

# The Big Problem with Meta-Learning and How Bayesians Can Fix It

Chelsea Finn



Stanford



training data

Braque



Cezanne



test datapoint



By Braque or Cezanne?



How did you accomplish this?

Through previous experience.

# How might you get a machine to accomplish this task?

Modeling image formation

Geometry

SIFT features, HOG features + SVM

Fine-tuning from ImageNet features

Domain adaptation from other painters

???

Fewer human priors,  
more data-driven priors

Greater success.



Can we explicitly **learn priors from previous experience**  
that lead to efficient downstream learning?

Can we learn to learn?

# Outline

1. Brief overview of meta-learning
2. The problem: peculiar, lesser-known, yet ubiquitous
3. Steps towards a solution

# How does meta-learning work? An example.

Given 1 example of 5 classes:



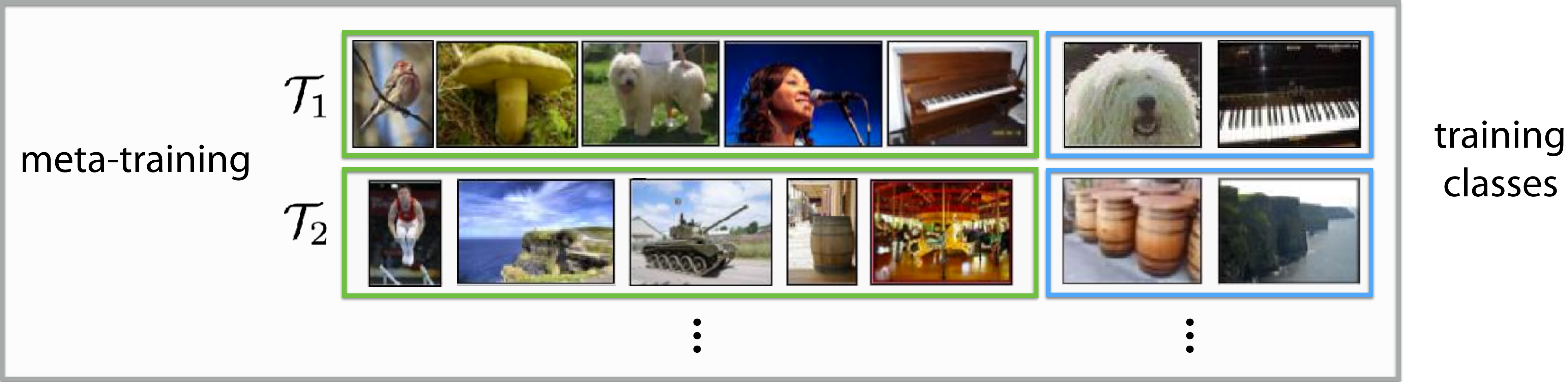
training data  $\mathcal{D}_{\text{train}}$

Classify new examples



test set  $\mathbf{X}_{\text{test}}$

# How does meta-learning work? An example.



Given 1 example of 5 classes:

