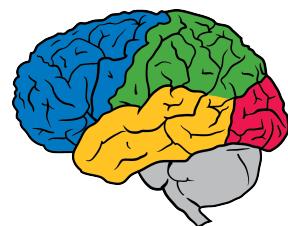
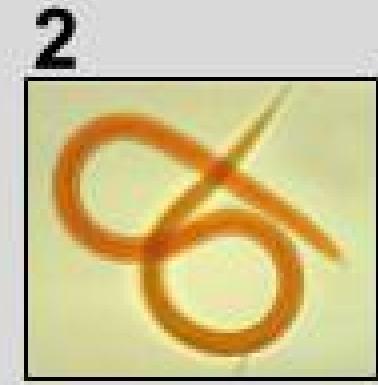


# THOUGHTS ON PROGRESS MADE AND CHALLENGES AHEAD IN FEW-SHOT LEARNING



Hugo Larochelle  
Google Brain



$D_{train}$



$D_{test}$

**People are  
good at it**



# Human-level concept learning through probabilistic program induction

Brenden M. Lake,<sup>1,\*</sup> Ruslan Salakhutdinov,<sup>2</sup> Joshua B. Tenenbaum<sup>3</sup>



**Machines are  
getting  
better at it**

പ	എ	ബ	ഡ	
കു	ന	നു	ബു	
ഭ	പ	ണ	തേ	ദ
ന	മ	ല	ക	ബു

## LEARNING

0 EXAMPLES

CONFIDENCE



TRAIN GREEN



0 EXAMPLES

CONFIDENCE



TRAIN PURPLE

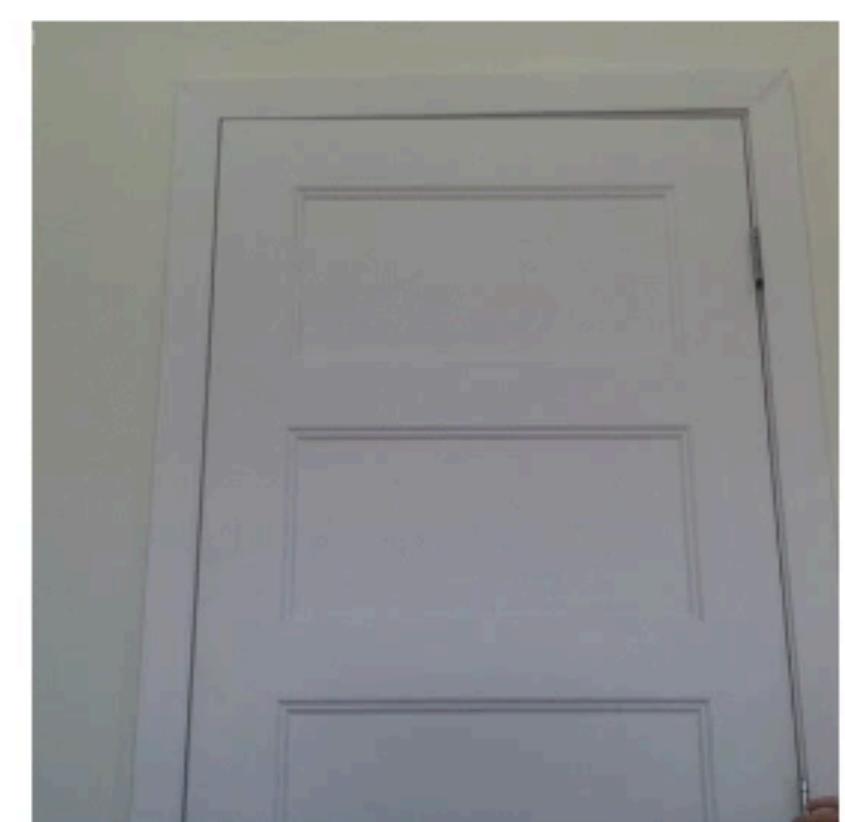
0 EXAMPLES

CONFIDENCE



TRAIN ORANGE

## INPUT



## OUTPUT

GIF

Sound

Speech



## LEARNING

0 EXAMPLES

CONFIDENCE



TRAIN GREEN



0 EXAMPLES

CONFIDENCE



TRAIN PURPLE

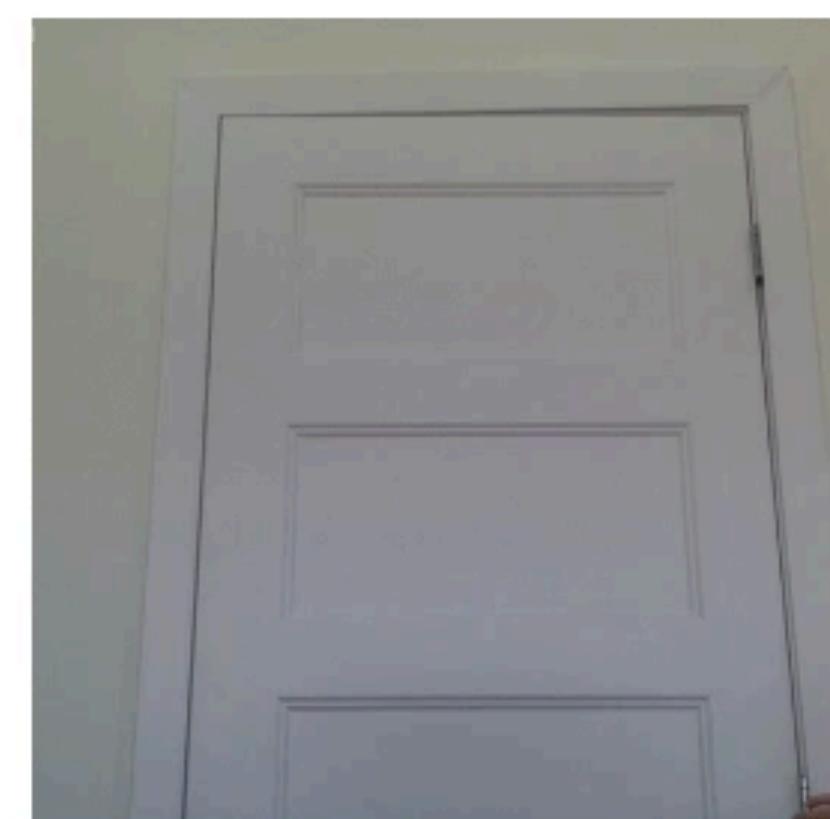
0 EXAMPLES

CONFIDENCE



TRAIN ORANGE

## INPUT

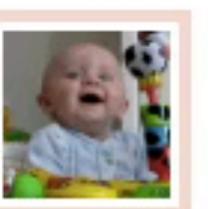


## OUTPUT

GIF

Sound

Speech



# RELATED WORK: ONE-SHOT LEARNING

- One-shot learning has been studied before
  - ▶ One-Shot learning of object categories (2006)  
*Fei-Fei Li, Rob Fergus and Pietro Perona*
  - ▶ Knowledge transfer in learning to recognize visual objects classes (2004)  
*Fei-Fei Li*
  - ▶ Object classification from a single example utilizing class relevance pseudo-metrics (2004)  
*Michael Fink*
  - ▶ Cross-generalization: learning novel classes from a single example by feature replacement (2005)  
*Evgeniy Bart and Shimon Ullman*
- These largely relied on hand-engineered features and algorithms
  - ▶ with recent progress in **end-to-end** deep learning, we hope to jointly learn a **representation** and **algorithm** better suited for few-shot learning

