

# Lecture 05: Meta Learning

Readings:

- Survey Paper (posted)

Large, diverse data  
(+ large models)



Broad generalization



Russakovsky et al. '14

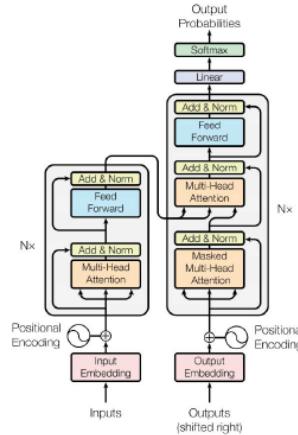


Figure 1: The Transformer - model architecture.

GPT-2

Radford et al. '19

Vaswani et al. '18

Under the paradigm of **supervised learning**.

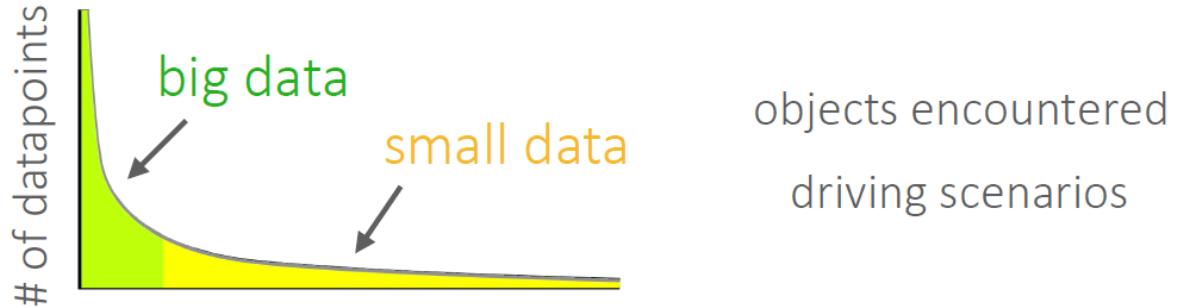
## What if you don't have a large dataset?

medical imaging      robotics      personalized education,  
translation for rare languages      recommendations

## What if you want a general-purpose AI system in the real world?

Need to continuously adapt and learn on the job.  
Learning each thing from scratch won't cut it.

## What if your data has a long tail?



These settings break the supervised learning paradigm.

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

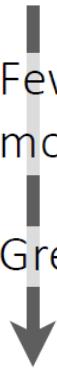
Fine-tuning from ImageNet features

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?

# Problem Statement

# Two ways to view meta-learning

## Mechanistic view

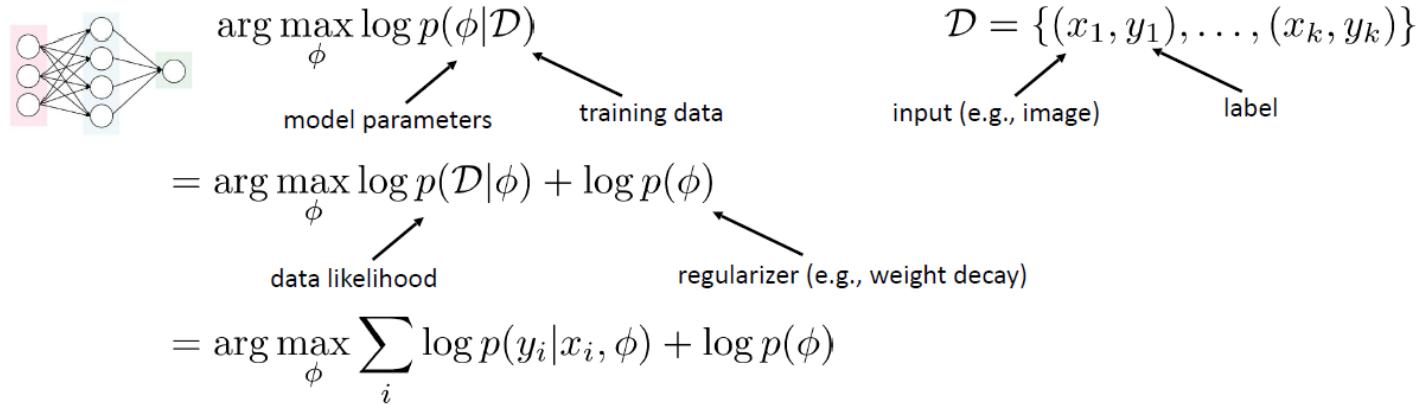
- Deep neural network model that can read in an entire dataset and make predictions for new datapoints
- Training this network uses a meta-dataset, which itself consists of many datasets, each for a different task
- This view makes it easier to implement meta-learning algorithms

## Probabilistic view

- Extract prior information from a set of (meta-training) tasks that allows efficient learning of new tasks
- Learning a new task uses this prior and (small) training set to infer most likely posterior parameters
- This view makes it easier to understand meta-learning algorithms

# Problem definitions

supervised learning:



What is wrong with this?

- The most powerful models typically require large amounts of labeled data
- Labeled data for some tasks may be very limited

# Problem definitions

supervised learning:

$$\arg \max_{\phi} \log p(\phi | \mathcal{D})$$

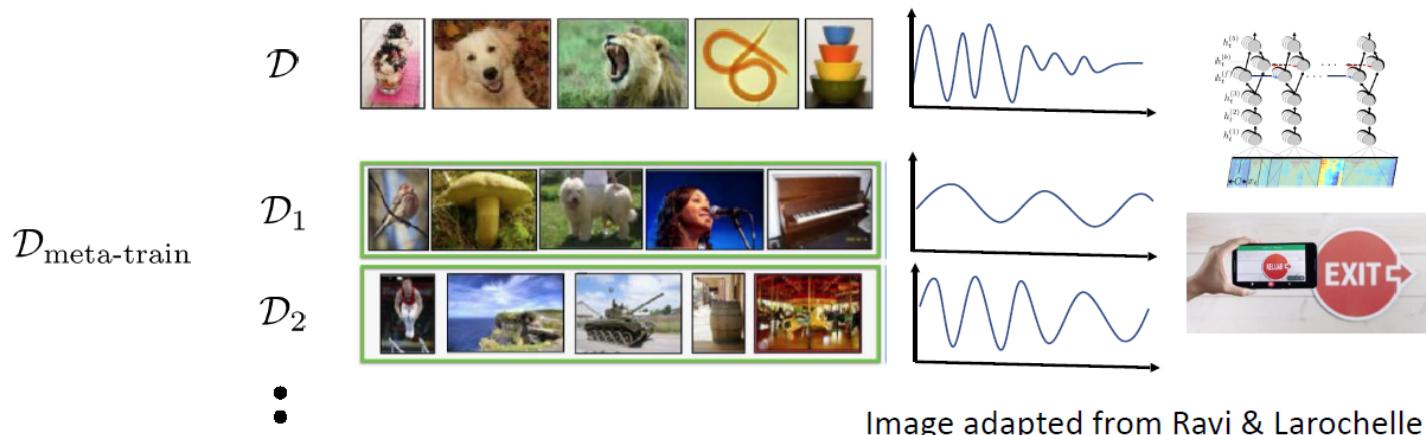
$$\mathcal{D} = \{(x_1, y_1), \dots, (x_k, y_k)\}$$

can we incorporate *additional* data?