Classify new examples

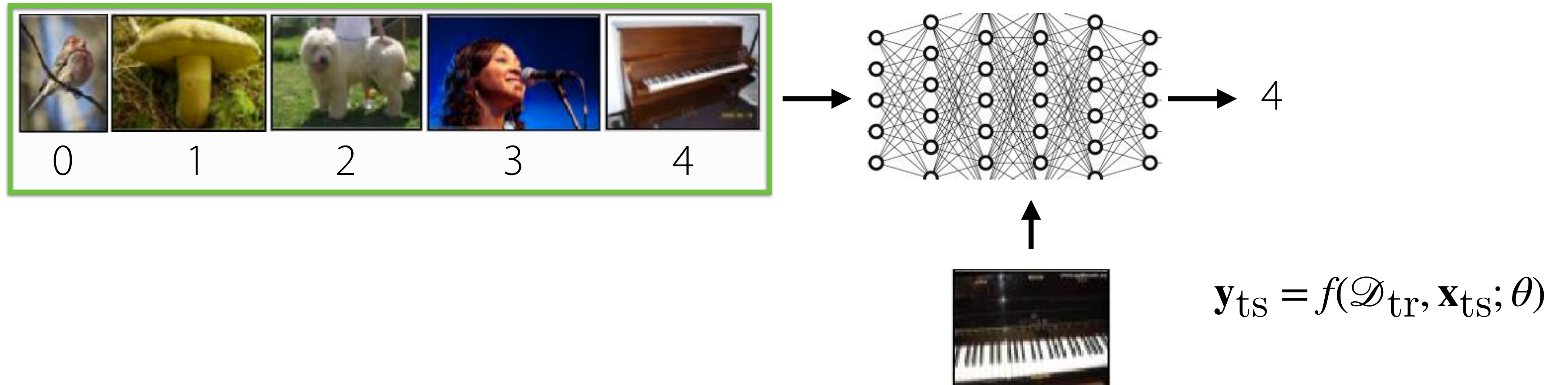




# How does meta-learning work?



One approach: parameterize learner by neural network



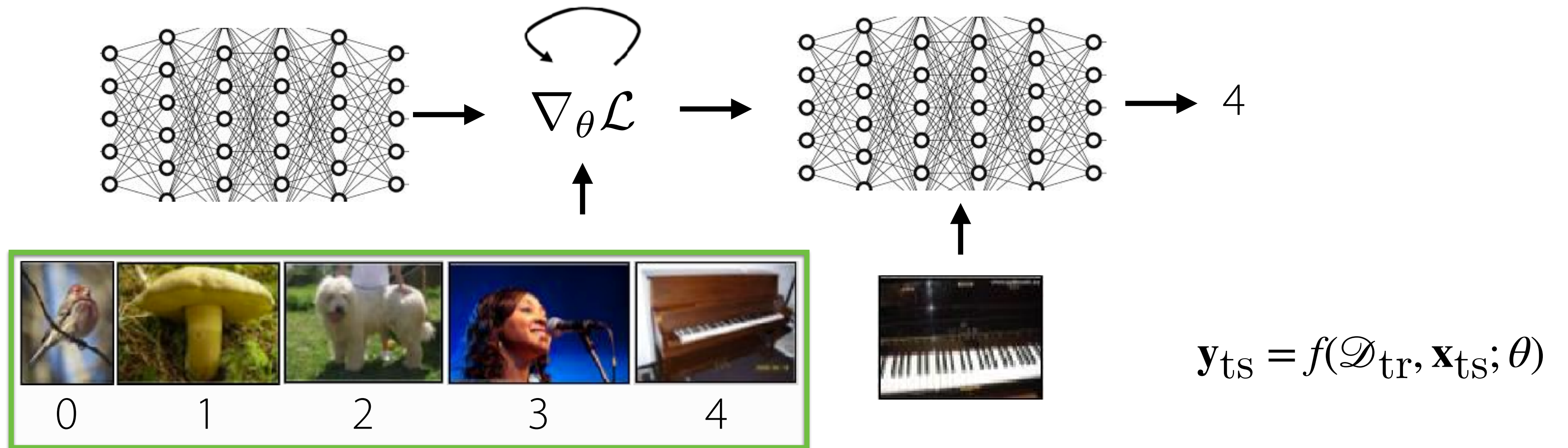
(Hochreiter et al. '91, Santoro et al. '16, many others)



