



Tencent  
AI Lab



PennState

# Hierarchically Structured Meta-learning

Huaxiu Yao<sup>1,2</sup>, Ying Wei<sup>2</sup>, Junzhou Huang<sup>1</sup>, Zhenhui Li<sup>2</sup>

<sup>1</sup>Pennsylvania State University

<sup>2</sup>Tencent AI Lab

**Oral: Thu Jun 13th 09:35 -- 09:40 AM @ Room 103**

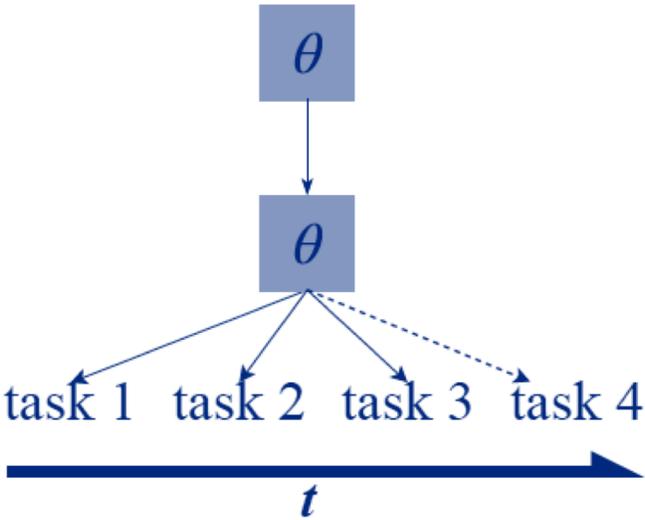
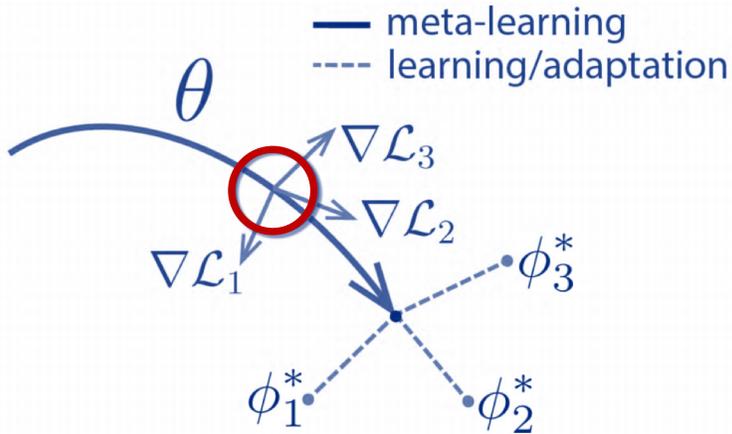
**Poster: Thu Jun 13th 06:30 -- 09:00 PM @ Pacific Ballroom #183**

## Is global initialization enough?

$$\min_{\theta} \sum_{\text{task } i} \mathcal{L}_{\text{test}}^i(\theta - \alpha \nabla_{\theta} \mathcal{L}_{\text{train}}^i(\theta))$$

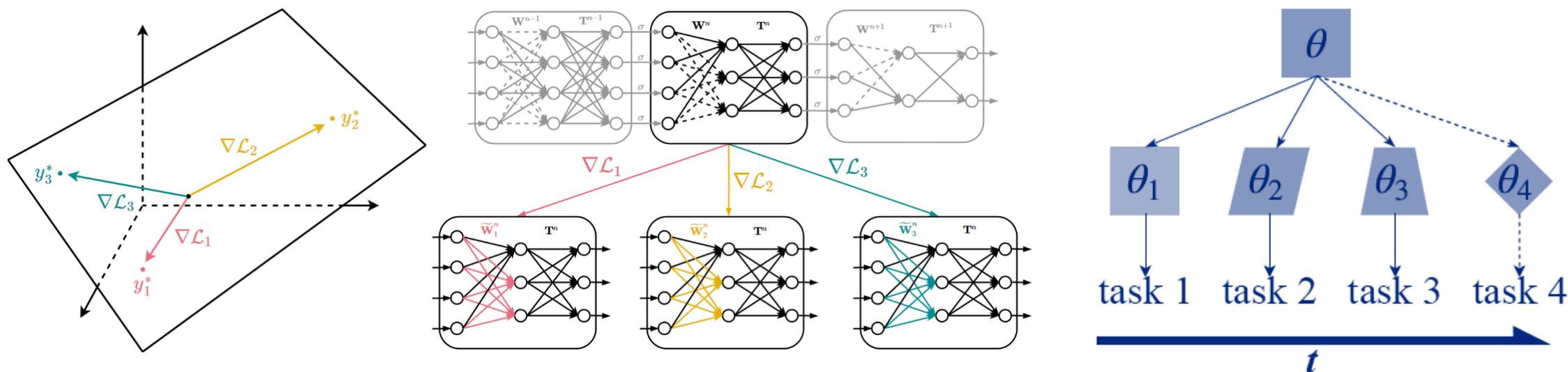
$\theta$  parameter vector being meta-learned