$$\arg \max_{\phi} \log p(\phi | \mathcal{D}, \mathcal{D}_{\text{meta-train}})$$

$$\mathcal{D}_{\text{meta-train}} = \{\mathcal{D}_1, \dots, \mathcal{D}_n\}$$

$$\mathcal{D}_i = \{(x_1^i, y_1^i), \dots, (x_k^i, y_k^i)\}$$



# The meta-learning problem

meta-learning:

$$\arg \max_{\phi} \log p(\phi | \mathcal{D}, \mathcal{D}_{\text{meta-train}})$$

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_k, y_k)\}$$

$$\mathcal{D}_{\text{meta-train}} = \{\mathcal{D}_1, \dots, \mathcal{D}_n\}$$

$$\mathcal{D}_i = \{(x_1^i, y_1^i), \dots, (x_k^i, y_k^i)\}$$

what if we don't want to keep  $\mathcal{D}_{\text{meta-train}}$  around forever?

learn *meta-parameters*  $\theta$ :  $p(\theta | \mathcal{D}_{\text{meta-train}})$

whatever we need to know about  $\mathcal{D}_{\text{meta-train}}$  to solve new tasks

$$\begin{aligned} \log p(\phi | \mathcal{D}, \mathcal{D}_{\text{meta-train}}) &= \log \int_{\Theta} p(\phi | \mathcal{D}, \theta) p(\theta | \mathcal{D}_{\text{meta-train}}) d\theta \\ &\approx \log p(\phi | \mathcal{D}, \theta^*) + \log p(\theta^* | \mathcal{D}_{\text{meta-train}}) \end{aligned}$$

$$\arg \max_{\phi} \log p(\phi | \mathcal{D}, \mathcal{D}_{\text{meta-train}}) \approx \arg \max_{\phi} \log p(\phi | \mathcal{D}, \theta^*)$$

this is the meta-learning problem

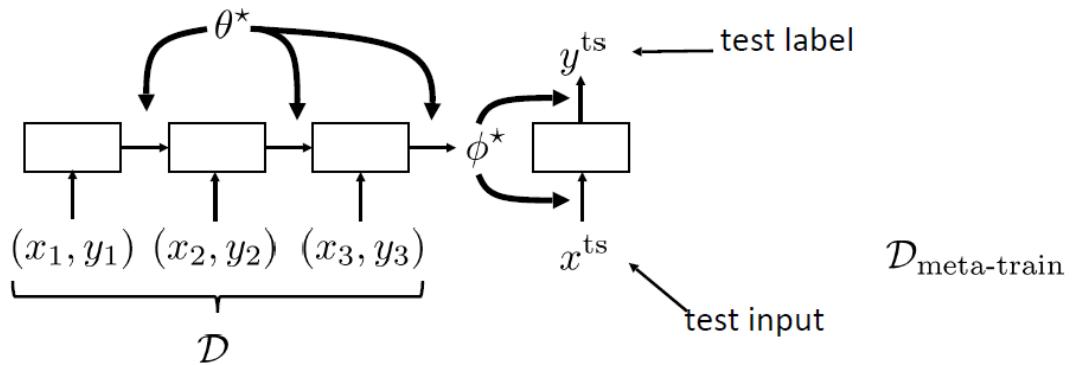
assume  $\phi \perp\!\!\!\perp \mathcal{D}_{\text{meta-train}} | \theta$

$$\theta^* = \arg \max_{\theta} \log p(\theta | \mathcal{D}_{\text{meta-train}})$$

# A Quick Example

meta-learning:  $\theta^* = \arg \max_{\theta} \log p(\theta | \mathcal{D}_{\text{meta-train}})$

adaptation:  $\phi^* = \arg \max_{\phi} \log p(\phi | \mathcal{D}, \theta^*)$



$$\mathcal{D} = \{(x_1, y_1), \dots, (x_k, y_k)\}$$

$$\mathcal{D}_{\text{meta-train}} = \{\mathcal{D}_1, \dots, \mathcal{D}_n\}$$

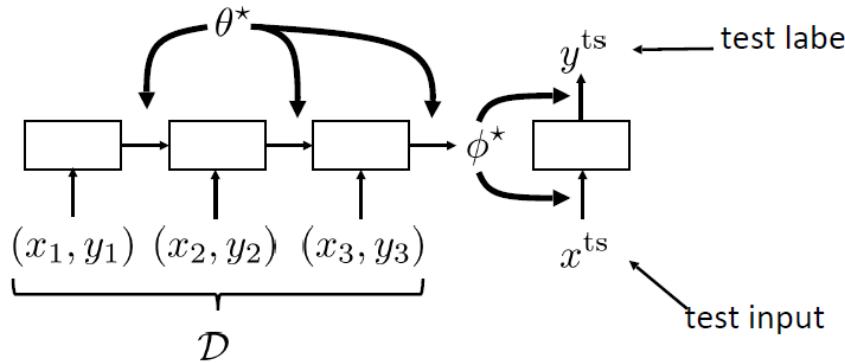
$$\mathcal{D}_i = \{(x_1^i, y_1^i), \dots, (x_k^i, y_k^i)\}$$



# How do we train this thing?

meta-learning:  $\theta^* = \arg \max_{\theta} \log p(\theta | \mathcal{D}_{\text{meta-train}})$

adaptation:  $\phi^* = \arg \max_{\phi} \log p(\phi | \mathcal{D}, \theta^*)$



Key idea:

"our training procedure is based on a simple machine learning principle: test and train conditions must match"