# How does meta-learning work?

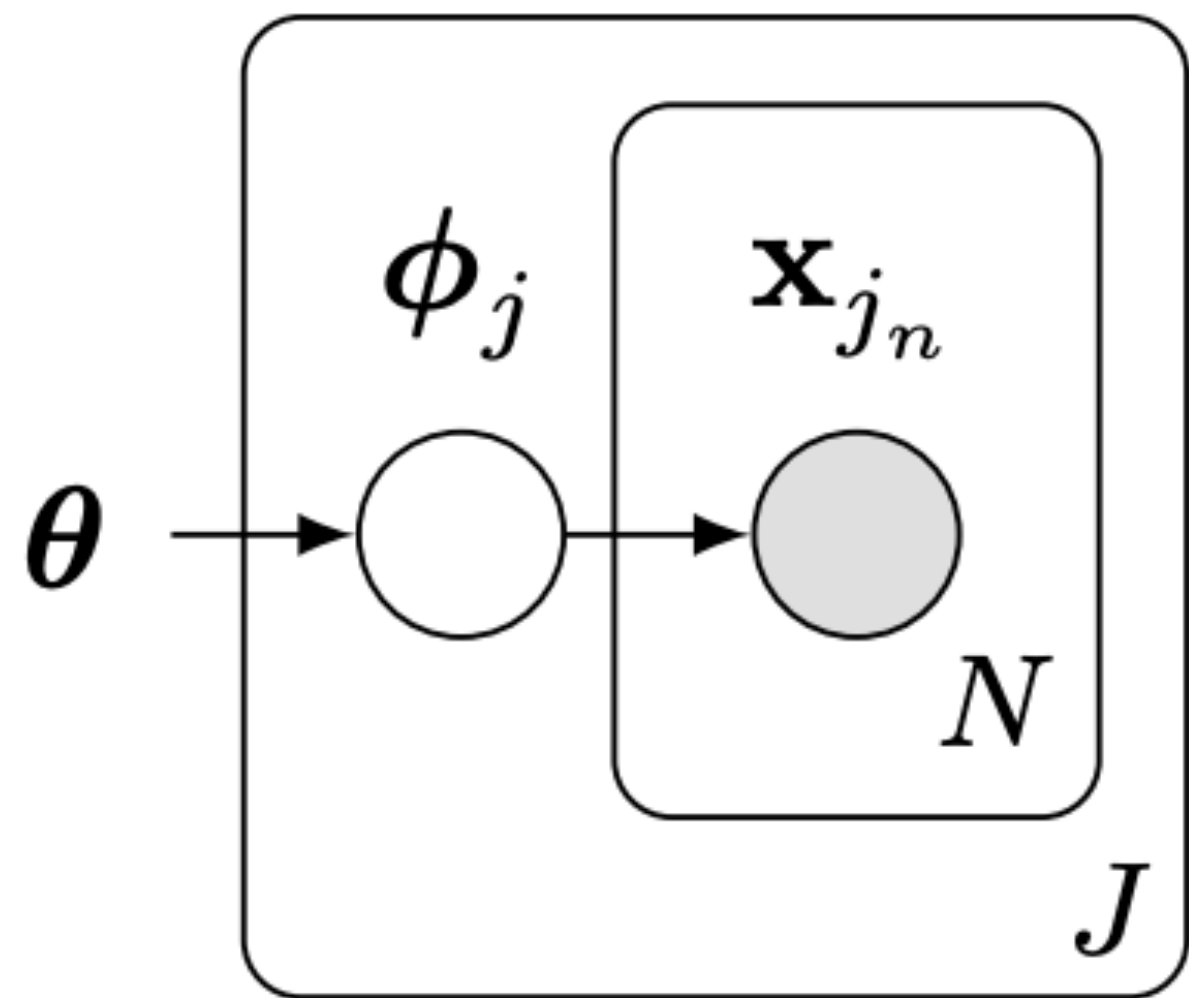


**Another approach:** embed optimization inside the learning process



(Maclaurin et al. '15, Finn et al. '17, many others)

# The Bayesian perspective



meta-learning  $\leftrightarrow$  learning priors  $p(\phi \mid \theta)$  from data

(Grant et al. '18, Gordon et al. '18, many others)



# Outline

1. Brief overview of meta-learning
- 2. The problem: peculiar, lesser-known, yet ubiquitous**
3. First steps towards a solution

# How we construct tasks for meta-learning.



Randomly assign class labels to image classes for each task  $\rightarrow$  Tasks are *mutually exclusive*.

Algorithms **must** use **training data** to infer label ordering.