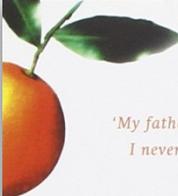
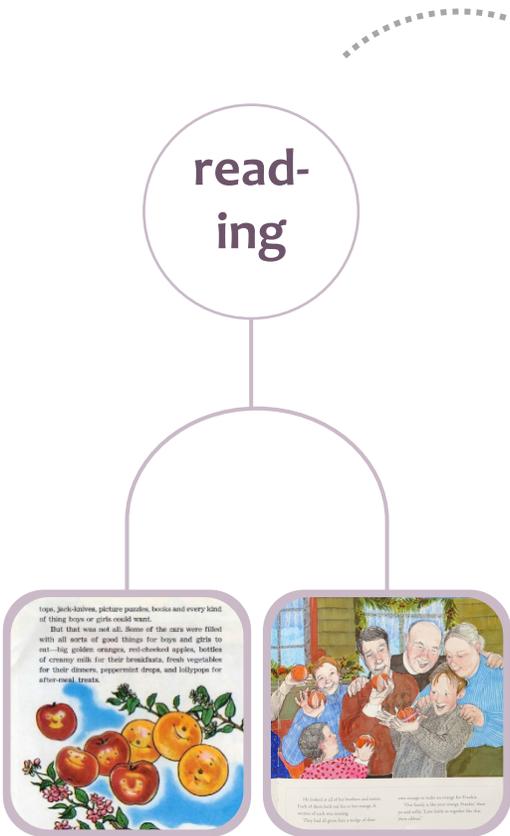
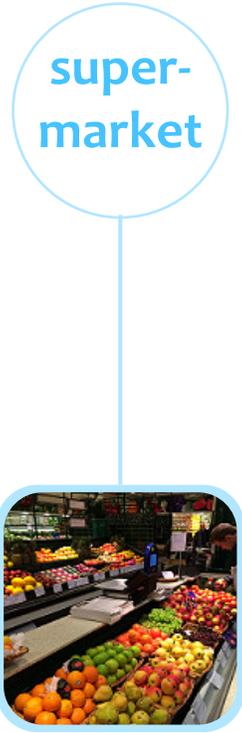
$\phi_i^*$  optimal parameter vector for task  $i$



[1] Finn, Chelsea, Pieter Abbeel, and Sergey Levine. "Model-agnostic meta-learning for fast adaptation of deep networks." Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR. org, 2017. [http://people.eecs.berkeley.edu/~cbfinn/\\_files/metalearning\\_frontiers\\_2018\\_small.pdf](http://people.eecs.berkeley.edu/~cbfinn/_files/metalearning_frontiers_2018_small.pdf)

## Should the initialization be tailored to each task?





*'My father died eleven years ago. I was only four then. I never thought I'd hear from him again, but now we're writing a book together'*

To Georg Røed, his father is no more than a shadow, a distant memory. But then one day his grandmother discovers some pages stuffed into the lining of an old red pushchair. The pages are a letter to Georg, written just before his father died, and a story, 'The Orange Girl'.