# META-LEARNING



# META-LEARNING

$D_{train}$  |  $D_{test}$

II  
episode

**Meta-  
Train**  
 $\mathcal{D}_{meta-train}$



$D_{train}$

$D_{test}$

$D_{train}$

$D_{test}$

⋮

⋮

**Meta-  
Test**  
 $\mathcal{D}_{meta-test}$



$D_{train}$

$D_{test}$

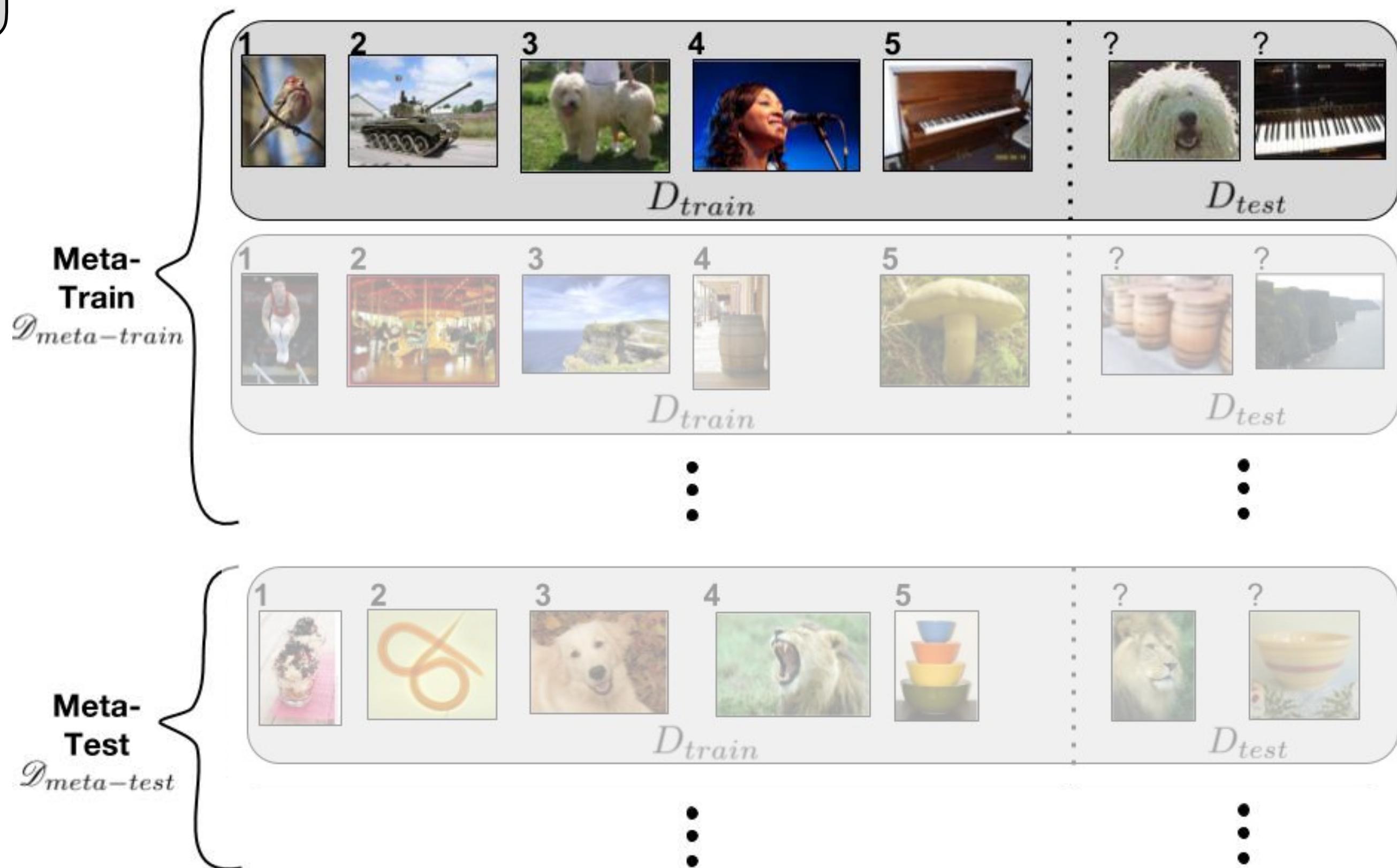
⋮

⋮

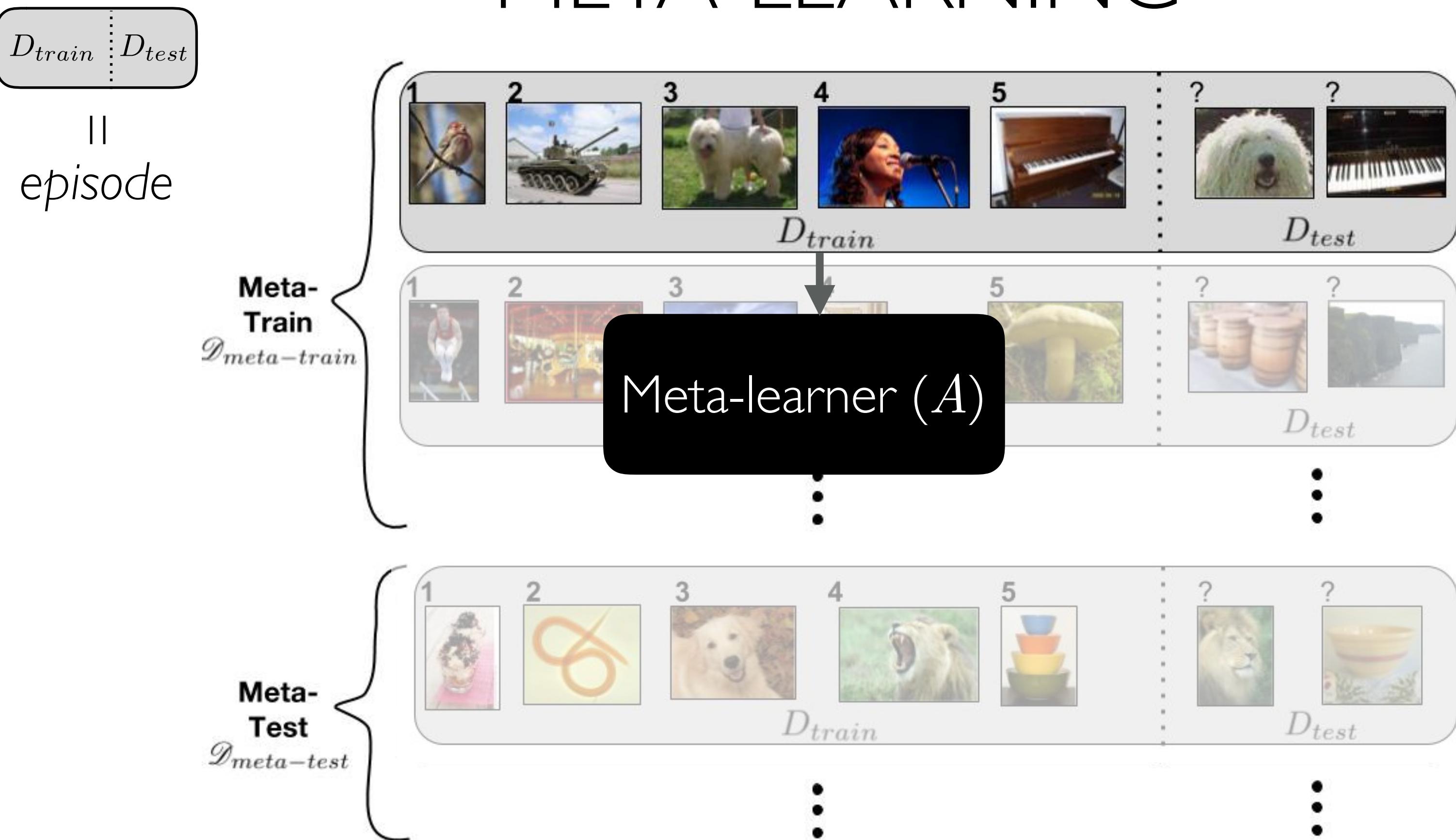
# META-LEARNING

$D_{train}$  |  $D_{test}$

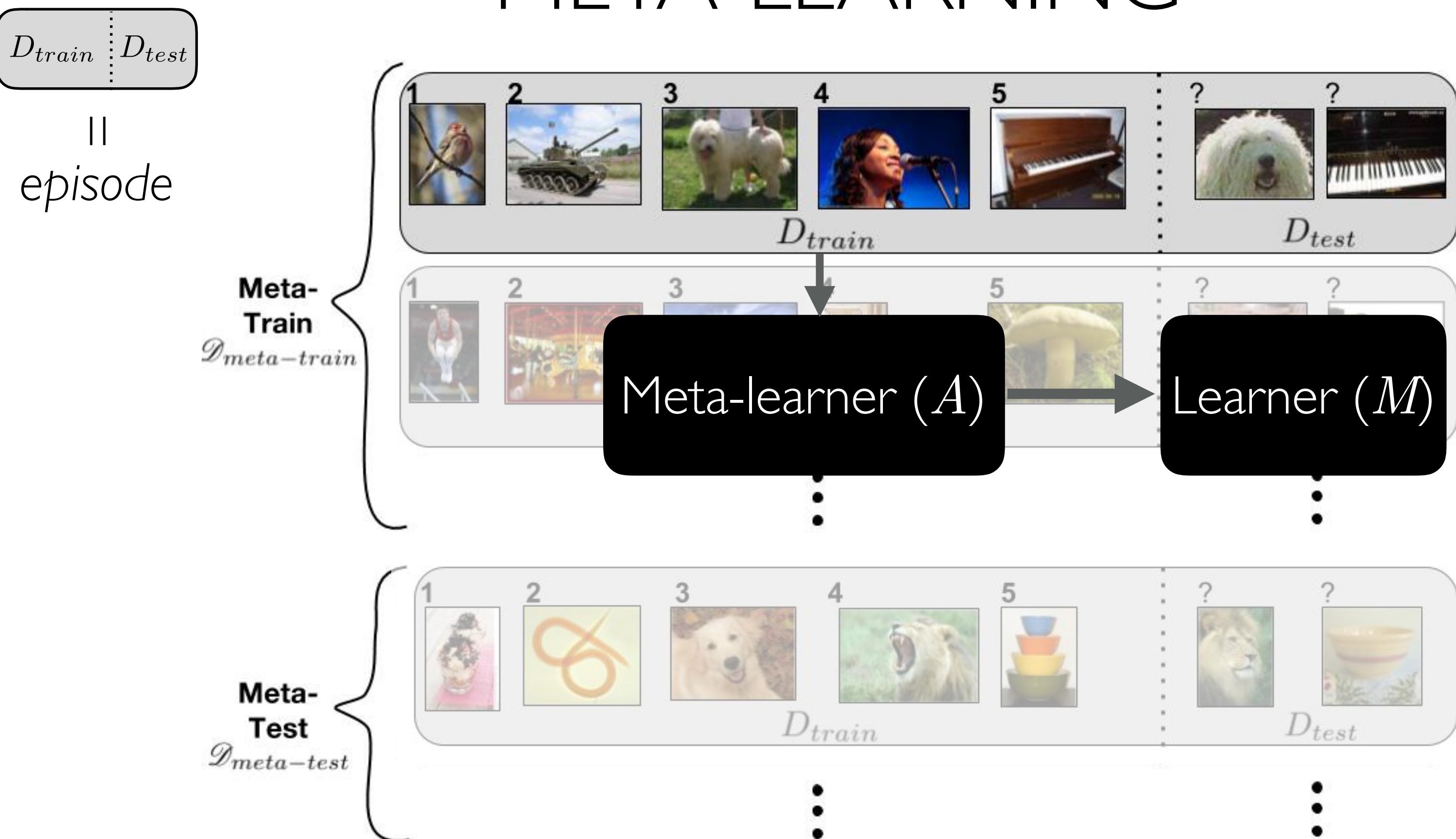
II  
episode



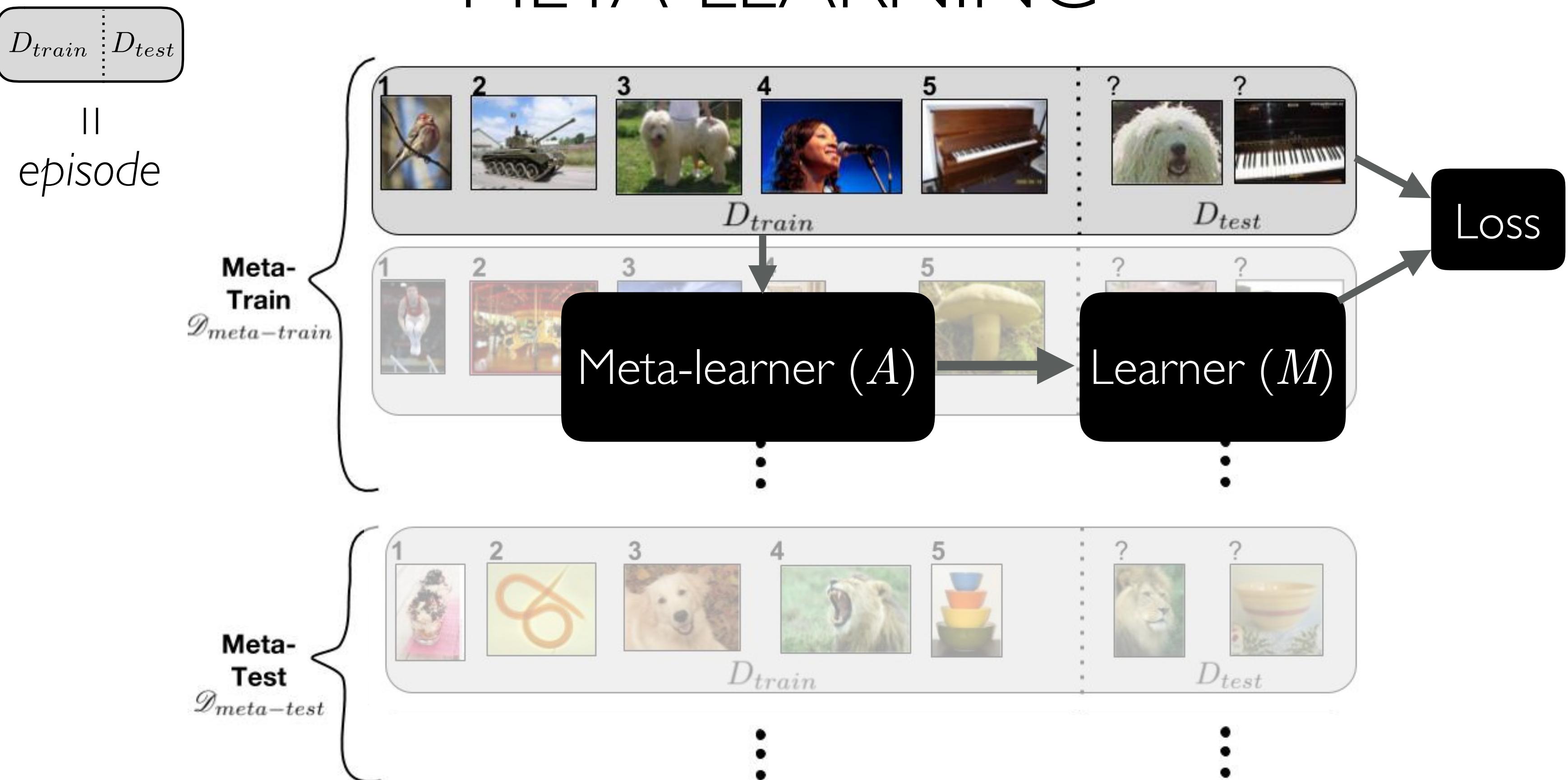
# META-LEARNING



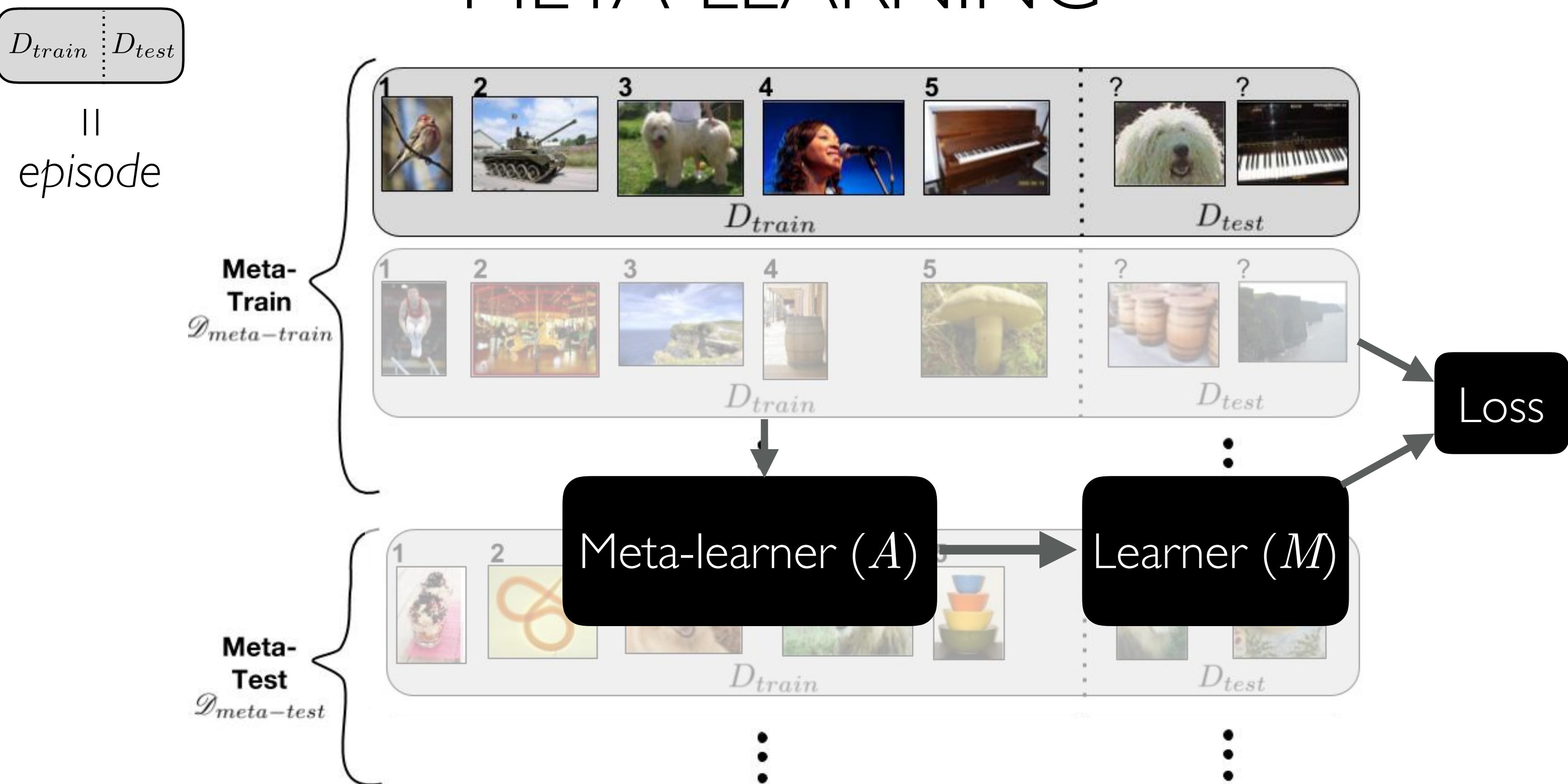
# META-LEARNING



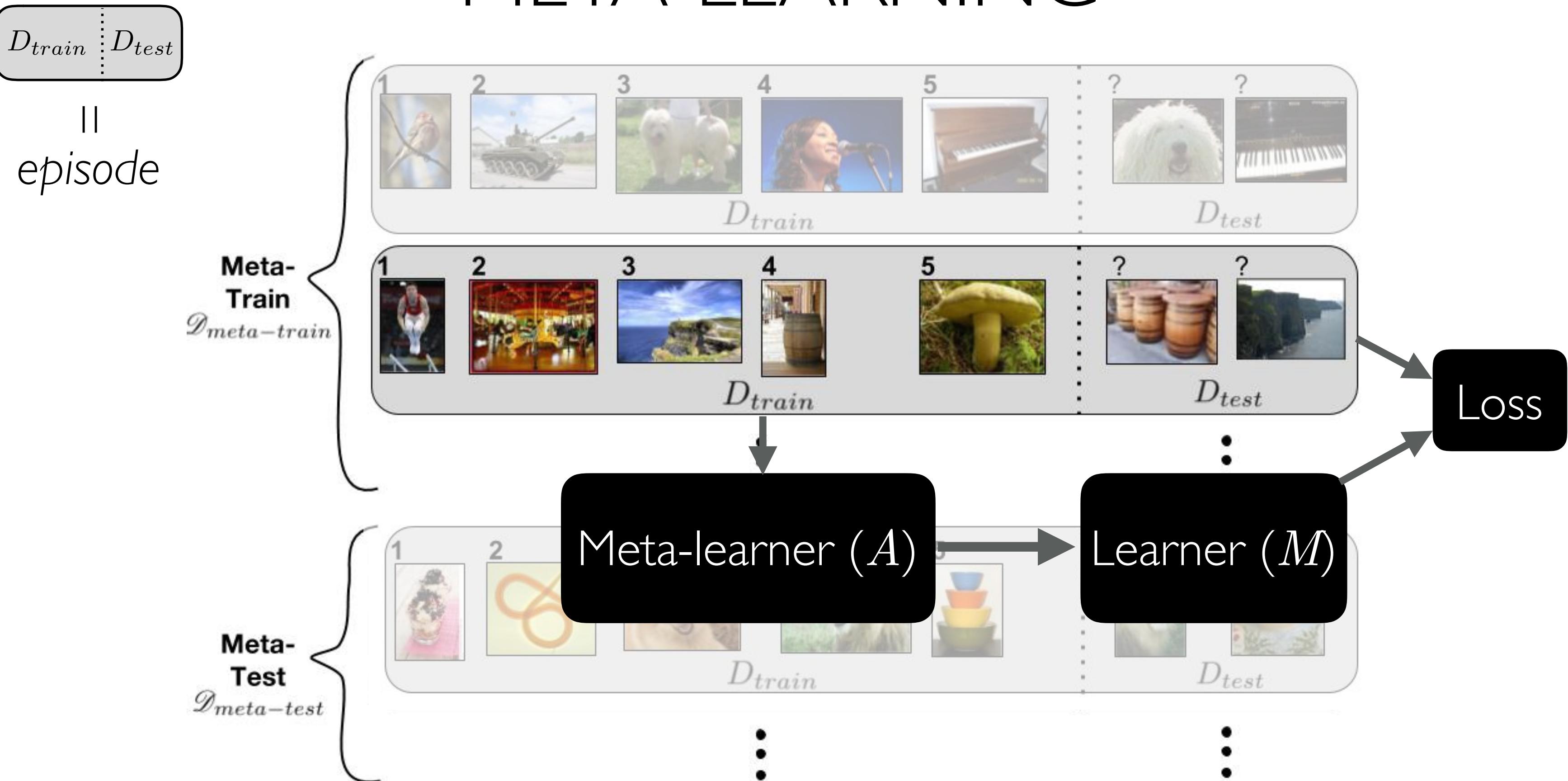
# META-LEARNING



# META-LEARNING



# META-LEARNING



If you don't evaluate on never-seen problems/datasets...

If you don't evaluate on never-seen problems/datasets...

... it's not meta-learning!

# LEARNING PROBLEM STATEMENT

- Assuming a probabilistic model  $M$  over labels, the cost per episode can be written as

$$C(D_{train}, D_{test}) = \frac{1}{|D_{test}|} \sum_{\substack{(\mathbf{x}_t, y_t) \\ \in D_{test}}} -\log p(y_t | \mathbf{x}_t, D_{train})$$

- Here  $p(y|\mathbf{x}, D_{train})$  jointly represents the meta-learner  $A$  (which processes  $D_{train}$ ) and the learner  $M$  (which processes  $\mathbf{x}$ )