# What if label order is consistent?

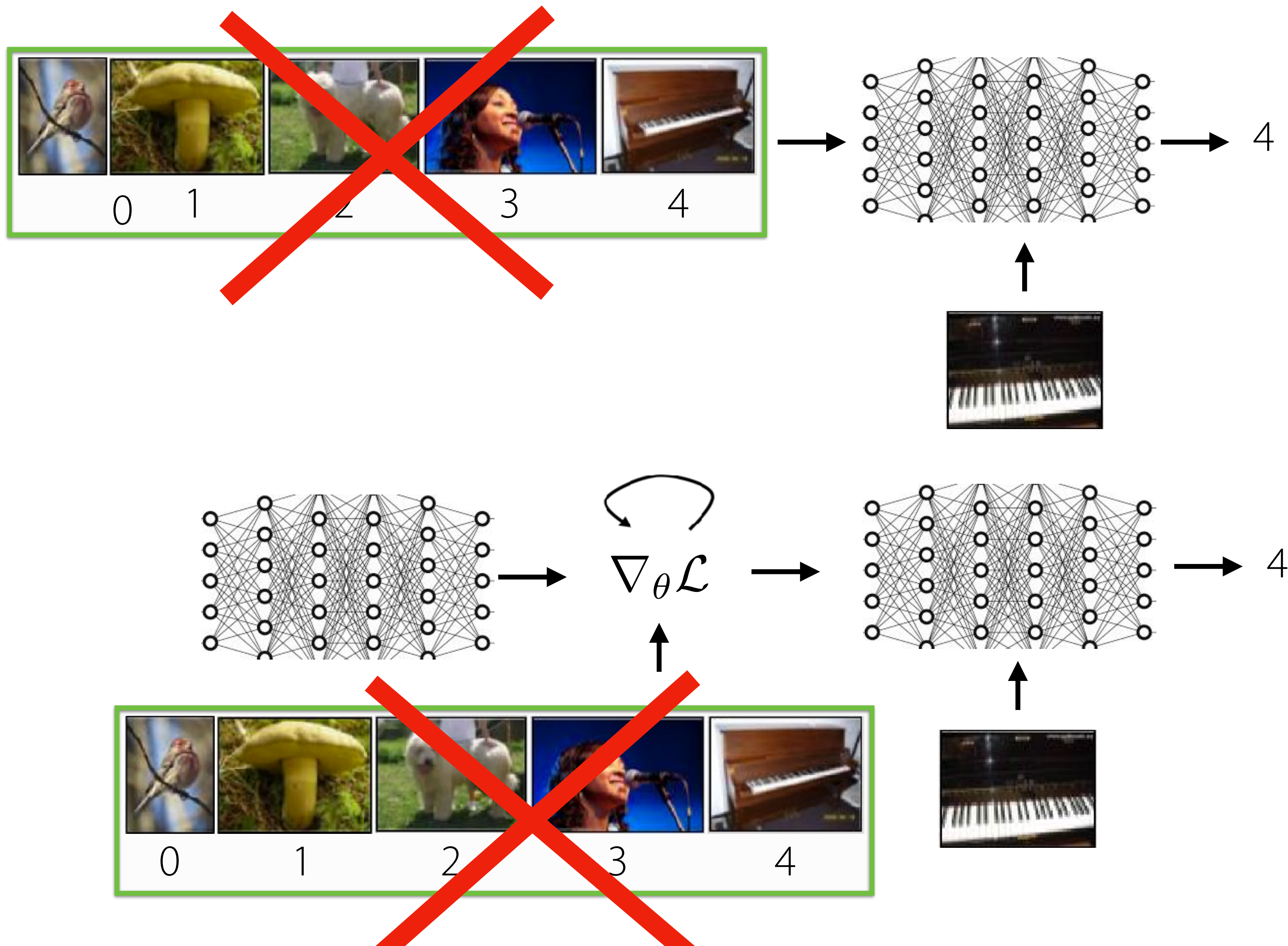


Tasks are **non-mutually exclusive**: a single function can solve all tasks.