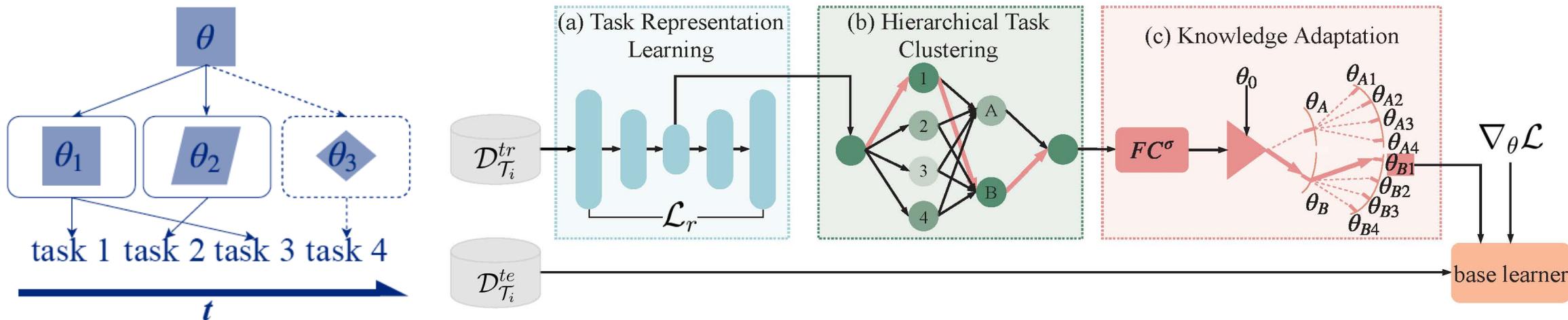
'The Orange Girl' begins when his father boarded a tram and was captivated by a girl standing in the aisle, clutching a huge paper bag of luscious-looking oranges.

[3] Gershman, Samuel J., David M. Blei, and Yael Niv. "Context, learning, and extinction." *Psychological review* 117.1 (2010): 197.  
[4] Gershman, Samuel J., et al. "Statistical computations underlying the dynamics of memory updating." *PLoS computational biology* 10.11 (2014): e1003939.



## Balance between generalization and customization

- Organize tasks by hierarchical clustering
- Adapt the global initialization to each cluster of tasks



## Overall optimization problem

$$\min_{\Theta} \sum_{i=1}^{N_t} \mathcal{L}(f_{\theta_{0i} - \alpha \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}_{T_i}^{tr})}, \mathcal{D}_{T_i}^{te}) + \xi \mathcal{L}_r(\mathcal{D}_{T_i}^{tr});$$

## Extension to continual adaptation

- Incrementally increase the clusters as tasks sequentially arrive.
- Criterion for adding a cluster—evaluate the average loss over Q epochs

$$\bar{\mathcal{L}}_{new} > \mu \bar{\mathcal{L}}_{old}$$

- For task  $\mathcal{T}_i \sim \mathcal{E}$ , training and testing samples are i.i.d. drawn from  $\mathcal{S}_i$
- The initialization of HSML (K clusters) can be represented as  $\theta_{0t} = \sum_{k=1}^K \widehat{\mathbf{B}}_k \theta_0$
- According to [5], the assumptions are  $\mathcal{L} \in [0, 1]$  is  $\eta$ -smooth and has a  $\rho$ -Lipschitz Hessian, step size at the  $u$ -step  $\alpha_u = c/u$  satisfying  $c \leq \min\left\{\frac{1}{\eta}, \frac{1}{4(2\eta \ln U)^2}\right\}$  with total steps  $U = n^{tr}$ .
- The generalization of base learner  $f_{\theta_{\mathcal{T}_i}}$  is bounded by  $\epsilon(\mathcal{S}_i, \theta_0)$ , where

$$\epsilon(\mathcal{S}_i, \theta_0) = \mathcal{O} \left( \left(1 + \frac{1}{c\widehat{\gamma}^-}\right) \widehat{R}_{\mathcal{D}_{\mathcal{T}_t}^{tr}}(\theta_{0t})^{\frac{c\widehat{\gamma}^+}{1+c\widehat{\gamma}^+}} \frac{1}{(n^{tr})^{\frac{1}{1+c\widehat{\gamma}^+}}} \right)$$