**Vinyals et al., Matching Networks for One-Shot Learning**

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_k, y_k)\}$$

$$\mathcal{D}_{\text{meta-train}} = \{\mathcal{D}_1, \dots, \mathcal{D}_n\}$$

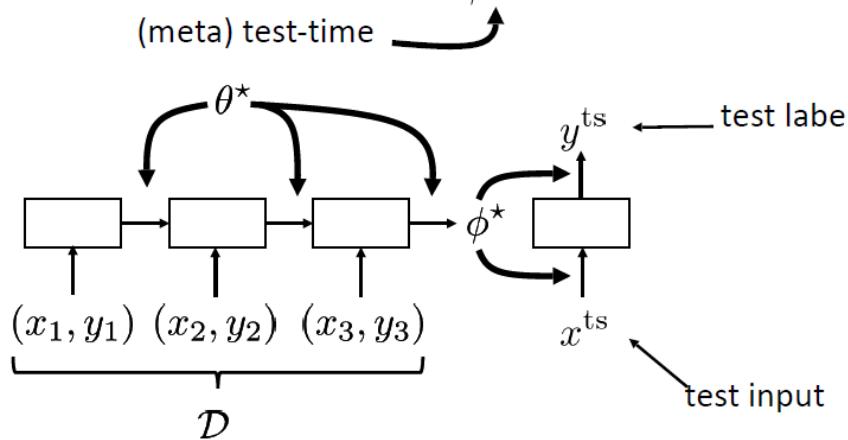
$$\mathcal{D}_i = \{(x_1^i, y_1^i), \dots, (x_k^i, y_k^i)\}$$



# How do we train this thing?

meta-learning:  $\theta^* = \arg \max_{\theta} \log p(\theta | \mathcal{D}_{\text{meta-train}})$

adaptation:  $\phi^* = \arg \max_{\phi} \log p(\phi | \mathcal{D}, \theta^*)$

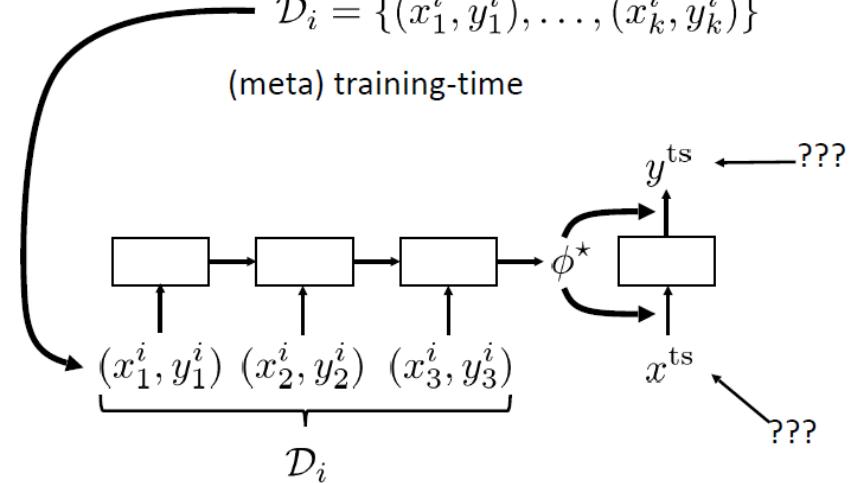


$$\mathcal{D} = \{(x_1, y_1), \dots, (x_k, y_k)\}$$

$$\mathcal{D}_{\text{meta-train}} = \{\mathcal{D}_1, \dots, \mathcal{D}_n\}$$

$$\mathcal{D}_i = \{(x_1^i, y_1^i), \dots, (x_k^i, y_k^i)\}$$

(meta) training-time



Key idea:

"our training procedure is based on a simple machine learning principle: test and train conditions must match"

**Vinyals et al., Matching Networks for One-Shot Learning**

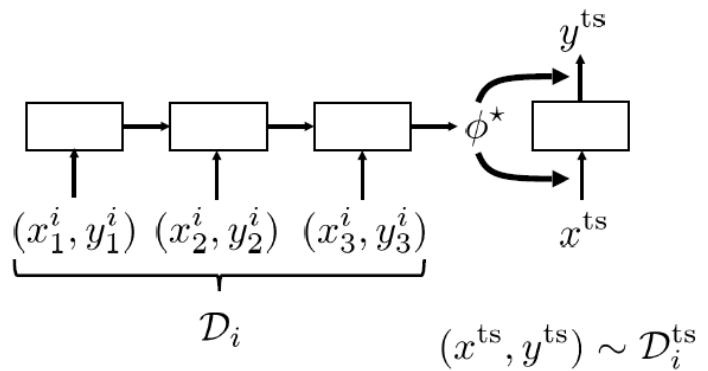
# Reserve a test set for each task!



$$\mathcal{D}_{\text{meta-train}} = \{(\mathcal{D}_1^{\text{tr}}, \mathcal{D}_1^{\text{ts}}), \dots, (\mathcal{D}_n^{\text{tr}}, \mathcal{D}_n^{\text{ts}})\}$$

$$\mathcal{D}_i^{\text{tr}} = \{(x_1^i, y_1^i), \dots, (x_k^i, y_k^i)\}$$

$$\mathcal{D}_i^{\text{ts}} = \{(x_1^i, y_1^i), \dots, (x_l^i, y_l^i)\}$$



Key idea:

"our training procedure is based on a simple machine learning principle: test and train conditions must match"

**Vinyals et al., Matching Networks for One-Shot Learning**

# The complete meta-learning optimization

meta-learning:  $\theta^* = \arg \max_{\theta} \log p(\theta | \mathcal{D}_{\text{meta-train}})$

$$\mathcal{D}_{\text{meta-train}} = \{(\mathcal{D}_1^{\text{tr}}, \mathcal{D}_1^{\text{ts}}), \dots, (\mathcal{D}_n^{\text{tr}}, \mathcal{D}_n^{\text{ts}})\}$$

adaptation:  $\phi^* = \arg \max_{\phi} \log p(\phi | \mathcal{D}^{\text{tr}}, \theta^*)$

$$\mathcal{D}_i^{\text{tr}} = \{(x_1^i, y_1^i), \dots, (x_k^i, y_k^i)\}$$



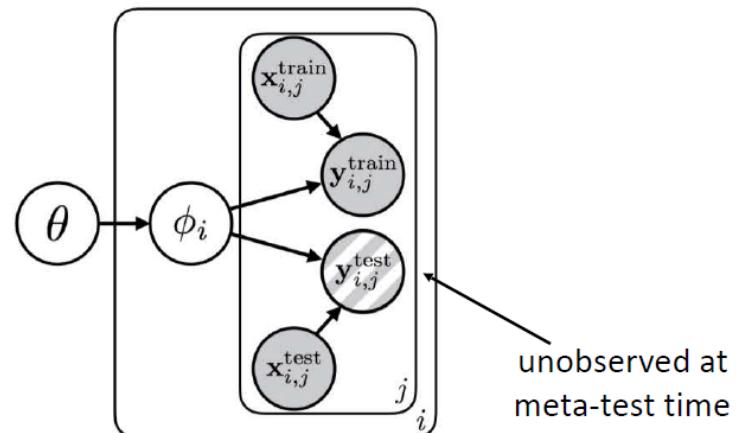
$$\phi^* = f_{\theta^*}(\mathcal{D}^{\text{tr}})$$

$$\mathcal{D}_i^{\text{ts}} = \{(x_1^i, y_1^i), \dots, (x_l^i, y_l^i)\}$$

learn  $\theta$  such that  $\phi = f_{\theta}(\mathcal{D}_i^{\text{tr}})$  is good for  $\mathcal{D}_i^{\text{ts}}$

$$\theta^* = \max_{\theta} \sum_{i=1}^n \log p(\phi_i | \mathcal{D}_i^{\text{ts}})$$

where  $\phi_i = f_{\theta}(\mathcal{D}_i^{\text{tr}})$



# Some meta-learning terminology

learn  $\theta$  such that  $\phi_i = f_\theta(\mathcal{D}_i^{\text{tr}})$  is good for  $\mathcal{D}_i^{\text{ts}}$