The network can simply learn to classify inputs, irrespective of  $\mathcal{D}_{\text{tr}}$



The network can simply learn to classify inputs, irrespective of  $\mathcal{D}_{tr}$





# What if label order is consistent?



For new image classes: can't make predictions w/o  $\mathcal{D}_{\text{tr}}$

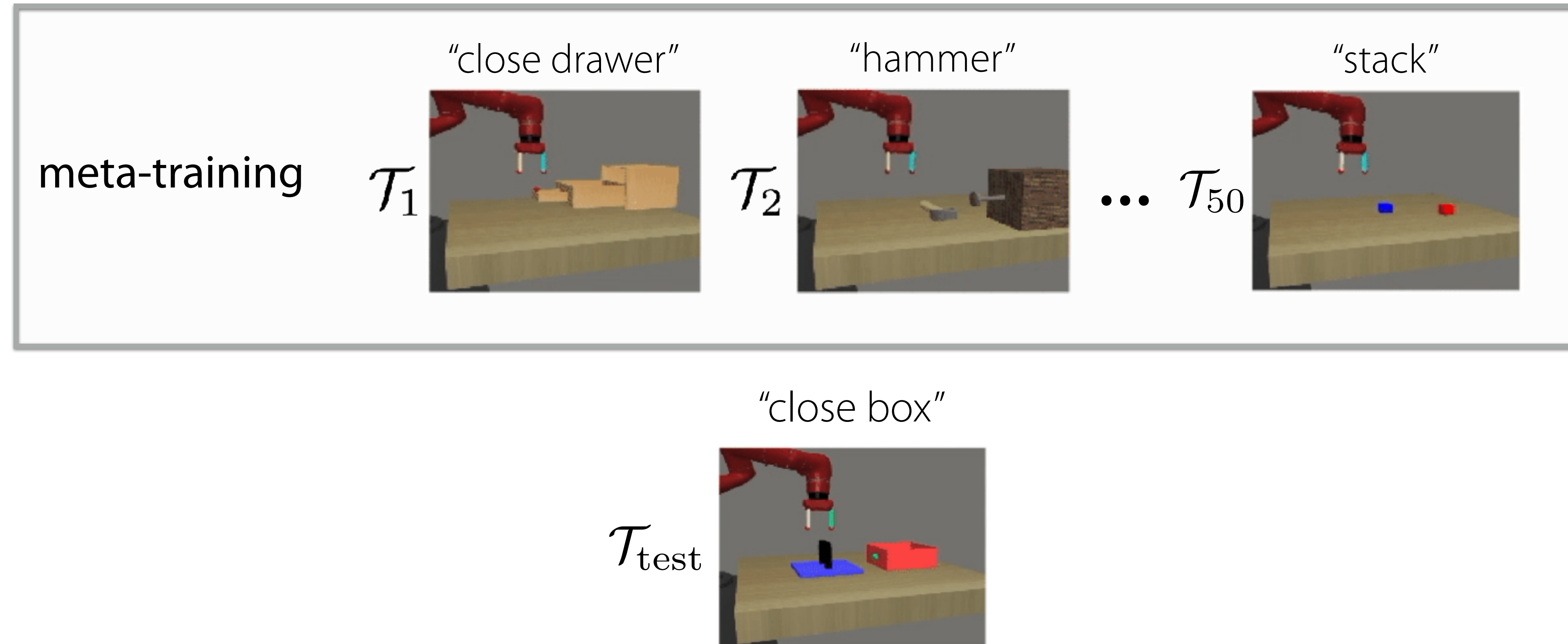
<i>NME Omniglot</i>	20-way 1-shot	20-way 5-shot
<b>MAML</b>	7.8 (0.2)%	50.7 (22.9)%

# Is this a problem?

- **No**: for image classification, we can just shuffle labels\*
- **No**, if we see the same image classes as training (& don't need to adapt at meta-test time)
- But, **yes**, if we want to be able to adapt with data for new tasks.