# CHOOSING A META-LEARNER

- How to parametrize learning algorithms (meta-learners  $p(y|\mathbf{x}, D_{train})$  )?
- Two approaches to defining a meta-learner
  - ▶ Take inspiration from a known learning algorithm
    - kNN/kernel machine: Matching networks (Vinyals et al. 2016)
    - Gaussian classifier: Prototypical Networks (Snell et al. 2017)
    - Gradient Descent: Meta-Learner LSTM (Ravi & Larochelle, 2017) , MAML (Finn et al. 2017)
  - ▶ Derive it from a black box neural network
    - SNAIL (Mishra et al. 2018)

# CHOOSING A META-LEARNER

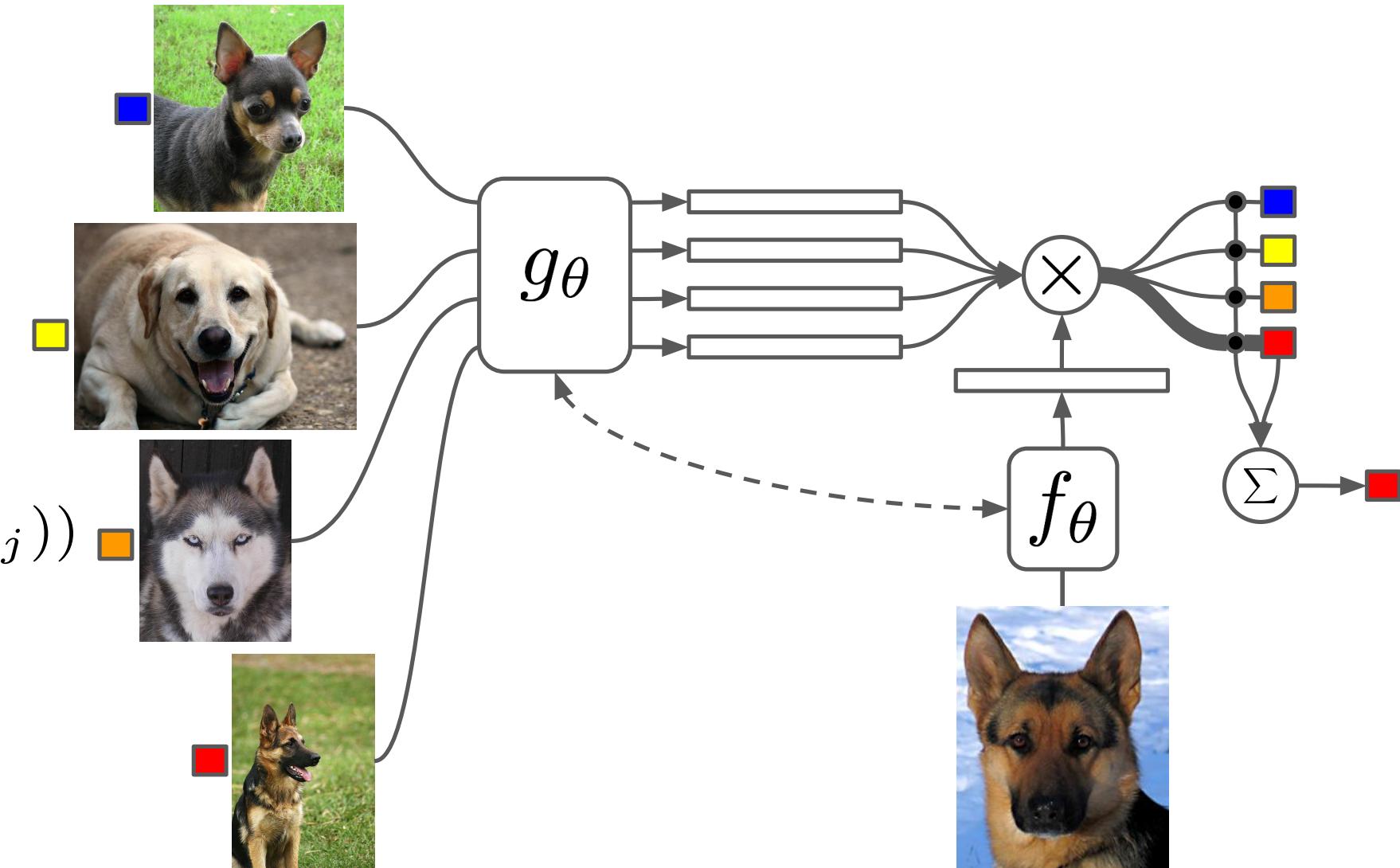
- How to parametrize learning algorithms (meta-learners  $p(y|\mathbf{x}, D_{train})$  )?
- Two approaches to defining a meta-learner
  - ▶ **Take inspiration from a known learning algorithm**
    - kNN/kernel machine: Matching networks (Vinyals et al. 2016)
    - Gaussian classifier: Prototypical Networks (Snell et al. 2017)
    - Gradient Descent: Meta-Learner LSTM (Ravi & Larochelle, 2017) , MAML (Finn et al. 2017)
  - ▶ Derive it from a black box neural network
    - SNAIL (Mishra et al. 2018)