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^n \log p(\phi_i | \mathcal{D}_i^{\text{ts}})$$

where  $\phi_i = f_\theta(\mathcal{D}_i^{\text{tr}})$

$$\mathcal{D}_{\text{meta-train}} = \{(\mathcal{D}_1^{\text{tr}}, \mathcal{D}_1^{\text{ts}}), \dots, (\mathcal{D}_n^{\text{tr}}, \mathcal{D}_n^{\text{ts}})\}$$

$$\mathcal{T}_i \left[ \begin{array}{l} \mathcal{D}_i^{\text{tr}} = \{(x_1^i, y_1^i), \dots, (x_k^i, y_k^i)\} \\ \mathcal{D}_i^{\text{ts}} = \{(x_1^i, y_1^i), \dots, (x_l^i, y_l^i)\} \end{array} \right]$$

shot  
(i.e., k-shot, 5-shot)

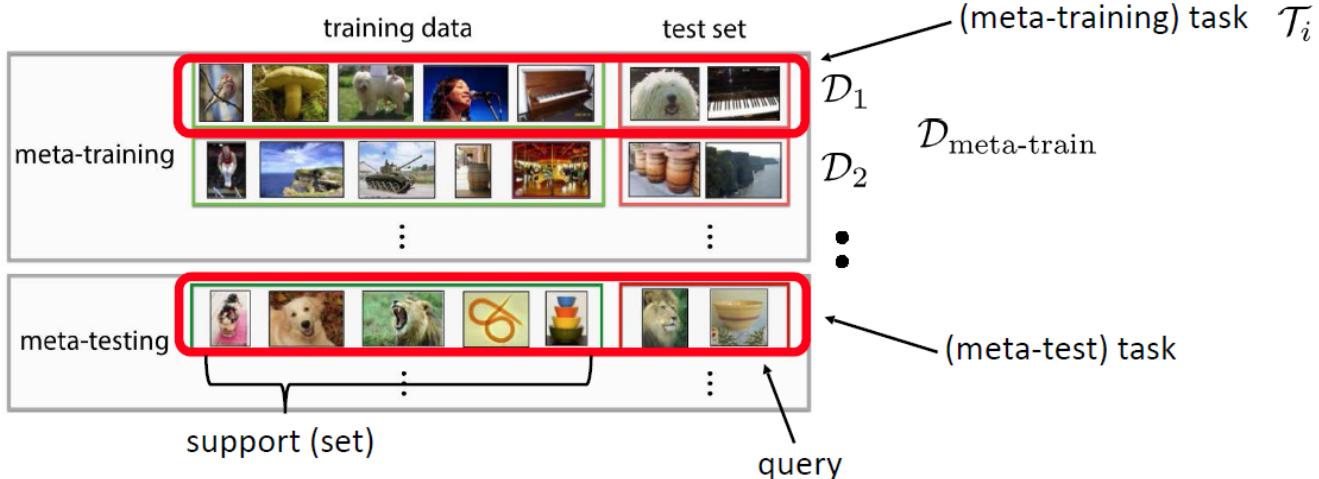


image credit: Ravi & Larochelle '17

# Closely related problem settings

meta-learning:

$$\theta^* = \max_{\theta} \sum_{i=1}^n \log p(\phi_i | \mathcal{D}_i^{\text{ts}})$$

$$\text{where } \phi_i = f_{\theta}(\mathcal{D}_i^{\text{tr}})$$

$$\mathcal{D}_{\text{meta-train}} = \{(\mathcal{D}_1^{\text{tr}}, \mathcal{D}_1^{\text{ts}}), \dots, (\mathcal{D}_n^{\text{tr}}, \mathcal{D}_n^{\text{ts}})\}$$

$$\mathcal{D}_i^{\text{tr}} = \{(x_1^i, y_1^i), \dots, (x_k^i, y_k^i)\}$$

$$\mathcal{D}_i^{\text{ts}} = \{(x_1^i, y_1^i), \dots, (x_l^i, y_l^i)\}$$

multi-task learning: learn model with parameters  $\theta^*$  that solves multiple tasks  $\theta^* = \arg \max_{\theta} \sum_{i=1}^n \log p(\theta | \mathcal{D}_i)$   
can be seen as special case where  $\phi_i = \theta$  (i.e.,  $f_{\theta}(\mathcal{D}_i) = \theta$ )

hyperparameter optimization & auto-ML: can be cast as meta-learning

hyperparameter optimization:  $\theta$  = hyperparameters,  $\phi$  = network weights

architecture search:  $\theta$  = architecture,  $\phi$  = network weights

very active area of research! but outside the scope of this tutorial

# Outline

- Motivation
- Problem Statement
- Meta-Learning Approaches
  - Black-Box Adaptation
  - Non-Parametric Methods
- Applications

# General recipe

# How to evaluate a meta-learning algorithm

the Omniglot dataset Lake et al. Science 2015

1623 characters from 50 different alphabets

Hebrew				
וְ	בַּ	כִּ	לְ	בָּ
אֵ	גַּ	מִ	תְּ	בָּ
חֶ	קְ	צָ	אַ	בָּ
בָּ	הַ	כָּ	לְ	בָּ
רַ				

ପ୍ରତ୍ଯେକିଆନ ତାଣ୍ଡଳୀ  
 ଓ କମ୍ପ୍ୟୁଟର ଏଟିବ  
 ଦଶ ମଧ୍ୟ ଏହି ଡାକ୍ତରି  
 ପଛଭାବୁ ମନ୍ଦିର  
 ଓ ତେଜି ଅଞ୍ଚଳ ଉଥି  
 ଚଗ୍ର ଚାଲ ଆମ୍ବାଦୁ  
 ଓ ଫର୍ମ ବାବା