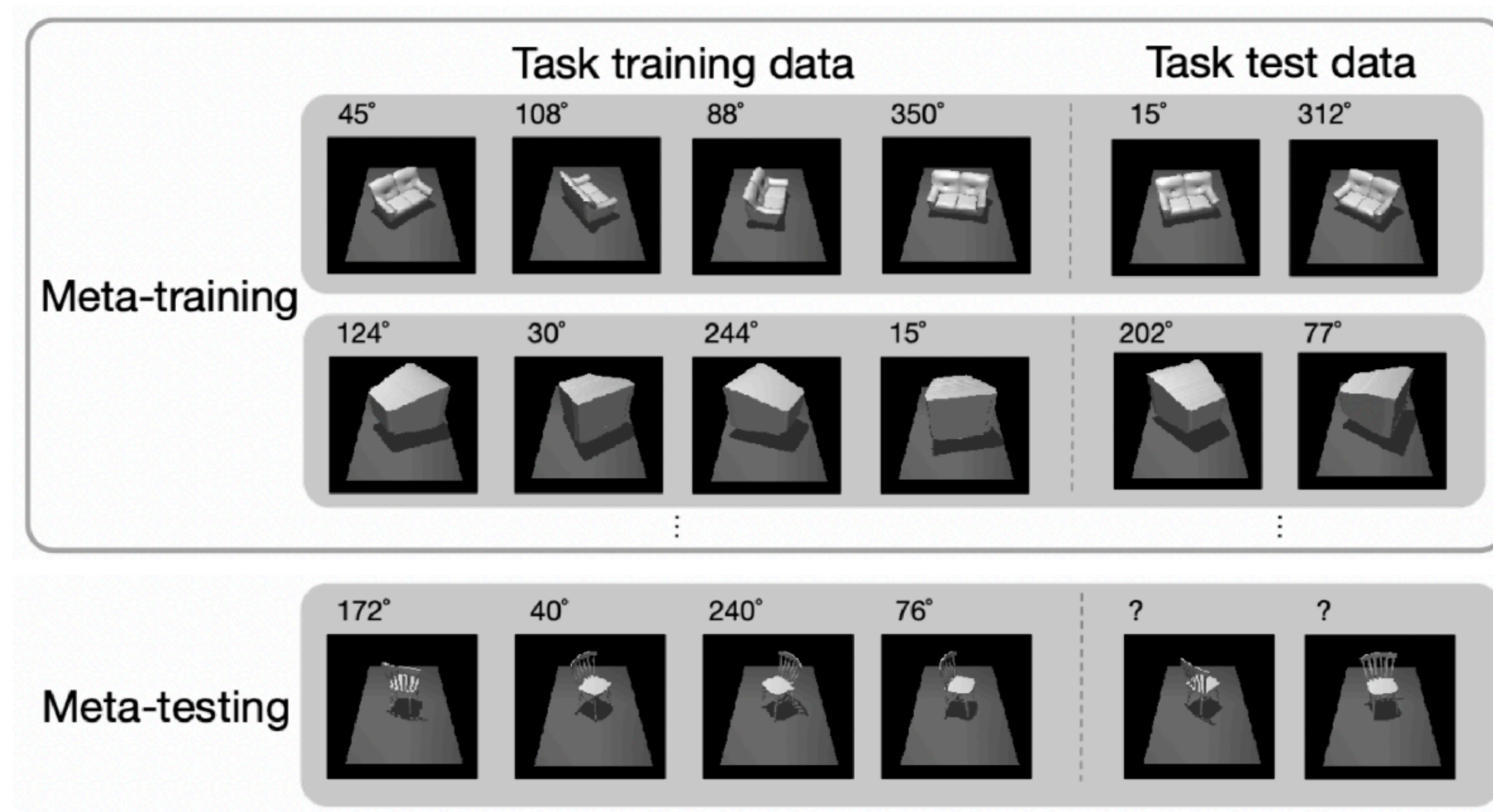


# Another example



If you tell the robot the task goal, the robot can **ignore** the trials.

# Another example



Model can memorize the canonical orientations of the training objects.



Can we do something about it?

If tasks *mutually exclusive*: single function cannot solve all tasks

(i.e. due to label shuffling, hiding information)

If tasks are *non-mutually exclusive*: single function can solve all tasks

*multiple solutions* to the  
meta-learning problem

$$y^{\text{ts}} = f_{\theta}(\mathcal{D}_i^{\text{tr}}, x^{\text{ts}})$$

**One solution:**

memorize canonical pose info in  $\theta$  & ignore  $\mathcal{D}_i^{\text{tr}}$

**Another solution:**

carry no info about canonical pose in  $\theta$ , acquire from  $\mathcal{D}_i^{\text{tr}}$

An entire **spectrum of solutions** based on how **information** flows.