# MATCHING NETWORKS

- Training a “**pattern matcher**” (kNN/kernel machine)

$$\hat{y} = \sum_{i=1}^k a(\hat{x}, x_i) y_i$$

$$a(\hat{x}, x_i) = e^{c(f(\hat{x}), g(x_i))} / \sum_{j=1}^k e^{c(f(\hat{x}), g(x_j))}$$



- Matching networks for one shot learning (2016)

Oriol Vinyals, Charles Blundell, Timothy P. Lillicrap, Koray Kavukcuoglu, and Daan Wierstra

# PROTOTYPICAL NETWORKS

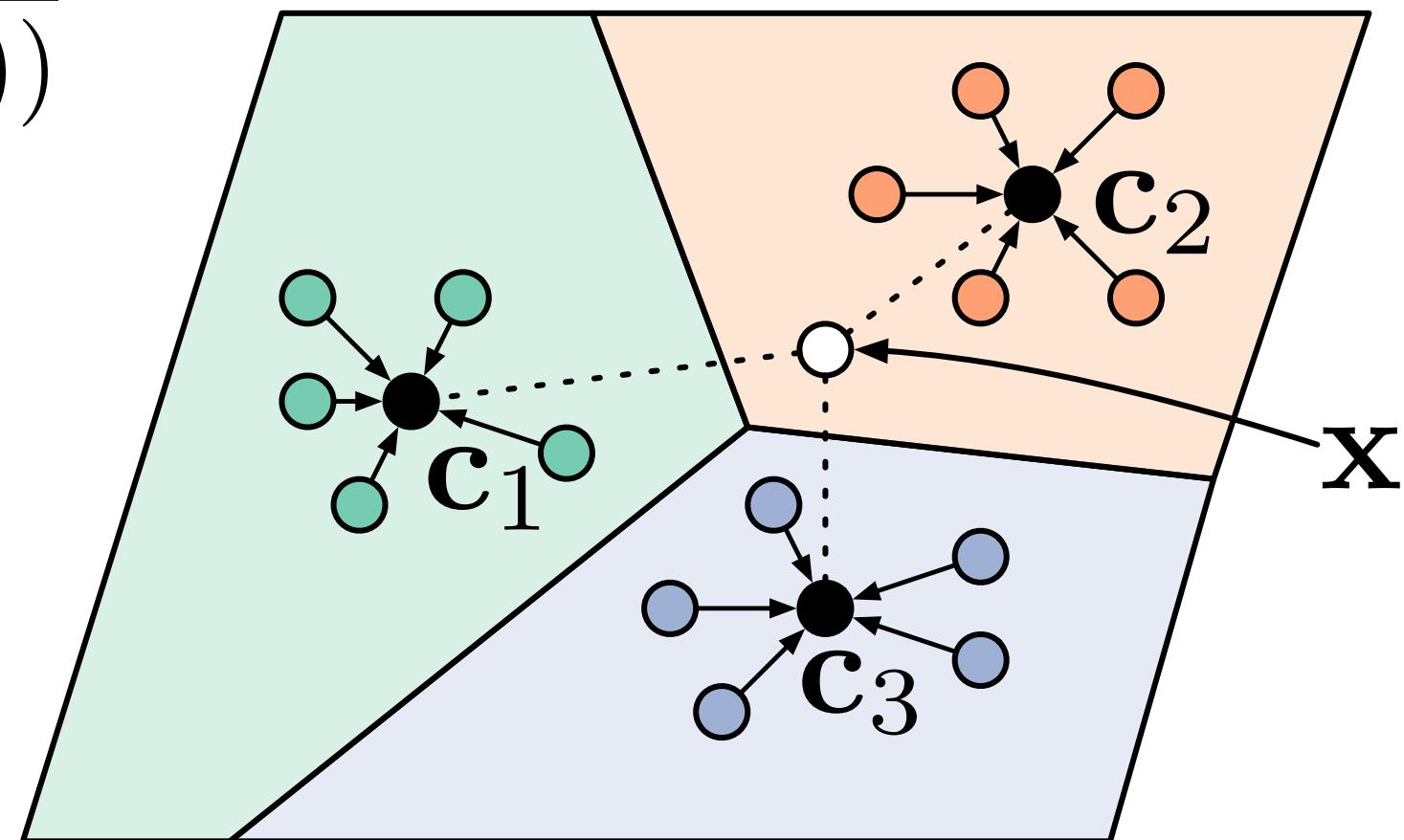
- Training a “**prototype extractor**” (Gaussian classifier)

$$p_{\phi}(y = k \mid \mathbf{x}) = \frac{\exp(-d(f_{\phi}(\mathbf{x}), \mathbf{c}_k))}{\sum_{k'} \exp(-d(f_{\phi}(\mathbf{x}), \mathbf{c}_{k'}))}$$

$$\mathbf{c}_k = \frac{1}{|S_k|} \sum_{(\mathbf{x}_i, y_i) \in S_k} f_{\phi}(\mathbf{x}_i)$$

$$S_k = \{(\mathbf{x}_i, y_i) \mid y_i = k, (\mathbf{x}_i, y_i) \in D_{train}\}$$

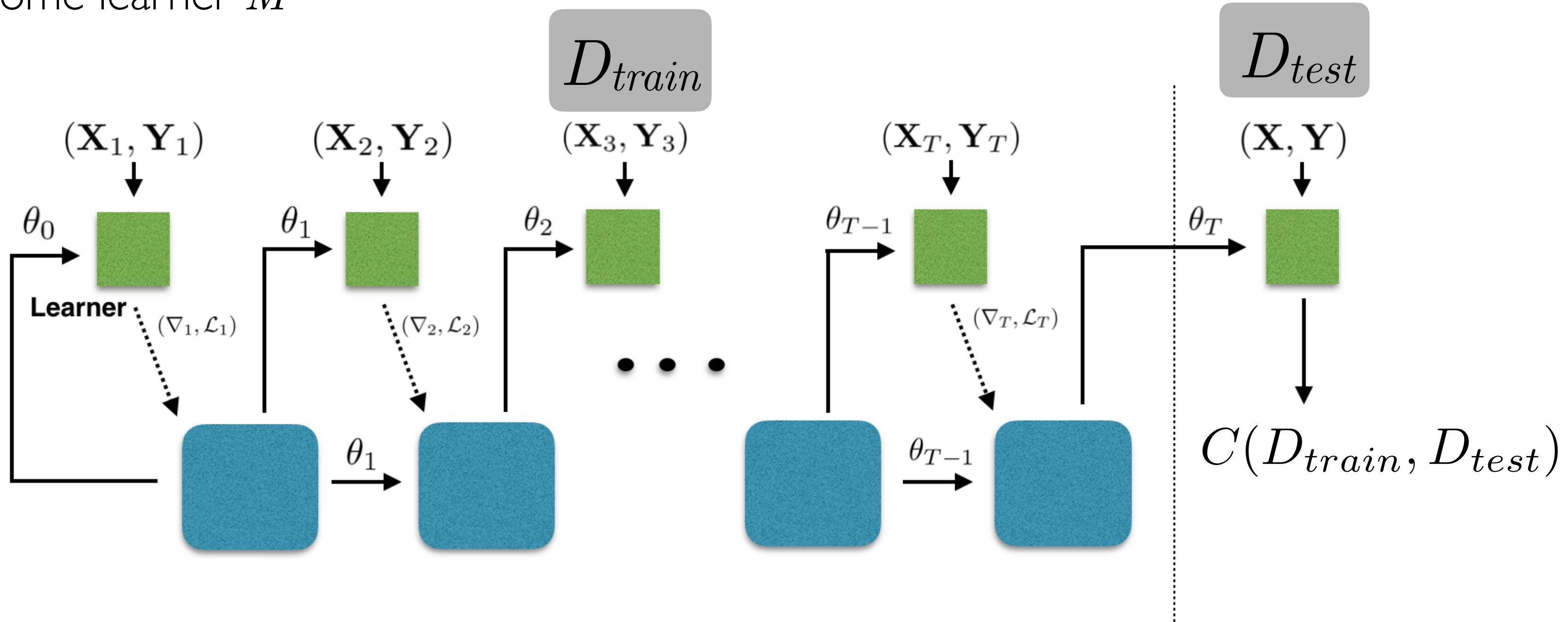
$$\phi \equiv \Theta$$



- Prototypical Networks for Few-shot Learning (2017)  
Jake Snell, Kevin Swersky and Richard Zemel