C	L	B	G	Z
H	A	K	X	V
U	Theta	Y	I	O
W	Pi	Nu	O	E
R	Epsilon	Zeta	Psi	

ବ୍ୟାକୁ ପାତା ମାତ୍ର ଏହିଲା  
କାହାର କାହାର କାହାର  
କାହାର କାହାର କାହାର  
କାହାର କାହାର କାହାର

many classes, few examples

the “transpose” of MNIST

statistics more reflective  
of the real world

20 instances of each character

Proposes both few-shot discriminative & few-shot generative problems

## Initial few-shot learning approaches w/ Bayesian models, non-parametrics

Fei-Fei et al. '03 Lake et al. '11 Salakhutdinov et al. '12 Lake et al. '13

Other datasets used for few-shot image recognition: Minilmagenet, CIFAR, CUB, CelebA, others

# General recipe

How to *evaluate* a meta-learning algorithm

5-way, 1-shot image classification (Minilmagenet)

Given 1 example of 5 classes:



Classify new examples



held-out classes

$\mathcal{T}_1$



meta-training

$\mathcal{T}_2$



training classes

⋮

⋮

any ML  
problem

, Can replace image classification with: regression, language generation, skill learning,

# General recipe

**How to *design* a meta-learning algorithm**

1. Choose a form of  $p(\phi_i | \mathcal{D}_i^{\text{tr}}, \theta)$
2. Choose how to optimize  $\theta$  w.r.t. max-likelihood objective using  $\mathcal{D}_{\text{meta-train}}$

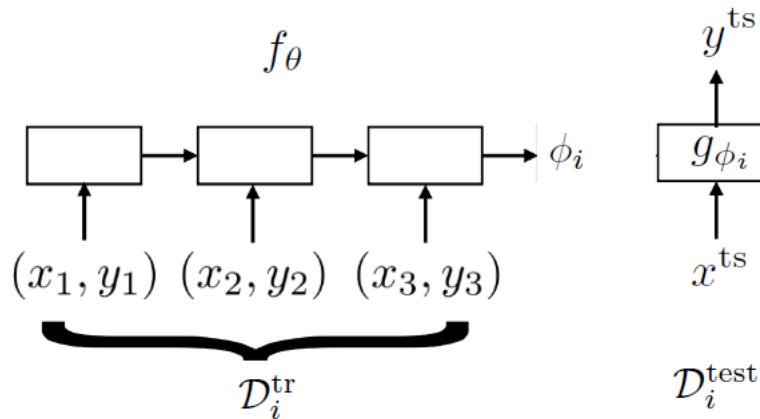
Can we treat  $p(\phi_i | \mathcal{D}_i^{\text{tr}}, \theta)$  as an **inference** problem?

Neural networks are good at inference.

# Black-Box Adaptation

**Key idea:** Train a neural network to represent  $p(\phi_i | \mathcal{D}_i^{\text{tr}}, \theta)$

For now: Use deterministic (point estimate)  $\phi_i = f_\theta(\mathcal{D}_i^{\text{tr}})$



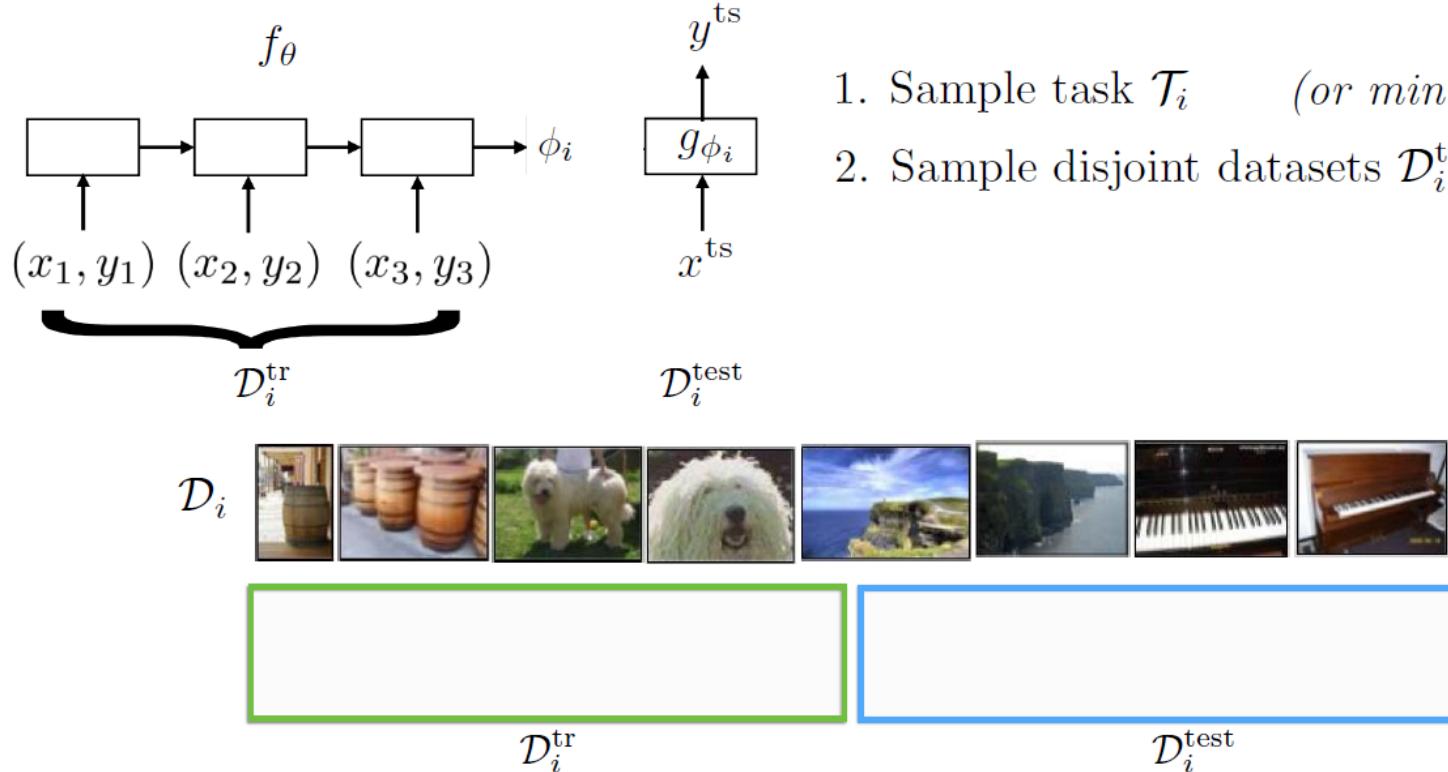
Train with standard supervised learning!

