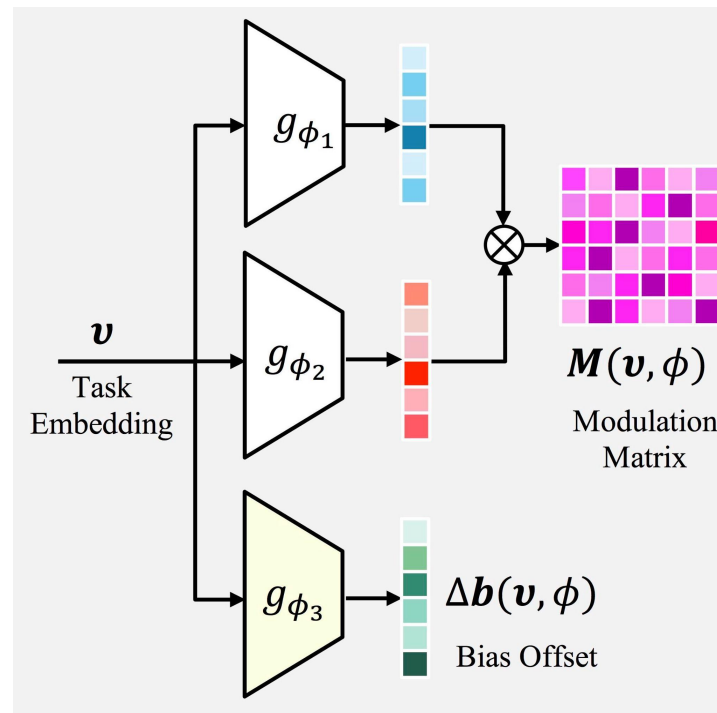
## Proposed Simplified Structure for Parameter Generator

Use multiple smaller MLPs instead of a large one

$$\mathbf{M}^l(\mathbf{v}_T, \phi) = \mathbf{g}_{\phi_1}^l(\mathbf{v}_T) \otimes \mathbf{g}_{\phi_2}^l(\mathbf{v}_T)$$

$$\Delta \mathbf{b}^l(\mathbf{v}_T, \phi) = \mathbf{g}_{\phi_3}^l(\mathbf{v}_T)$$

- Reduces the required parameters by a factor of **152**
- Improves generalization performance



# Experimental Results

## Datasets

### mini-ImageNet:

natural objects

### FC100:

natural objects with  
lower resolution

### Omniglot:

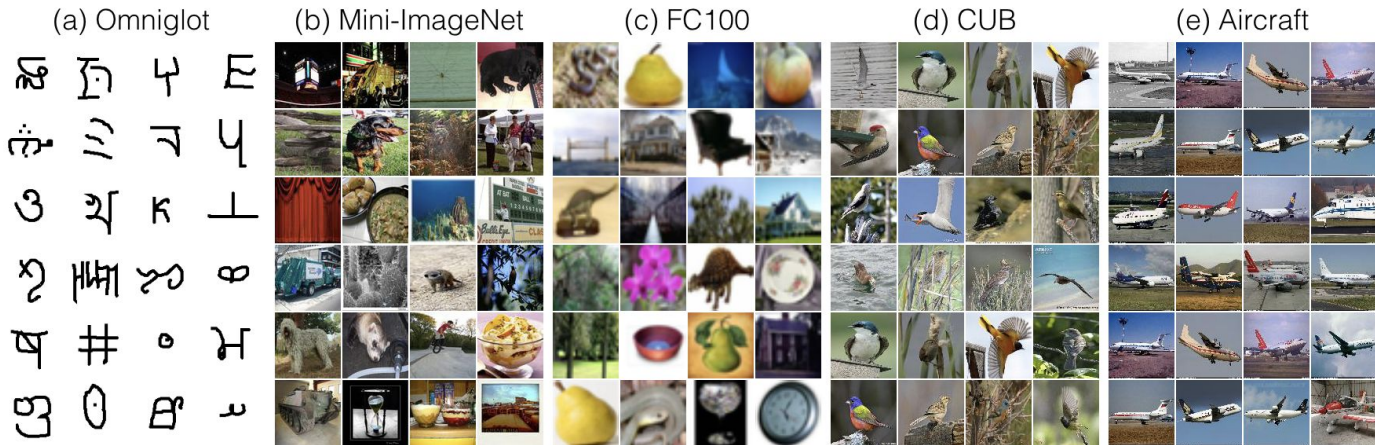
handwritten characters

### CUB:

fine-grained bird  
classification

### Aircraft:

fine-grained aircraft  
classification



Datasets used to generate meta-dataset for multimodal meta-learning following MMAML

**2Mode<sup>†</sup>**: mini-ImageNet and FC100; **2Mode**: mini-ImageNet and Omniglot

**3Mode**: mini-ImageNet, Omniglot and FC100; **5Mode**: mini-ImageNet, Omniglot, FC100, AIRCRAFT and CUB

# Experimental Results

## Multimodal Few-Shot Classification Results

Setup		Method			
		MAML [19]*	Multi-MAML	MMAML [1]*	MMAML+KML (ours)
<b>2Mode<sup>†</sup></b>	1-shot	40.53±68%	39.27±0.76%	39.11±0.62%	<b>40.73±0.66%</b>
	5-shot	<b>54.11±0.63%</b>	53.51±0.72%	52.02±0.63%	53.72±0.60%
<b>2Mode</b>	1-shot	65.18±0.61%	66.77±0.68%	67.67±0.63%	<b>68.01±0.59%</b>
	5-shot	74.18±0.57%	73.07±0.61%	73.52±0.71%	<b>77.02±0.66%</b>
<b>3 Mode</b>	1-shot	54.40±0.56%	56.01±0.66%	57.35±0.61%	<b>57.68±0.59%</b>
	5-shot	66.51±0.54%	65.92±0.62%	64.21±0.57%	<b>67.12±0.55%</b>
<b>5Mode</b>	1-shot	47.19±0.49%	48.33±0.58%	49.53±0.50%	<b>50.31±0.49%</b>
	5-shot	58.13±0.48%	59.20±0.52%	58.89±0.47%	<b>60.51±0.47%</b>