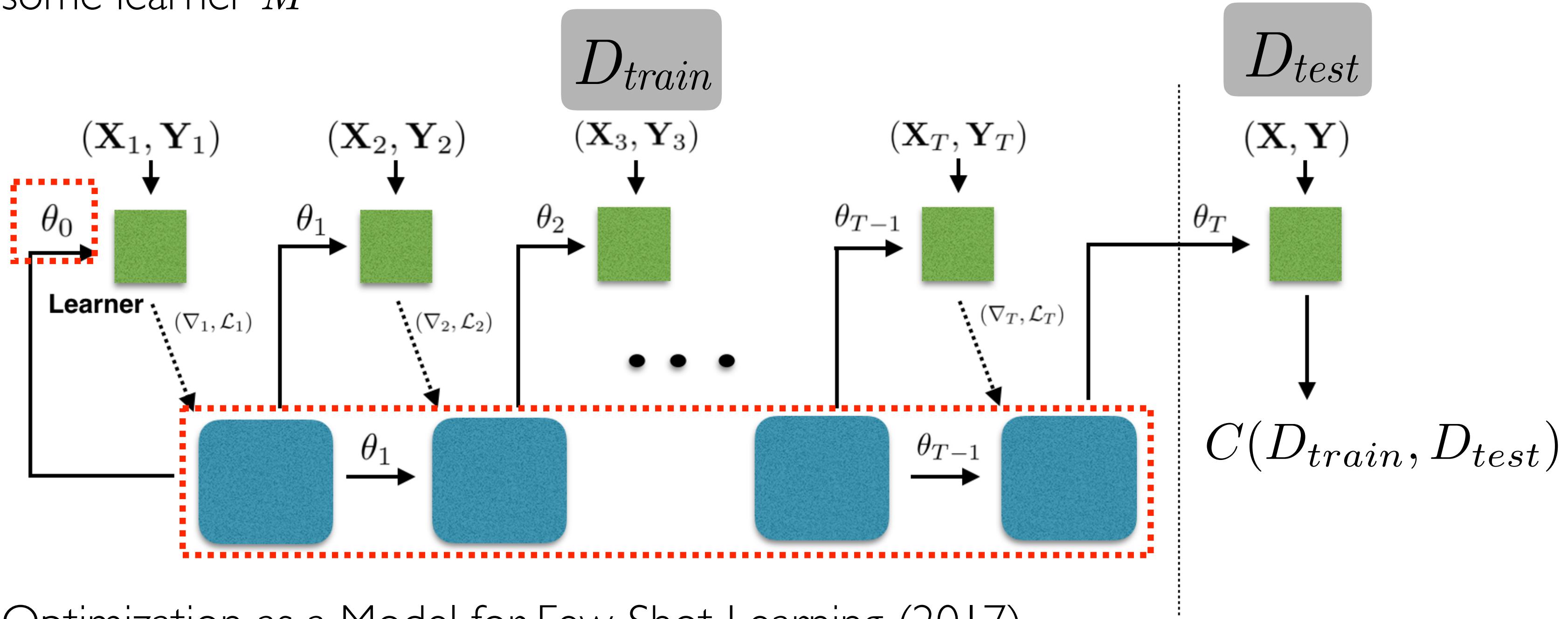
# META-LEARNER LSTM

- Training an “**initialize and gradient descent procedure**” applied on some learner  $M$



# META-LEARNER LSTM

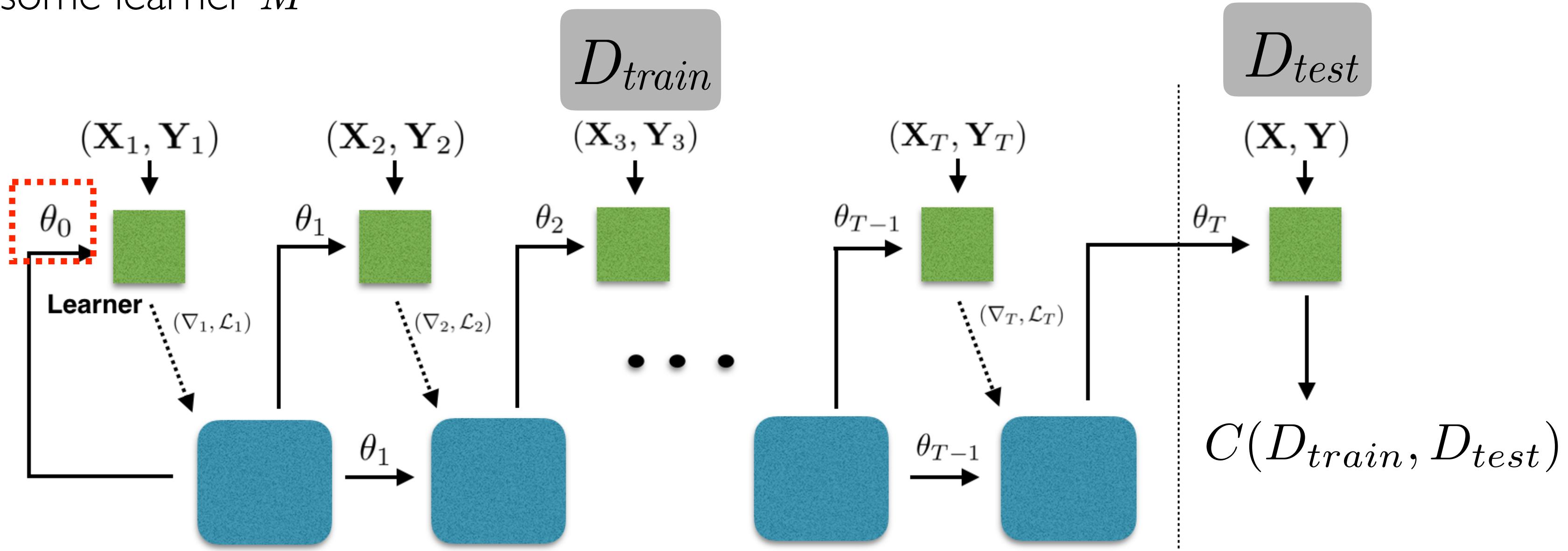
- Training an “**initialize and gradient descent procedure**” applied on some learner  $M$



- Optimization as a Model for Few-Shot Learning (2017)  
*Sachin Ravi and Hugo Larochelle*

# META-LEARNER LSTM

- Training an “**initialize and gradient descent procedure**” applied on some learner  $M$



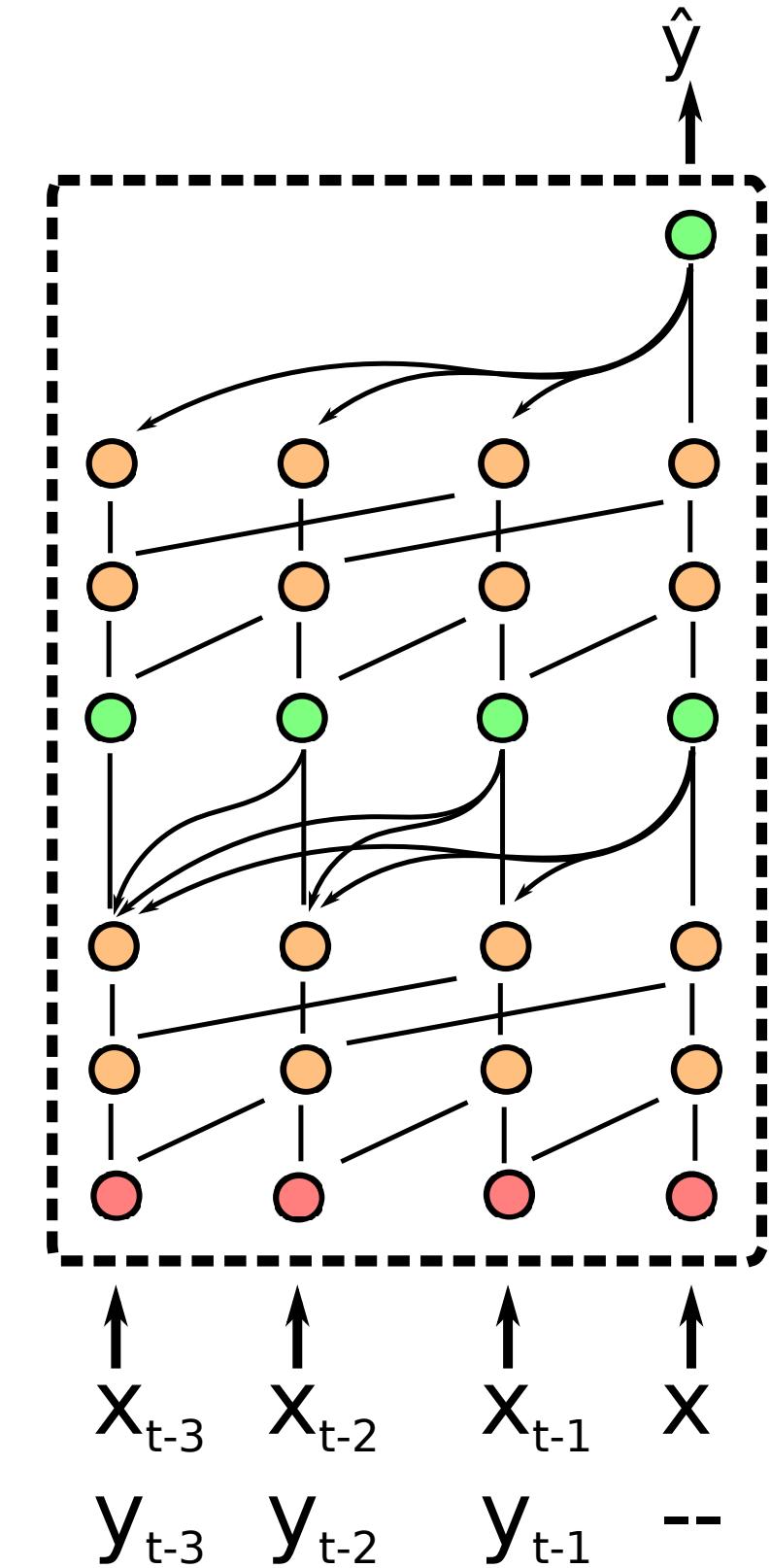
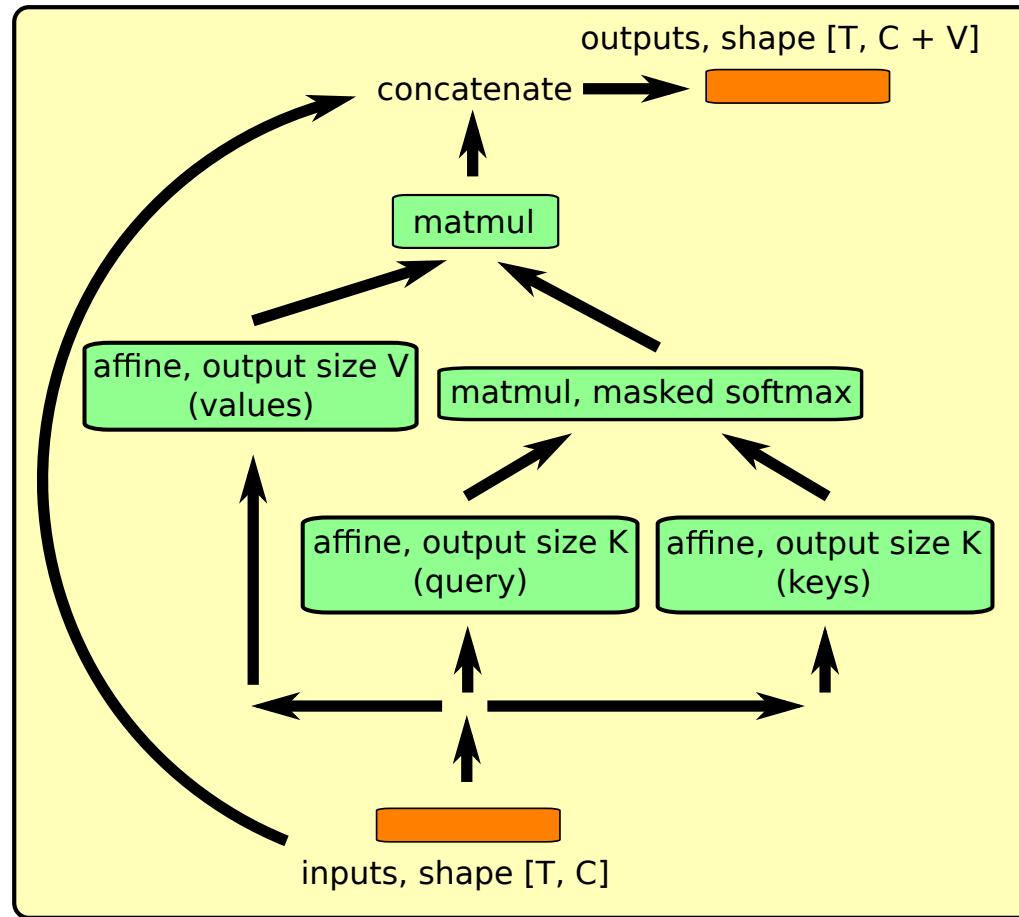
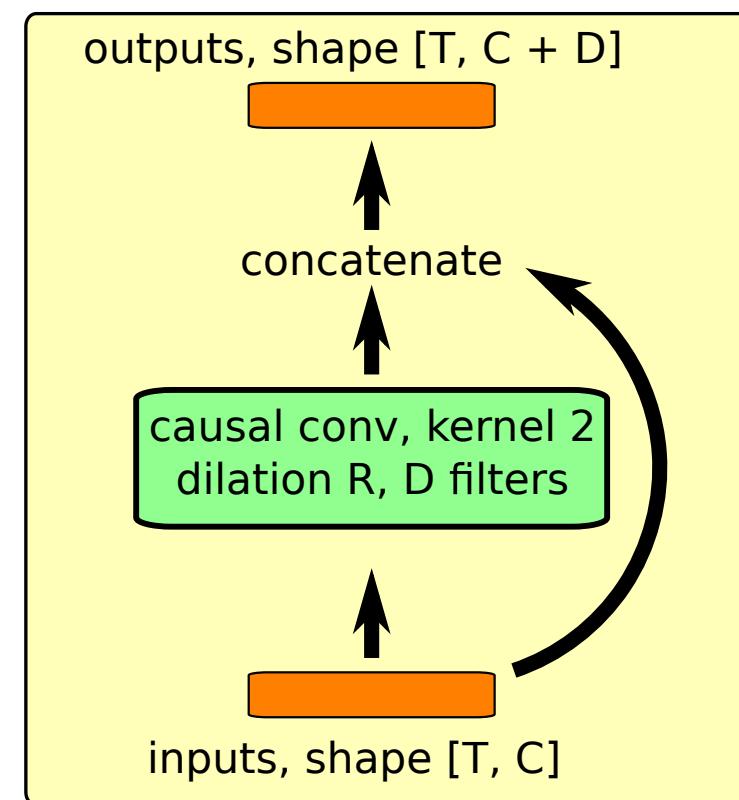
- Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks (2017)  
*Chelsea Finn, Pieter Abbeel and Sergey Levine*