- MAML can be regarded as a special case of HSML, i.e.,  $\forall k, \widehat{\mathbf{B}}_k = \mathbf{I}$
- After proving  $\exists \{\widehat{\mathbf{B}}_k\}_{k=1}^K$ , s.t.,  $\widehat{R}_{\mathcal{D}_{\mathcal{T}_t}^{tr}}(\theta_{0t}) \leq \widehat{R}_{\mathcal{D}_{\mathcal{T}_t}^{tr}}(\theta_0)$ , we conclude that HSML achieves a tighter generalization bound than MAML

## Data

- 4 sync family functions—Sin, Line, Cubic, Quadratic
- K-shot: K samples are used as training (each task)

## Base model

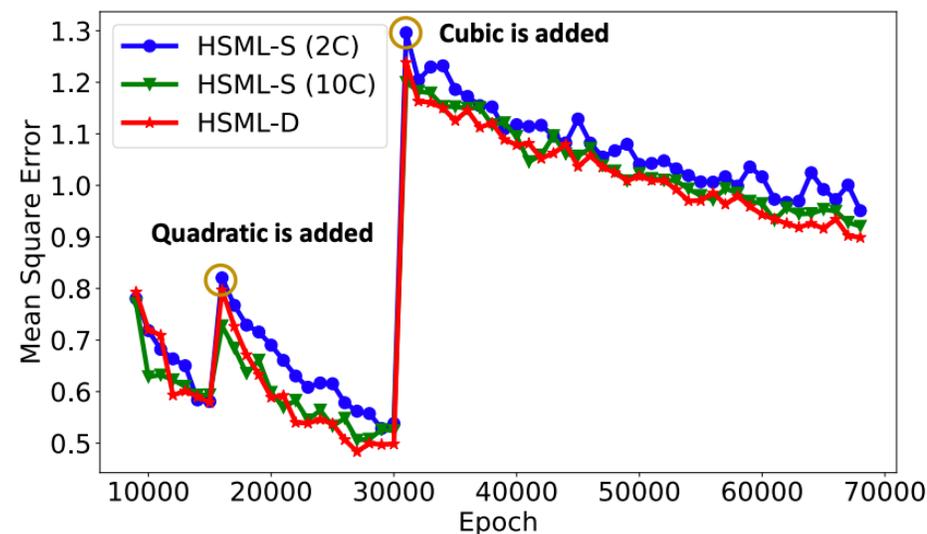
- 2 layers FC with 40 neurons each

## Quantitative results

- Comparison on regression MSEs

Method	5-shot	10-shot
Global shared (MAML)	$2.205 \pm 0.121$	$0.761 \pm 0.068$
Task-specific (MUMOMAML[6])	$1.096 \pm 0.085$	$0.256 \pm 0.028$
<b>Our method (HSML)</b>	<b><math>0.856 \pm 0.073</math></b>	<b><math>0.161 \pm 0.021</math></b>

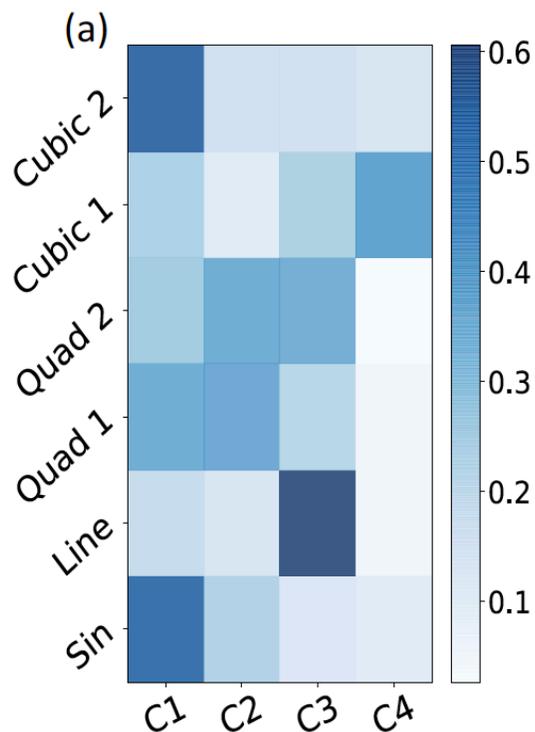
- Comparison in the continual adaptation scenario