$$\max_{\theta} \sum_{\mathcal{T}_i} \underbrace{\sum_{(x,y) \sim \mathcal{D}_i^{\text{test}}} \log g_{\phi_i}(y|x)}_{\mathcal{L}(\phi_i, \mathcal{D}_i^{\text{test}})}$$

$$\max_{\theta} \sum_{\mathcal{T}_i} \mathcal{L}(f_\theta(\mathcal{D}_i^{\text{tr}}), \mathcal{D}_i^{\text{test}})$$

# Black-Box Adaptation

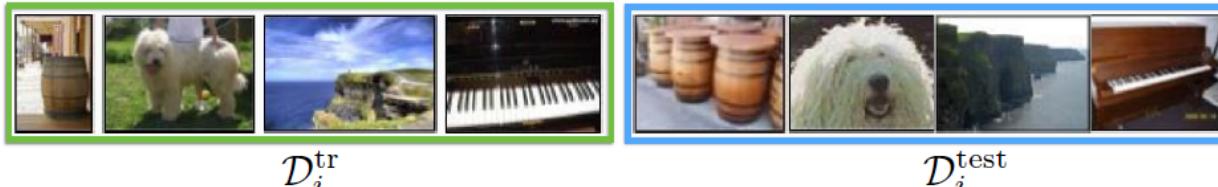
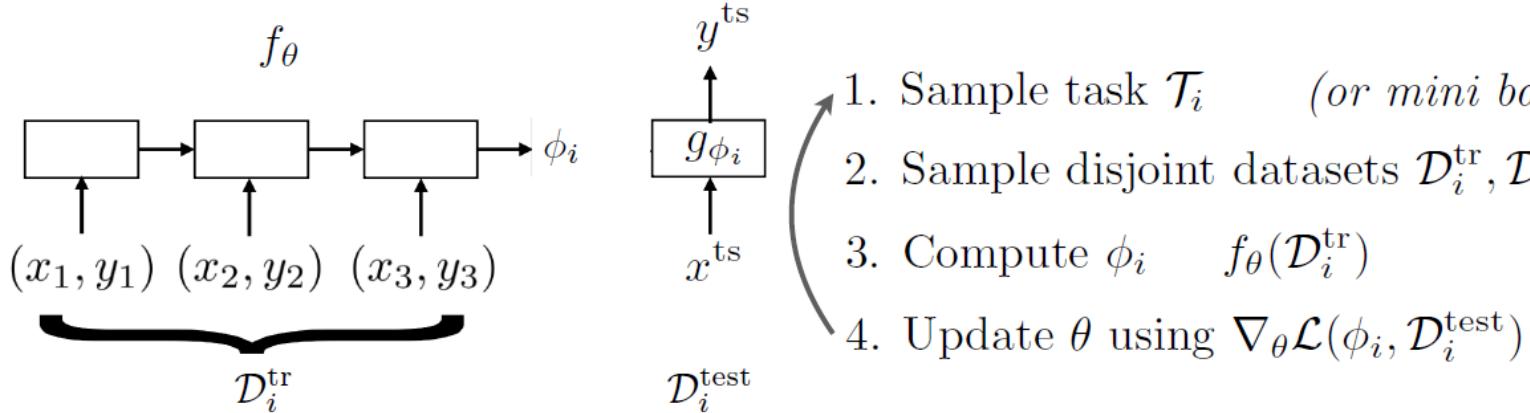
**Key idea:** Train a neural network to represent  $p(\phi_i | \mathcal{D}_i^{\text{tr}}, \theta)$



1. Sample task  $\mathcal{T}_i$  (or mini batch of tasks)
2. Sample disjoint datasets  $\mathcal{D}_i^{\text{tr}}, \mathcal{D}_i^{\text{test}}$  from  $\mathcal{D}_i$

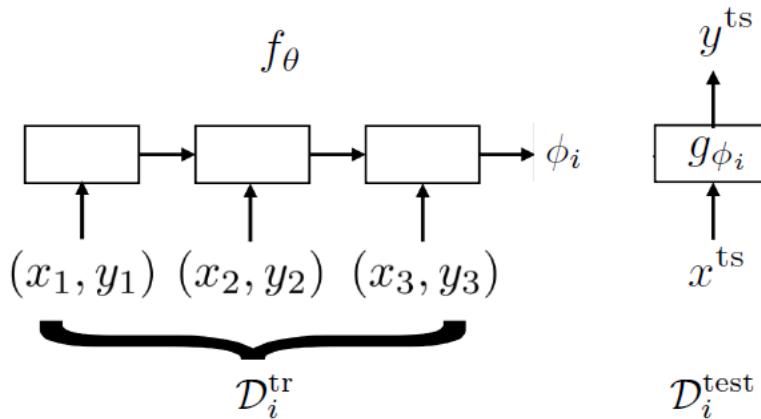
# Black-Box Adaptation

**Key idea:** Train a neural network to represent  $p(\phi_i | \mathcal{D}_i^{\text{tr}}, \theta)$



# Black-Box Adaptation

**Key idea:** Train a neural network to represent  $p(\phi_i | \mathcal{D}_i^{\text{tr}}, \theta)$



**Form of  $f_\theta$ ?**

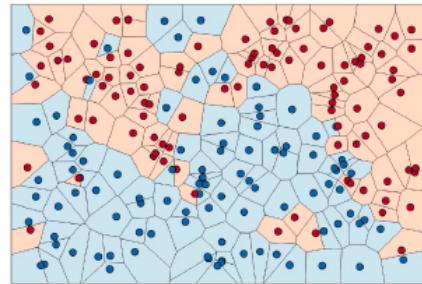
- LSTM
- Neural turing machine (NTM)
- Self-attention
- 1D convolutions
- feedforward + average

# Outline

- Motivation
- Problem Statement
- Meta-Learning Approaches
  - Black-Box Adaptation
  - Non-Parametric Methods
- Applications

**So far:** Learning parametric models.

In low data regimes, **non-parametric** methods are simple, work well.