# CHOOSING A META-LEARNER

- How to parametrize learning algorithms (meta-learners  $p(y|\mathbf{x}, D_{train})$  )?
- Two approaches to defining a meta-learner
  - ▶ Take inspiration from a known learning algorithm
    - kNN/kernel machine: Matching networks (Vinyals et al. 2016)
    - Gaussian classifier: Prototypical Networks (Snell et al. 2017)
    - Gradient Descent: Meta-Learner LSTM (Ravi & Larochelle, 2017) , MAML (Finn et al. 2017)
  - ▶ **Derive it from a black box neural network**
    - SNAIL (Mishra et al. 2018)

# SIMPLE NEURAL ATTENTIVE LEARNER

- Using a **convolutional/attentional network** to represent  $p(y|\mathbf{x}, D_{train})$ 
  - alternates between **dilated convolutional layers** and **attentional layers**
  - when inputs are images, an convolutional embedding network is used to map to a vector space



- A Simple Neural Attentive Meta-Learner (2018)  
Nikhil Mishra, Mostafa Rohaninejad, Xi Chen and Pieter Abbeel

# AND SO MUCH MORE!!!



Hugo Larochelle - Few-shot Learning with Meta-Learning:  
Progress Made and Challenges Ahead

[bit.ly/2PikS82](https://bit.ly/2PikS82)

# EXPERIMENT

- Mini-ImageNet (split used in Ravi & Larochelle, 2017)
  - ▶ random subset of 100 classes (64 training, 16 validation, 20 testing)
  - ▶ random sets  $D_{train}$  are generated by randomly picking 5 classes from class subset

<b>Model</b>	<b>5-class</b>	
	<b>1-shot</b>	<b>5-shot</b>
<b>Baseline-finetune</b>	$28.86 \pm 0.54\%$	$49.79 \pm 0.79\%$
<b>Baseline-nearest-neighbor</b>	$41.08 \pm 0.70\%$	$51.04 \pm 0.65\%$
<b>Matching Network</b>	$43.40 \pm 0.78\%$	$51.09 \pm 0.71\%$
<b>Matching Network FCE</b>	$43.56\% \pm 0.84\%$	$55.31\% \pm 0.73\%$
<b>Meta-Learner LSTM (OURS)</b>	$43.44\% \pm 0.77\%$	$60.60\% \pm 0.71\%$

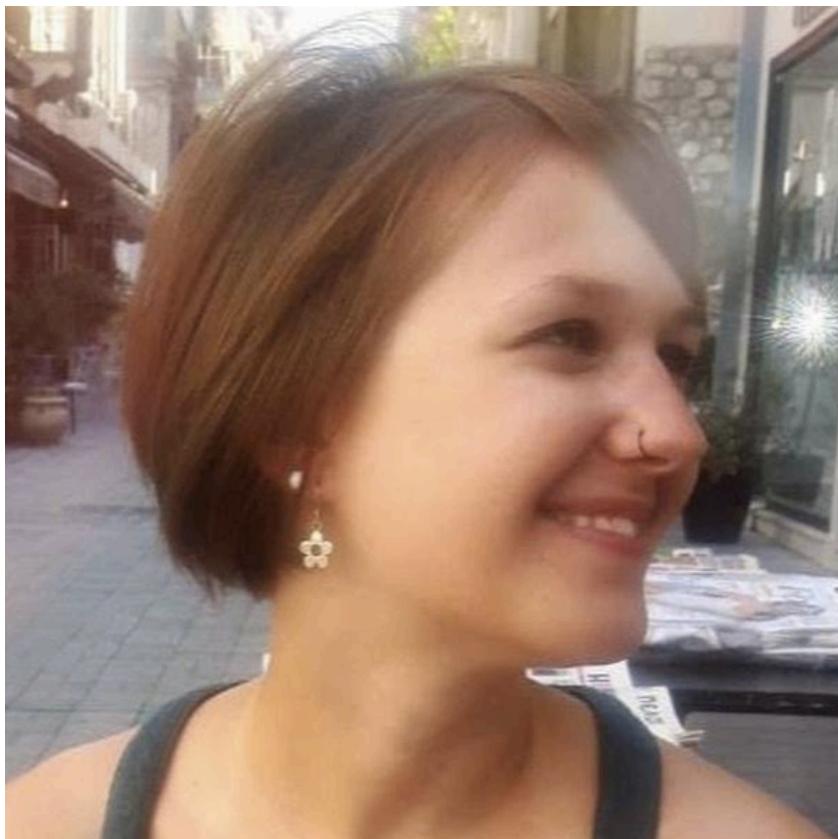
# EXPERIMENT

- Mini-ImageNet (split used in Ravi & Larochelle, 2017)
  - ▶ random subset of 100 classes (64 training, 16 validation, 20 testing)
  - ▶ random sets  $D_{train}$  are generated by randomly picking 5 classes from class subset

<b>Model</b>	<b>5-class</b>	
	<b>1-shot</b>	<b>5-shot</b>
<b>Prototypical Nets</b> (Snell et al.)	<b><math>49.42\% \pm 0.78\%</math></b>	<b><math>68.20\% \pm 0.66\%</math></b>
<b>MAML</b> (Finn et al.)	$48.70\% \pm 1.84\%$	$63.10\% \pm 0.92\%$
<b>SNAIL</b> (Mishra et al.)	<b><math>55.71\% \pm 0.99\%</math></b>	<b><math>68.88\% \pm 0.98\%</math></b>
<b>Matching Network FCE</b>	$43.56\% \pm 0.84\%$	$55.31\% \pm 0.73\%$
<b>Meta-Learner LSTM (OURS)</b>	$43.44\% \pm 0.77\%$	$60.60\% \pm 0.71\%$

# REMAINING CHALLENGES

- Going beyond supervised classification
  - ▶ unsupervised learning, structured output, interactive learning
- Going beyond Mini-ImageNet
  - ▶ coming up with a realistic definition of distributions over problems/datasets



---

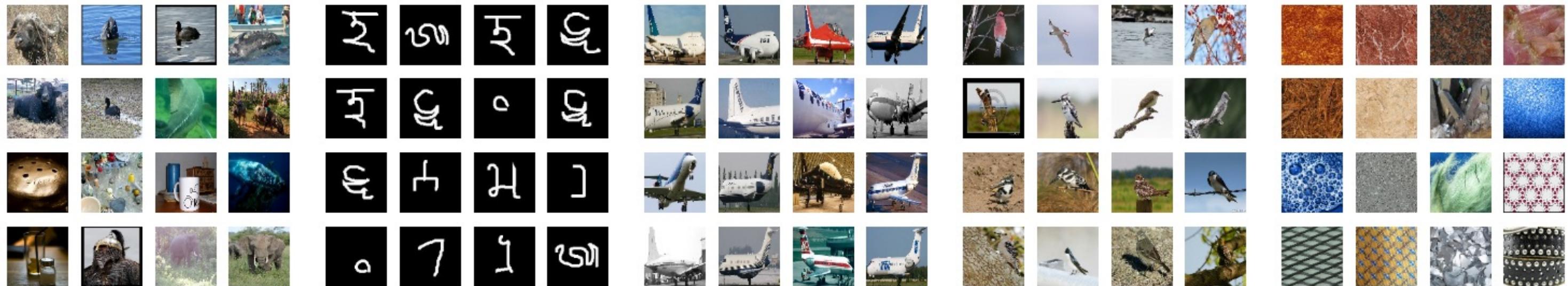
## **Meta-Dataset: A Dataset of Datasets for Learning to Learn from Few Examples**

---

**Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Kelvin Xu, Ross Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, Hugo Larochelle**

# META-DATASET

- To learn across many tasks requires learning over **many datasets**



(a) ImageNet

(b) Omniglot

(c) Aircraft

(d) Birds

(e) DTD



(f) Quick Draw

(g) Fungi

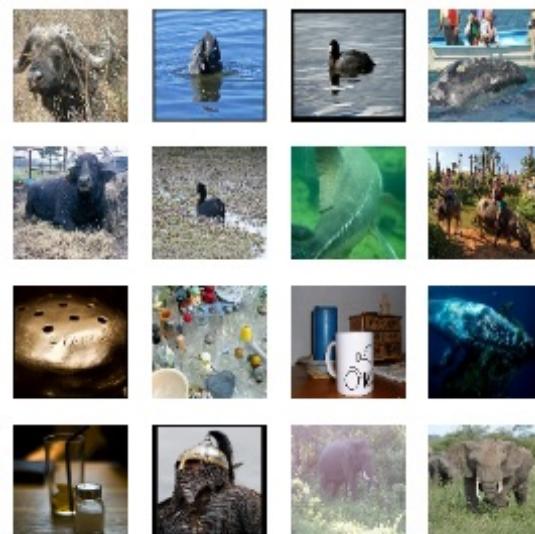
(h) VGG Flower

(i) Traffic Signs

(j) MSCOCO

# META-DATASET

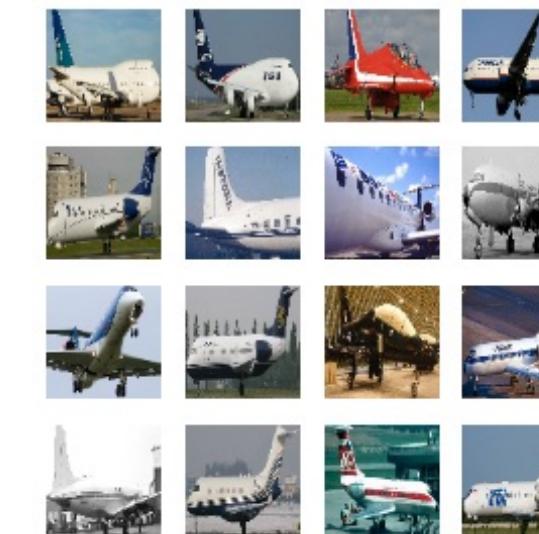
- To learn across many tasks requires learning over **many datasets**