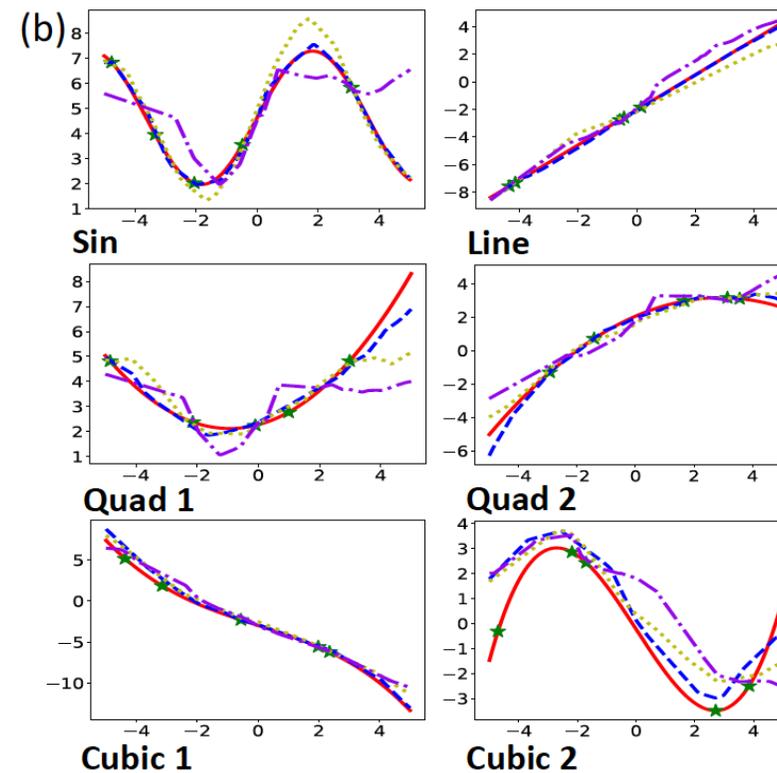
Model	HSML-S (2C)	HSML-S (10C)	HSML-D
MSE $\pm$ 95% CI	$0.933 \pm 0.074$	$0.889 \pm 0.071$	<b><math>0.869 \pm 0.072</math></b>

## Qualitative results

- Cluster assignment interpretation



- Regression results



— Ground Truth ★ Selected Point - - MAML ··· MUMOMAML - - HSML

## Data

- 4 image classification datasets—Bird, Texture, Aircraft, Fungi
- 5-way, 1-shot

## Base model

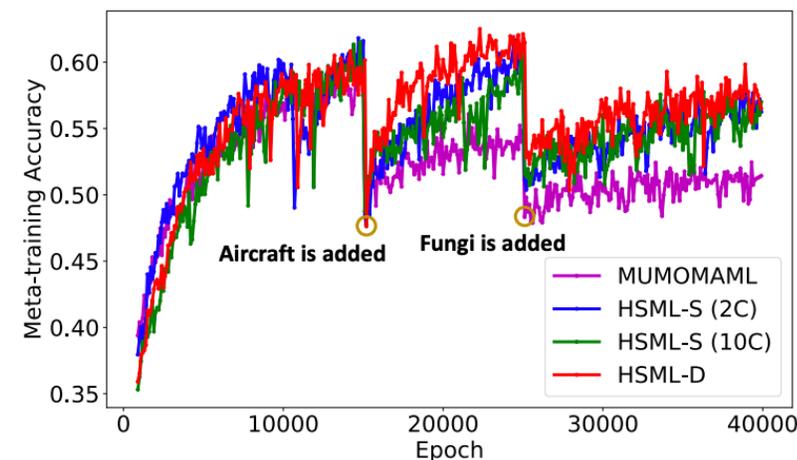
- a convolutional network with 4 convolution blocks

## Quantitative results

- Comparison on accuracy

Method	Bird	Texture	Aircraft	Fungi
Global shared (MAML)	53.94 %	31.66 %	51.37 %	42.12 %
Task-specific (MUMOMAML[6])	56.82 %	33.81 %	53.14 %	42.22 %
<b>Our method (HSML)</b>	<b>60.98 %</b>	<b>35.01%</b>	<b>57.38 %</b>	<b>44.02 %</b>