Suggests a potential approach: control information flow.



If tasks are *non-mutually exclusive*: single function can solve all tasks  
*multiple solutions* to the meta-learning problem

$$y^{\text{ts}} = f_{\theta}(\mathcal{D}_i^{\text{tr}}, x^{\text{ts}})$$

**One solution:** memorize canonical pose info in  $\theta$  & ignore  $\mathcal{D}_i^{\text{tr}}$

**Another solution:** carry no info about canonical pose in  $\theta$ , acquire from  $\mathcal{D}_i^{\text{tr}}$

An entire **spectrum of solutions** based on how **information** flows.

---

**Meta-regularization** one option:  $\max I(\hat{\mathbf{y}}_{\text{ts}}, \mathcal{D}_{\text{tr}} | \mathbf{x}_{\text{ts}})$

minimize meta-training loss + information in  $\theta$

$$\mathcal{L}(\theta, \mathcal{D}_{\text{meta-train}}) + \beta D_{KL}(q(\theta; \theta_{\mu}, \theta_{\sigma}) || p(\theta))$$

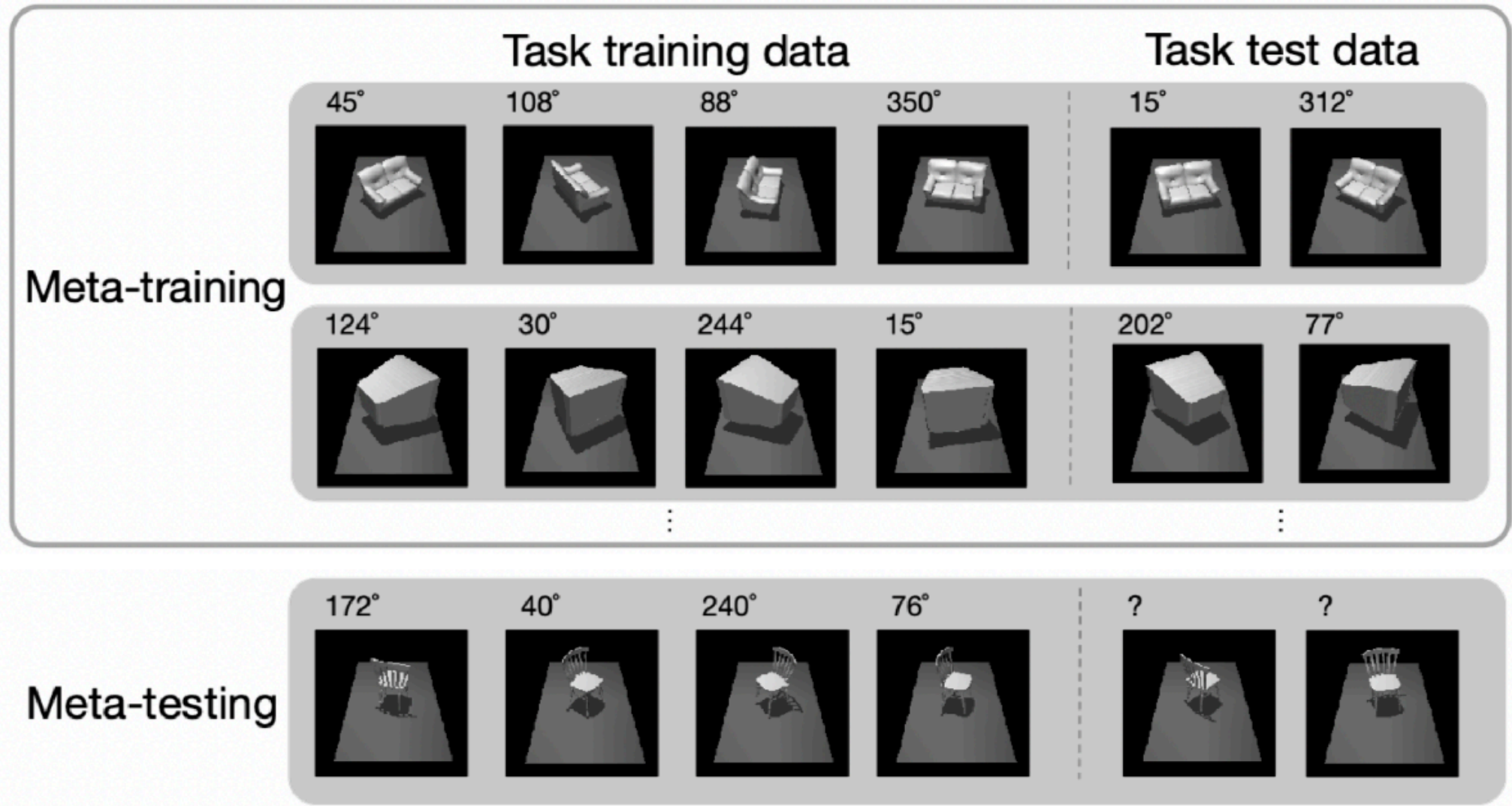
Places precedence on using information from  $\mathcal{D}_{\text{tr}}$  over storing info in  $\theta$ .