During **meta-test time**: few-shot learning <-> low data regime

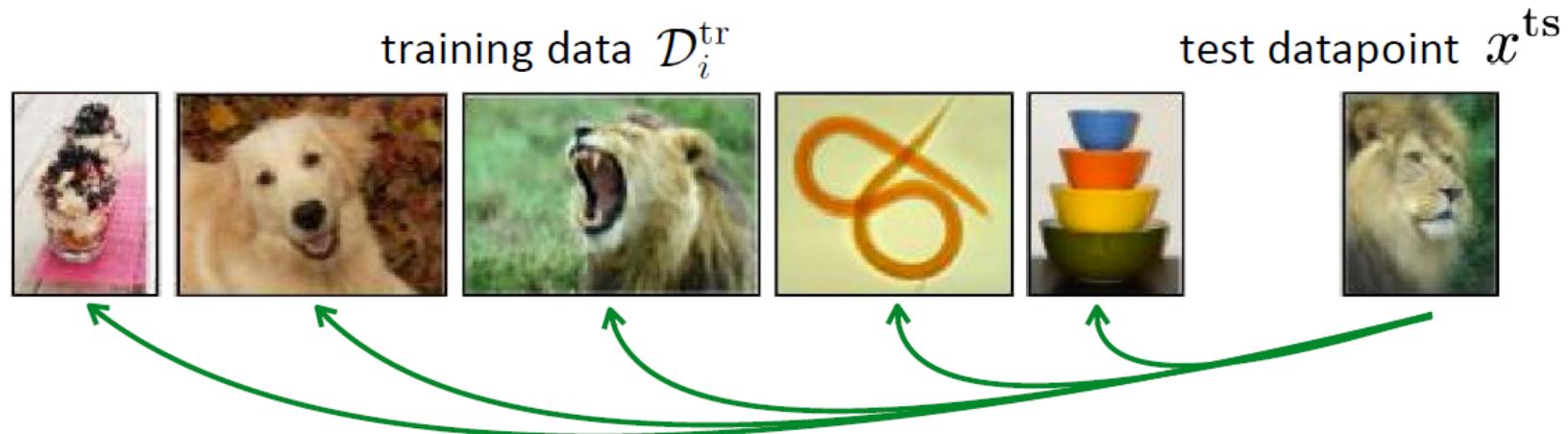
During **meta-training**: still want to be parametric

Can we use **parametric meta-learners** that produce effective **non-parametric learners**?

Note: some of these methods precede parametric approaches

# Non-parametric methods

**Key Idea:** Use non-parametric learner.



In what space do you compare? With what distance metric?

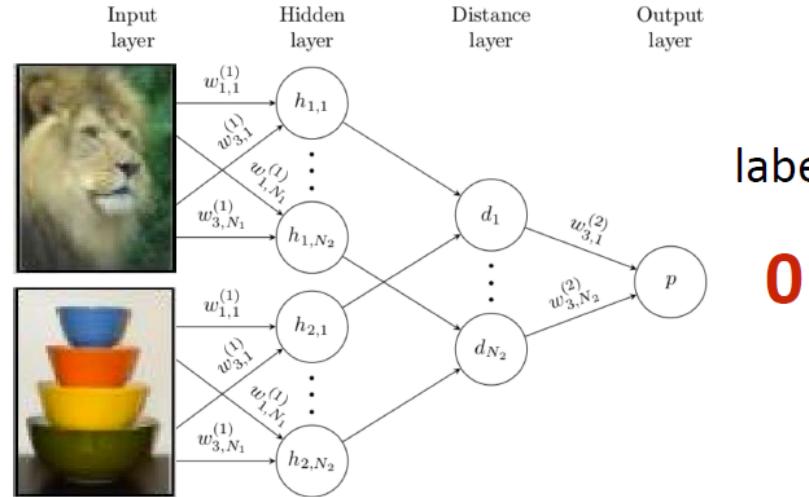
~~pixel space,  $\ell_2$ -distance?~~

Learn to compare using data!

# Non-parametric methods

**Key Idea:** Use non-parametric learner.

train Siamese network to predict whether or not two images are the same class

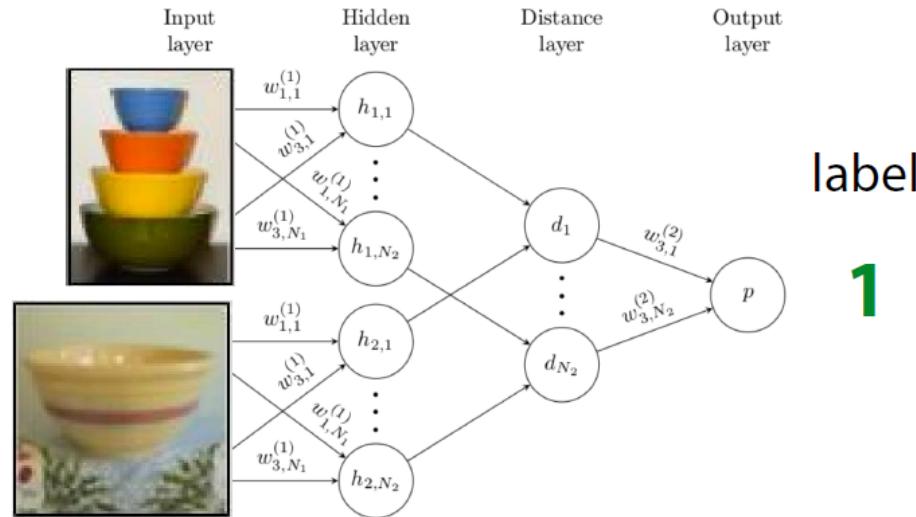


Koch et al., ICML '15

# Non-parametric methods

**Key Idea:** Use non-parametric learner.

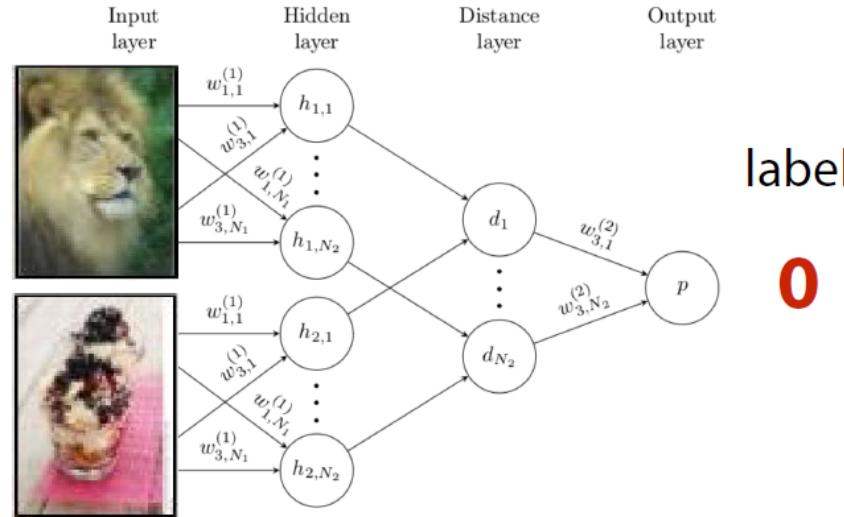
train Siamese network to predict whether or not two images are the same class