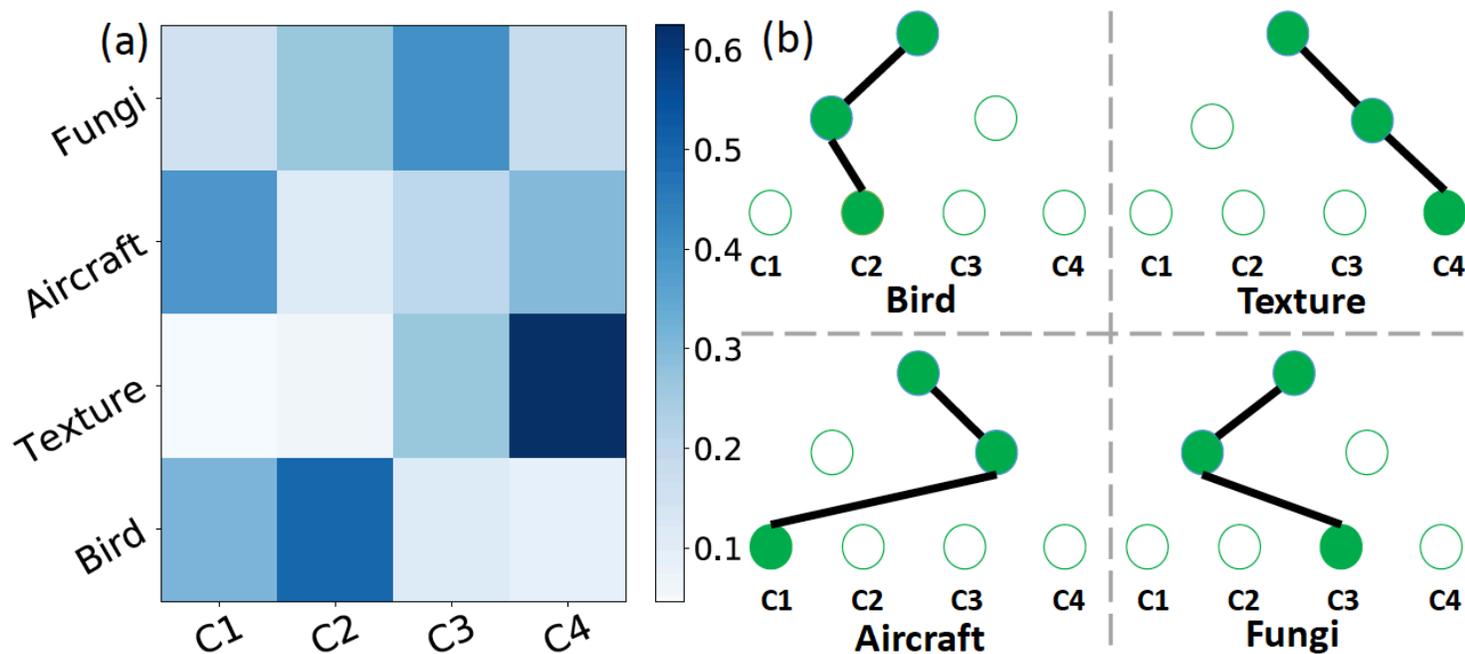
- Comparison in the continual adaptation scenario



Model	Bird	Texture	Aircraft	Fungi
MUMOMAML	56.66%	33.68%	45.73%	40.38%
HSML-S (2C)	60.77%	33.41%	51.28%	40.78%
HSML-S (10C)	59.16%	34.48%	52.30%	40.56%
HSML-D	<b>61.16%</b>	<b>34.53%</b>	<b>54.50%</b>	<b>41.66%</b>

## Qualitative results

- Cluster assignment interpretation



- HSML simultaneously customizes task knowledge and preserves knowledge generalization via the hierarchical clustering structure.
- Experiments demonstrate the effectiveness and interpretability of HSML in both toy regression and few-shot classification problems.



PennState

# THANK YOU

**Oral: Thu Jun 13th 09:35 -- 09:40 AM @ Room 103**

**Poster: Thu Jun 13th 06:30 -- 09:00 PM @ Pacific Ballroom #183**