# Experimental Results

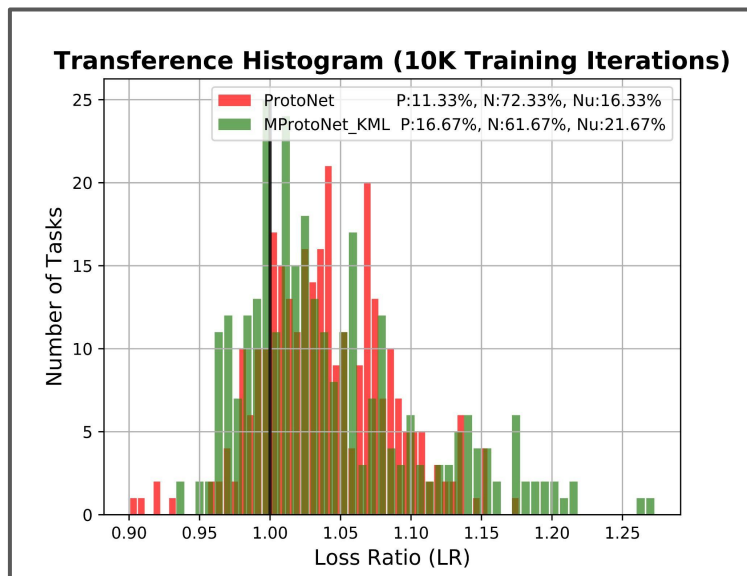
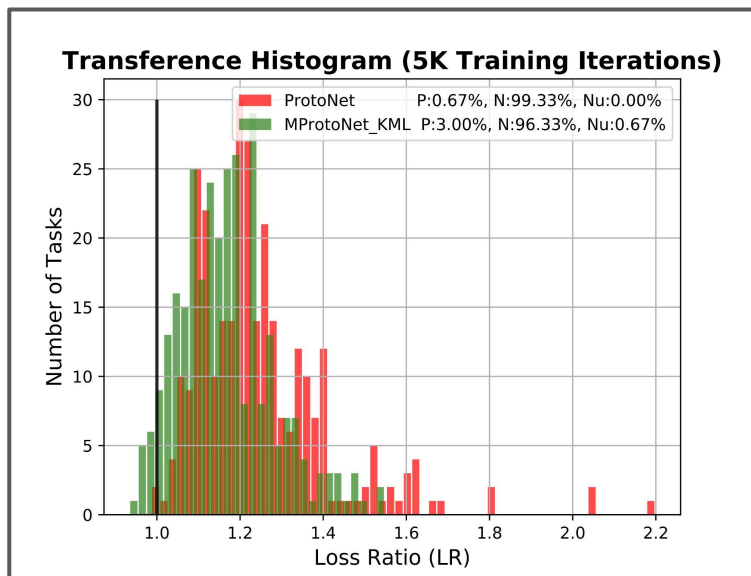
## Multimodal Few-Shot Classification Results

Setup		Method			
		ProtoNet [11]**	Multi-ProtoNet	MProtoNet [1]**	MProtoNet+KML (ours)
<b>2Mode<sup>†</sup></b>	1-shot	43.05±0.58%	43.42±0.56%	43.57±0.59%	<b>44.40±0.65%</b>
	5-shot	57.70±0.59%	56.73±0.64%	56.03±0.64%	<b>59.31±0.62%</b>
<b>2Mode</b>	1-shot	69.55±0.54%	70.17±0.61%	70.60±0.56%	<b>73.69±0.52%</b>
	5-shot	75.12±0.41%	75.33±0.46%	75.72±0.47%	<b>79.82±0.40%</b>
<b>3 Mode</b>	1-shot	58.14±0.49%	59.89±0.50%	59.62±0.54%	<b>62.08±0.54%</b>
	5-shot	66.84±0.44%	67.03±0.44%	67.51±0.47%	<b>70.03±0.43%</b>
<b>5Mode</b>	1-shot	49.31±0.53%	50.69±0.57%	51.75±0.52%	<b>56.72±0.46%</b>
	5-shot	58.91±0.51%	59.88±0.54%	59.95±0.42%	<b>64.91±0.38%</b>

# Experimental Results

## Transference Results (mini-ImageNet $\Rightarrow$ FC100)

$$LR_{i \rightarrow j} = \frac{\mathcal{L}_{\mathcal{T}_j}(Q_j; \theta_i^{t+1}, \mathcal{S}_j)}{\mathcal{L}_{\mathcal{T}_j}(Q_j; \theta^t, \mathcal{S}_j)} \quad \longrightarrow \quad LR < 1: \text{Positive Knowledge Transfer}$$



# Summary

## **Research gaps in multimodal meta-learning**

- How can we measure interaction between few-shot tasks?
- How can we improve the generalization performance?

## **Proposed work**

- Adapt transference idea from MTL to quantify interaction between few-shot tasks
- A new interpretation of FiLM scheme
- Kernel modulation to improve generalization

## **Experimental Results**

- Significant improvement over previous state-of-the art in both micro and macro-level

Transference analysis and proposed KML can be extended to **conventional meta-learning** (Supplementary Material).