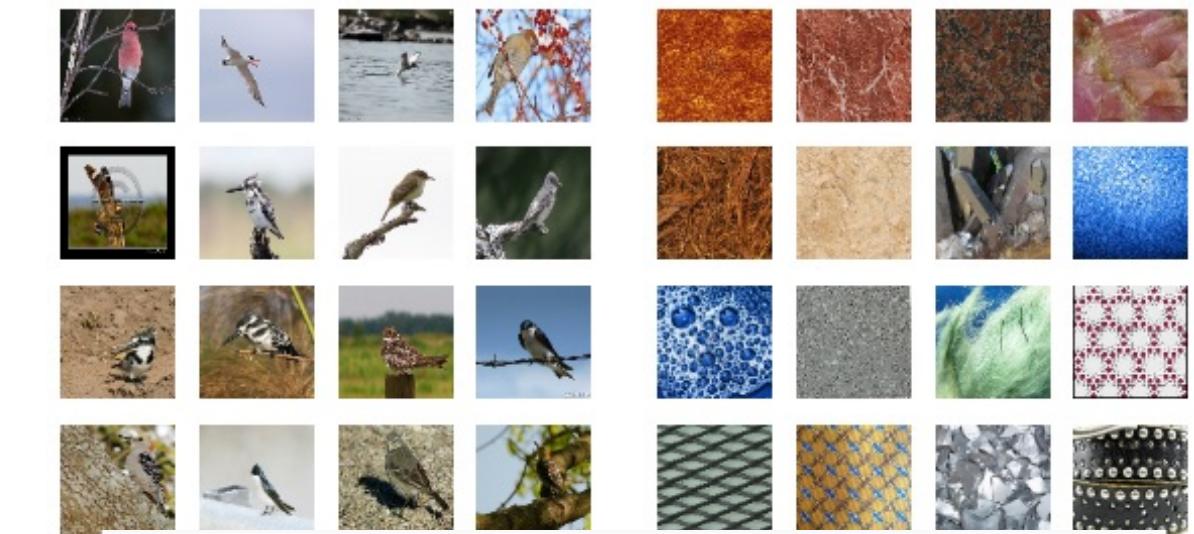
(a) ImageNet



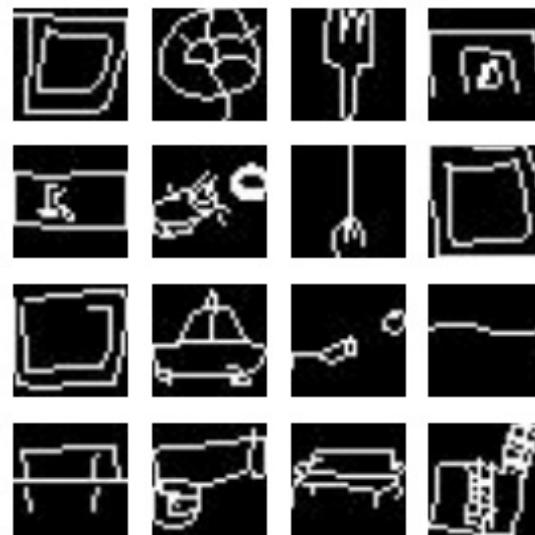
(b) Omniglot



(c) Aircraft



Held out for testing



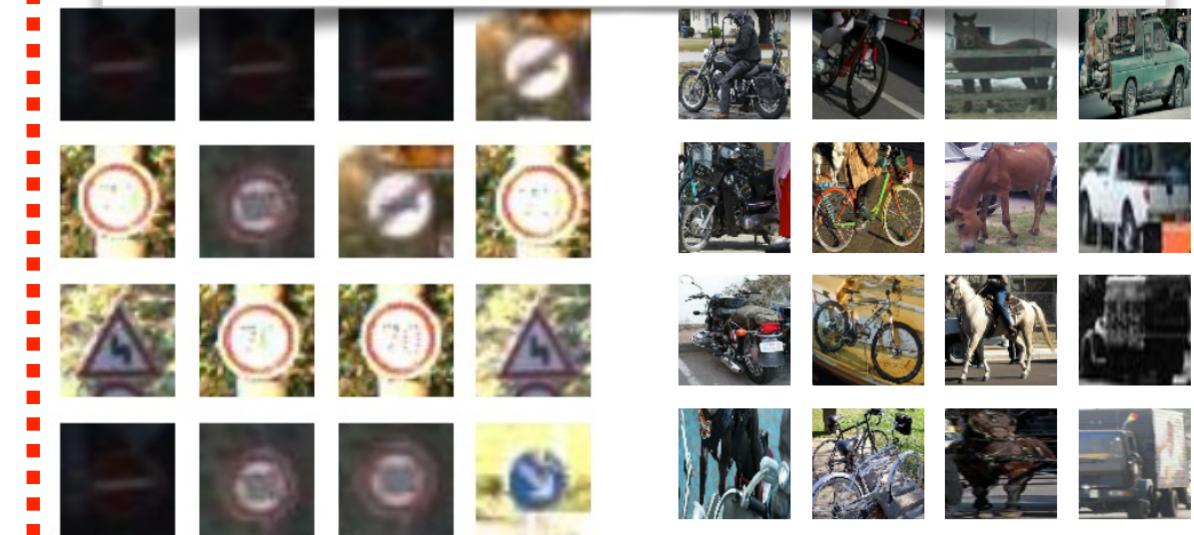
(f) Quick Draw



(g) Fungi



(h) VGG Flower



(i) Traffic Signs (j) MSCOCO

# META-DATASET

- Meta-training only on ImageNet

Test Source	Method: Accuracy $\pm$ confidence				
	$k$ -NN	Finetune	MatchingNet	ProtoNet	MAML
ILSVRC	34.70 $\pm$ 0.95	38.34 $\pm$ 1.12	40.89 $\pm$ 1.08	<b>43.37<math>\pm</math>1.17</b>	38.10 $\pm$ 1.13
Omniglot	59.84 $\pm$ 0.96	59.19 $\pm$ 1.18	61.85 $\pm$ 1.00	<b>66.18<math>\pm</math>1.12</b>	54.00 $\pm$ 1.47
Aircraft	36.47 $\pm$ 0.93	<b>41.18<math>\pm</math>1.07</b>	<b>41.91<math>\pm</math>0.96</b>	<b>42.14<math>\pm</math>0.97</b>	<b>42.52<math>\pm</math>1.16</b>
Birds	40.38 $\pm$ 1.09	45.82 $\pm$ 1.25	54.26 $\pm$ 1.16	<b>57.85<math>\pm</math>1.23</b>	50.78 $\pm$ 1.32
Textures	56.45 $\pm$ 0.78	58.06 $\pm$ 0.88	<b>61.70<math>\pm</math>0.84</b>	<b>60.95<math>\pm</math>0.80</b>	<b>61.26<math>\pm</math>0.93</b>
Quick Draw	36.09 $\pm$ 1.19	38.43 $\pm$ 1.39	38.52 $\pm$ 1.12	<b>44.02<math>\pm</math>1.35</b>	30.71 $\pm$ 1.51
Fungi	23.70 $\pm$ 0.97	22.20 $\pm$ 0.92	27.21 $\pm$ 0.97	<b>31.18<math>\pm</math>1.15</b>	20.35 $\pm$ 0.87
VGG Flower	66.16 $\pm$ 0.99	69.32 $\pm$ 1.13	75.05 $\pm$ 0.91	<b>79.89<math>\pm</math>0.90</b>	65.12 $\pm$ 1.15
Traffic Signs	<b>44.81<math>\pm</math>1.47</b>	39.36 $\pm$ 1.28	<b>45.36<math>\pm</math>1.31</b>	<b>44.04<math>\pm</math>1.24</b>	31.10 $\pm$ 1.20
MSCOCO	29.69 $\pm$ 1.00	30.25 $\pm$ 1.17	32.32 $\pm$ 1.08	<b>36.44<math>\pm</math>1.23</b>	25.17 $\pm$ 1.15
<b>Avg. rank</b>	4	3.4	2.2	<b>1.35</b>	4.05

# META-DATASET

- Meta-training on all training datasets

Test Source	Method: Accuracy $\pm$ confidence				
	<i>k</i> -NN	Finetune	MatchingNet	ProtoNet	MAML
ILSVRC	25.88 $\pm$ 0.83	25.84 $\pm$ 0.83	35.88 $\pm$ 0.98	<b>38.51<math>\pm</math>1.01</b>	30.56 $\pm$ 1.00
Omniglot	<b>92.45<math>\pm</math>0.41</b>	85.20 $\pm$ 0.73	90.21 $\pm$ 0.46	91.32 $\pm$ 0.50	78.05 $\pm$ 0.98
Aircraft	54.60 $\pm$ 0.97	58.22 $\pm$ 1.02	<b>70.71<math>\pm</math>0.78</b>	<b>71.54<math>\pm</math>0.84</b>	68.62 $\pm$ 0.90
Birds	36.74 $\pm$ 1.01	38.56 $\pm$ 1.08	59.28 $\pm$ 1.06	<b>61.81<math>\pm</math>1.13</b>	54.59 $\pm$ 1.24
Textures	50.06 $\pm$ 0.77	48.37 $\pm$ 0.82	<b>60.61<math>\pm</math>0.82</b>	<b>59.31<math>\pm</math>0.75</b>	<b>59.25<math>\pm</math>0.80</b>
Quick Draw	<b>59.54<math>\pm</math>1.08</b>	54.05 $\pm$ 1.30	57.44 $\pm$ 1.17	<b>60.99<math>\pm</math>1.21</b>	44.48 $\pm$ 1.41
Fungi	24.60 $\pm$ 0.95	22.90 $\pm$ 0.95	31.10 $\pm$ 1.04	<b>35.96<math>\pm</math>1.25</b>	21.12 $\pm$ 0.88
VGG Flower	62.49 $\pm$ 0.91	59.72 $\pm$ 1.17	76.72 $\pm$ 0.83	<b>81.06<math>\pm</math>0.87</b>	66.05 $\pm$ 1.09
Traffic Signs	<b>41.68<math>\pm</math>1.46</b>	30.02 $\pm$ 1.13	<b>43.20<math>\pm</math>1.33</b>	39.95 $\pm$ 1.18	30.23 $\pm$ 1.24
MSCOCO	23.55 $\pm$ 0.99	23.01 $\pm$ 0.96	26.87 $\pm$ 1.00	<b>30.81<math>\pm</math>1.13</b>	21.13 $\pm$ 1.06
<b>Avg. rank</b>	3.4	4.3	2.15	1.4	3.75