Can combine with your favorite meta-learning algorithm.

Omniglot without label shuffling: “non-mutually-exclusive” Omniglot

<i>NME Omniglot</i>	20-way 1-shot	20-way 5-shot
MAML	7.8 (0.2)%	50.7 (22.9)%
TAML	9.6 (2.3)%	67.9 (2.3)%
MR-MAML (W) (ours)	<b>83.3 (0.8)%</b>	<b>94.1 (0.1)%</b>

On pose prediction task:



Method	MAML	MR-MAML(W) (ours)	CNP	MR-CNP(W) (ours)
MSE	5.39 (1.31)	<b>2.26 (0.09)</b>	8.48 (0.12)	2.89 (0.18)

(and it’s not just as simple as standard regularization)

CNP	CNP + Weight Decay	CNP + BbB	MR-CNP (W) (ours)
8.48 (0.12)	6.86 (0.27)	7.73 (0.82)	<b>2.89 (0.18)</b>

TAML: Jamal & Qi. Task-Agnostic Meta-Learning for Few-Shot Learning. CVPR’19

Yin, Tucker, Yuan, Levine, Finn. Meta-Learning without Memorization. ’19



# Does meta-regularization lead to better generalization?

Let  $P(\theta)$  be an arbitrary distribution over  $\theta$  that doesn't depend on the meta-training data.

(e.g.  $P(\theta) = \mathcal{N}(\theta; \mathbf{0}, \mathbf{I})$ )

For MAML, with probability at least  $1 - \delta$ ,

$$\underbrace{er(\theta_\mu, \theta_\sigma)}_{\text{generalization error}} \leq \underbrace{\frac{1}{n} \sum_{i=1}^n \hat{e}r(\theta_\mu, \theta_\sigma, \mathcal{D}_i, \mathcal{D}_i^*)}_{\text{error on the meta-training set}} + \left( \sqrt{\frac{1}{2(K-1)}} + \sqrt{\frac{1}{2(n-1)}} \right) \underbrace{\sqrt{D_{KL}(\mathcal{N}(\theta; \theta_\mu, \theta_\sigma) \| P) + \log \frac{n(K+1)}{\delta}}}_{\text{meta-regularization}}, \quad \forall \theta_\mu, \theta_\sigma$$

With a Taylor expansion of the RHS + a particular value of  $\beta \rightarrow$  recover the MR MAML objective.

Proof: draws heavily on Amit & Meier '18

# Want to Learn More?

## CS330: Deep Multi-Task & Meta-Learning

Lecture videos coming out soon!

# Working on Meta-RL?



Try out the Meta-World benchmark

# Collaborators



T Yu, D Quillen, Z He, R Julian, K Hausman, C Finn, S Levine. *Meta-World*. CoRL '19

Yin, Tucker, Yuan, Levine, Finn. *Meta-Learning without Memorization*. '19