# META-DATASET

- Difference in performance when meta-training on all datasets

Test Source	Method: Accuracy $\pm$ confidence				
	$k$ -NN	Finetune	MatchingNet	ProtoNet	MAML
ILSVRC	-8.82 $\pm$ 1.26	-12.5 $\pm$ 1.39	-5.01 $\pm$ 1.46	-4.86 $\pm$ 1.55	-7.54 $\pm$ 1.51
Omniglot	32.61 $\pm$ 1.04	26.01 $\pm$ 1.39	28.36 $\pm$ 1.1	25.14 $\pm$ 1.23	24.05 $\pm$ 1.77
Aircraft	18.13 $\pm$ 1.34	17.04 $\pm$ 1.48	28.8 $\pm$ 1.24	29.4 $\pm$ 1.28	26.1 $\pm$ 1.47
Birds	-3.64 $\pm$ 1.49	-7.26 $\pm$ 1.65	5.02 $\pm$ 1.57	3.96 $\pm$ 1.67	3.81 $\pm$ 1.81
Textures	-6.39 $\pm$ 1.1	-9.69 $\pm$ 1.2	-1.09 $\pm$ 1.17	-1.64 $\pm$ 1.1	-2.01 $\pm$ 1.23
Quick Draw	23.45 $\pm$ 1.61	15.62 $\pm$ 1.9	18.92 $\pm$ 1.62	16.97 $\pm$ 1.81	13.77 $\pm$ 2.07
Fungi	0.9 $\pm$ 1.36	0.7 $\pm$ 1.32	3.89 $\pm$ 1.42	4.78 $\pm$ 1.7	0.77 $\pm$ 1.24
VGG Flower	-3.67 $\pm$ 1.34	-9.6 $\pm$ 1.63	1.67 $\pm$ 1.23	1.17 $\pm$ 1.25	0.93 $\pm$ 1.58
Traffic Signs	-3.13 $\pm$ 2.07	-9.34 $\pm$ 1.71	-2.16 $\pm$ 1.87	-4.09 $\pm$ 1.71	-0.87 $\pm$ 1.73
MSCOCO	-6.14 $\pm$ 1.41	-7.24 $\pm$ 1.51	-5.45 $\pm$ 1.47	-5.63 $\pm$ 1.67	-4.04 $\pm$ 1.56

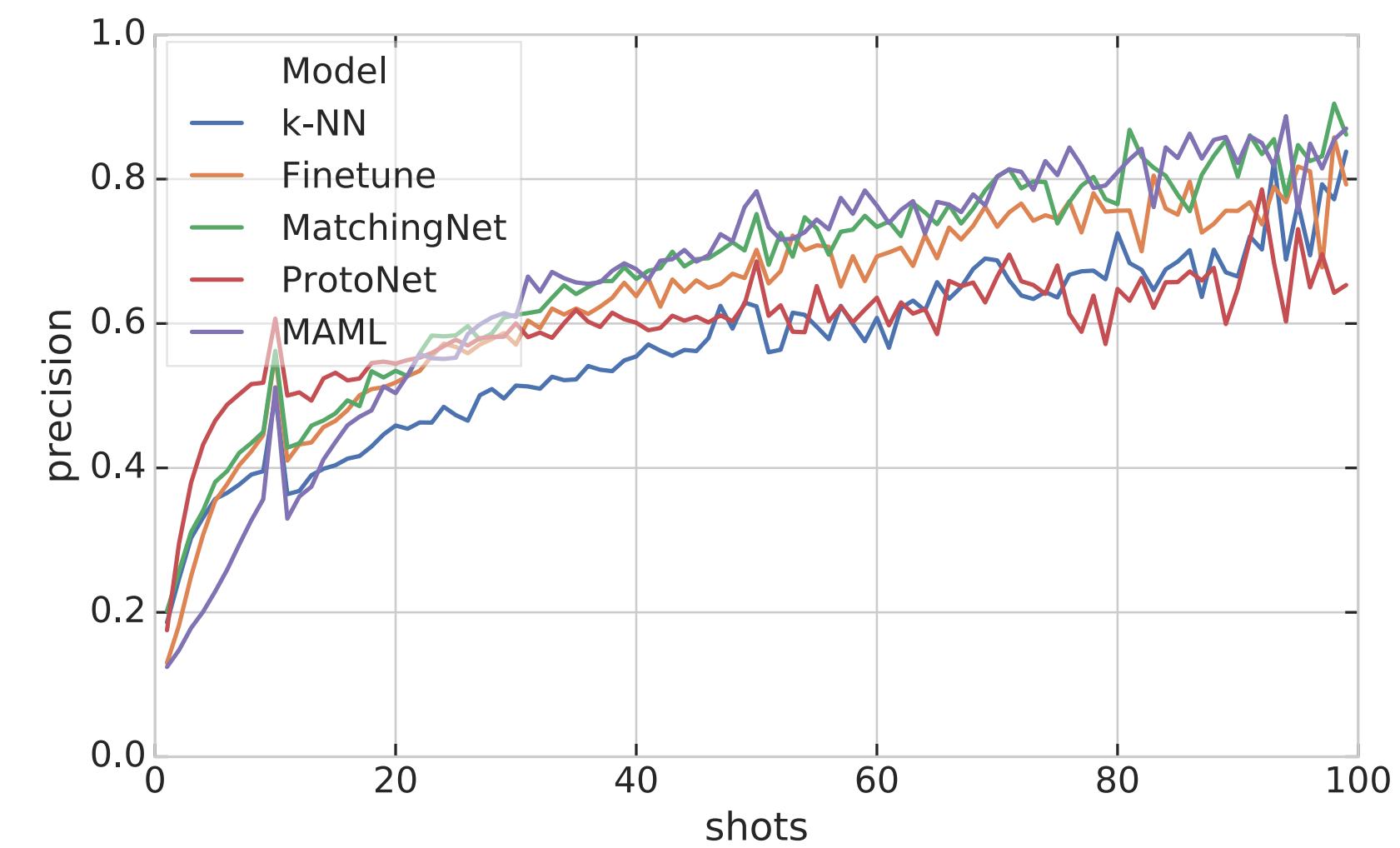
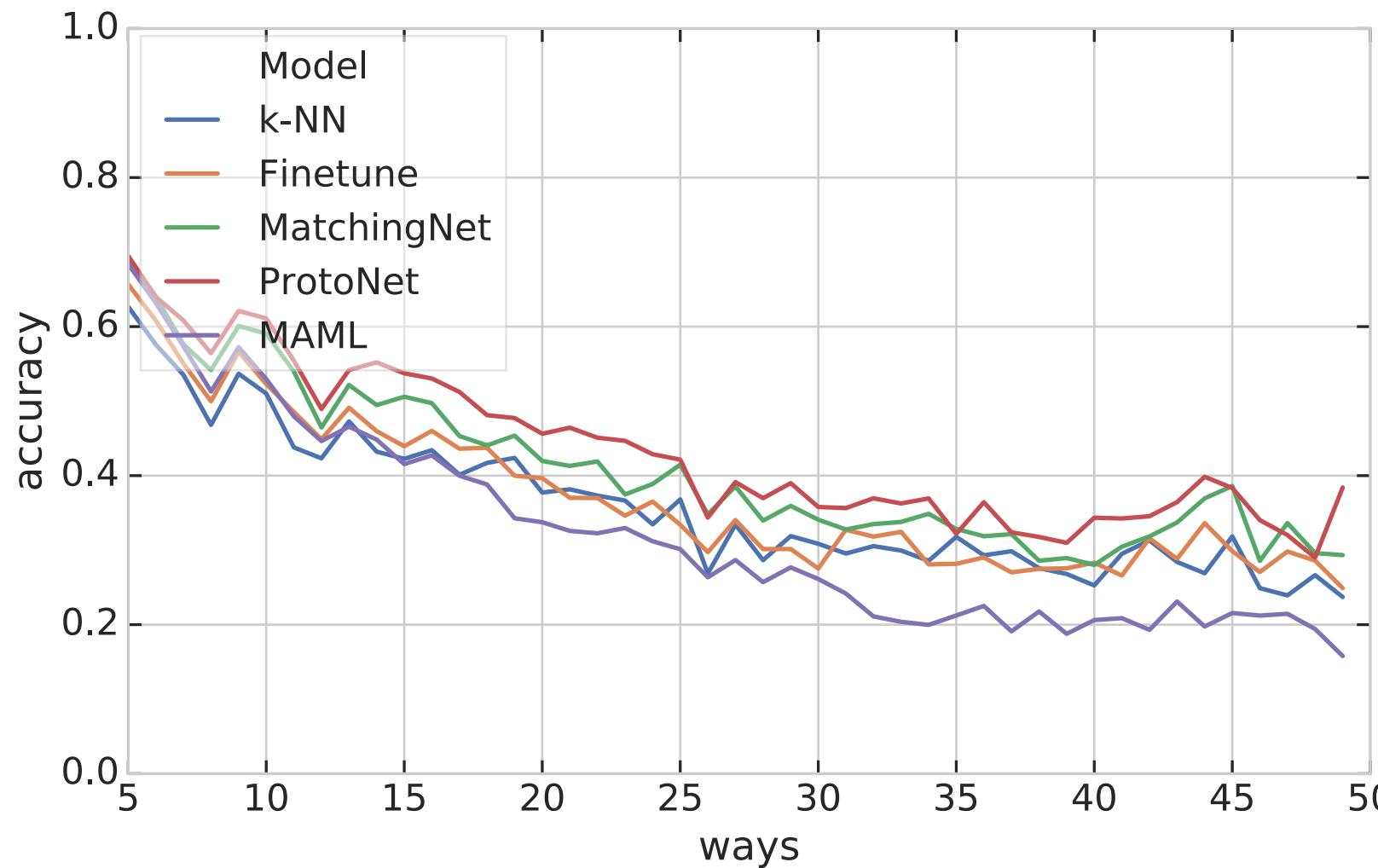
# META-DATASET

- Difference in performance when meta-training on all datasets

Test Source	Method: Accuracy $\pm$ confidence				
	$k$ -NN	Finetune	MatchingNet	ProtoNet	MAML
ILSVRC	-8.82 $\pm$ 1.26	-12.5 $\pm$ 1.39	-5.01 $\pm$ 1.46	-4.86 $\pm$ 1.55	-7.54 $\pm$ 1.51
Omniglot	32.61 $\pm$ 1.04	26.01 $\pm$ 1.39	28.36 $\pm$ 1.1	25.14 $\pm$ 1.23	24.05 $\pm$ 1.77
Aircraft	18.13 $\pm$ 1.34	17.04 $\pm$ 1.48	28.8 $\pm$ 1.24	29.4 $\pm$ 1.28	26.1 $\pm$ 1.47
Birds	-3.64 $\pm$ 1.49	-7.26 $\pm$ 1.65	5.02 $\pm$ 1.57	3.96 $\pm$ 1.67	3.81 $\pm$ 1.81
Textures	-6.39 $\pm$ 1.1	-9.69 $\pm$ 1.2	-1.09 $\pm$ 1.17	-1.64 $\pm$ 1.1	-2.01 $\pm$ 1.23
Quick Draw	23.45 $\pm$ 1.61	15.62 $\pm$ 1.9	18.92 $\pm$ 1.62	16.97 $\pm$ 1.81	13.77 $\pm$ 2.07
Fungi	0.9 $\pm$ 1.36	0.7 $\pm$ 1.32	3.89 $\pm$ 1.42	4.78 $\pm$ 1.7	0.77 $\pm$ 1.24
VGG Flower	-3.67 $\pm$ 1.34	-9.6 $\pm$ 1.63	1.67 $\pm$ 1.23	1.17 $\pm$ 1.25	0.93 $\pm$ 1.58
Traffic Signs	-3.13 $\pm$ 2.07	-9.34 $\pm$ 1.71	-2.16 $\pm$ 1.87	-4.09 $\pm$ 1.71	-0.87 $\pm$ 1.73
MSCOCO	-6.14 $\pm$ 1.41	-7.24 $\pm$ 1.51	-5.45 $\pm$ 1.47	-5.63 $\pm$ 1.67	-4.04 $\pm$ 1.56

# META-DATASET

- Varying the number of shots and ways



# TAKE AWAYS (SO FAR)

- Meta-training distribution of episodes can make a big difference  
(at least for current methods)
- Using “regular training” as initialization makes a big difference
- MAML needs to be adjusted to be more robust

# DISCUSSION

- Now is time to move beyond our current simple benchmarks
- What is the “right” meta-training distribution?
- How should we be increasing the size of the benchmark (what should be V2)?
- What are the properties of the optimization landscape of the episodic framework?
- What fairness-related questions does meta-learning pose?

MERCI !