# Non-parametric methods

**Key Idea:** Use non-parametric learner.

train Siamese network to predict whether or not two images are the same class

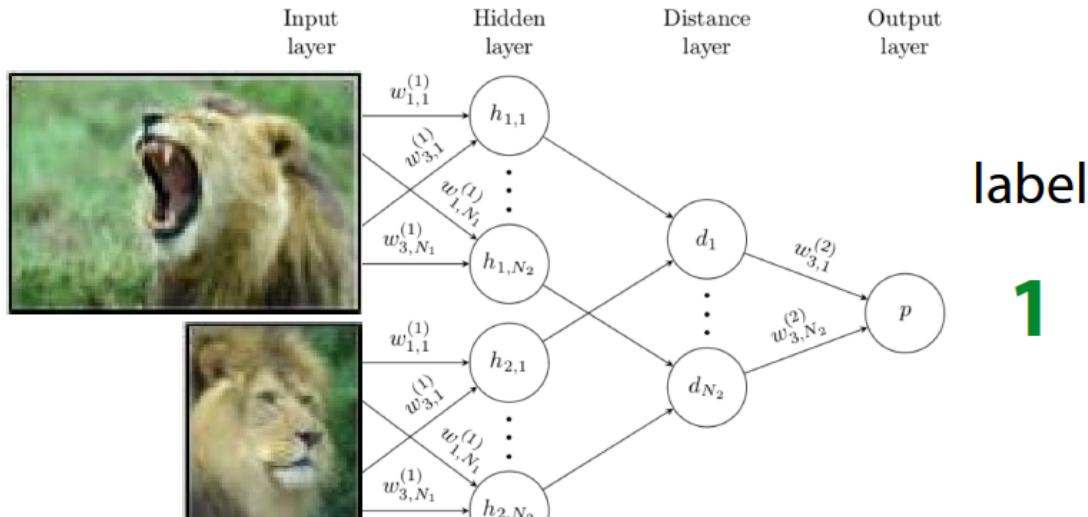


Koch et al., ICML '15

# Non-parametric methods

**Key Idea:** Use non-parametric learner.

train Siamese network to predict whether or not two images are the same class



Meta-test time: compare image  $\mathbf{x}_{\text{test}}$  to each image in  $\mathcal{D}_j^{\text{tr}}$

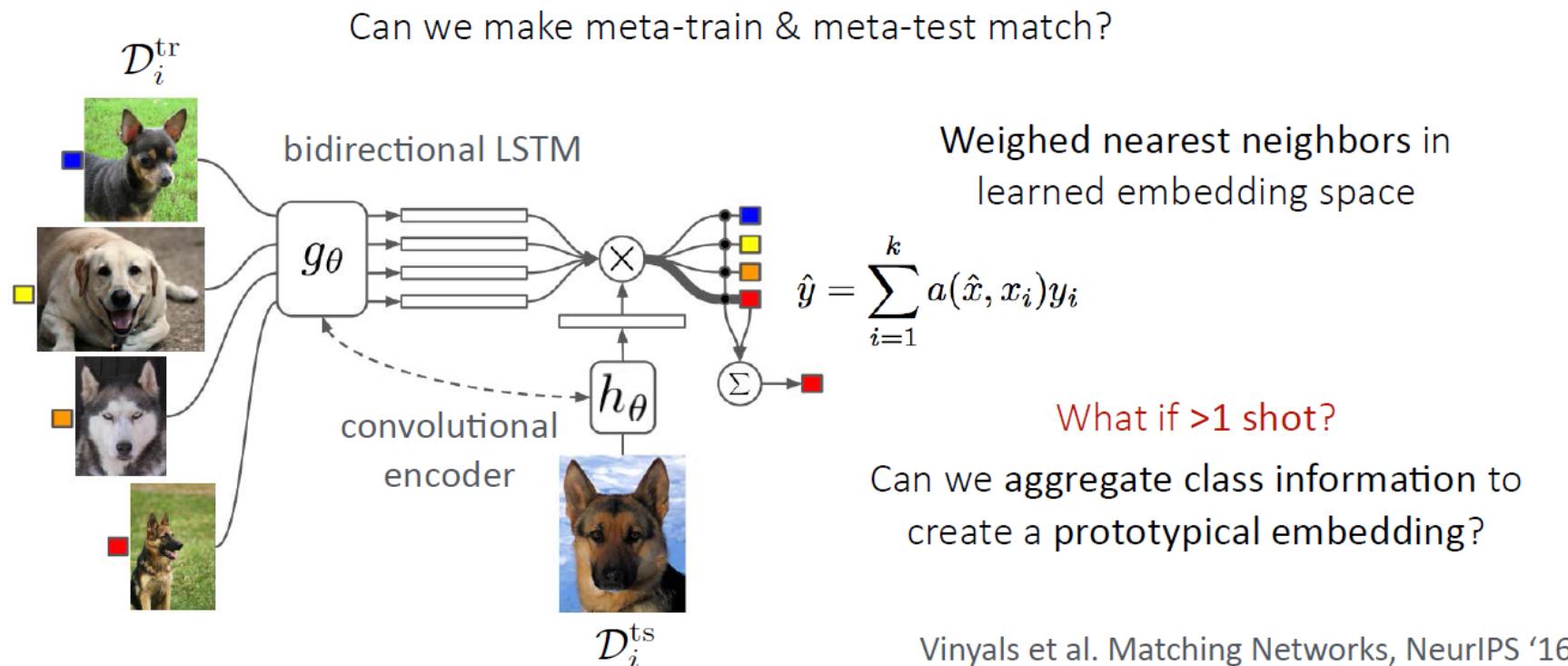
Meta-training: 2-way classification

Can we **match** meta-train & meta-test?

Meta-test: N-way classification

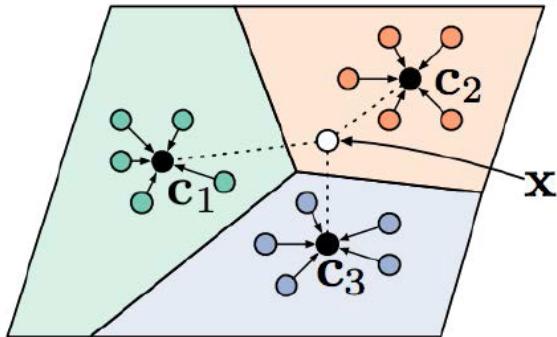
# Non-parametric methods

**Key Idea:** Use non-parametric learner.



# Non-parametric methods

**Key Idea:** Use non-parametric learner.



(a) Few-shot

$$\mathbf{c}_k = \frac{1}{|\mathcal{D}_i^{\text{tr}}|} \sum_{(x,y) \in \mathcal{D}_i^{\text{tr}}} f_\theta(x)$$

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

$d$ : Euclidean, or cosine distance

Snell et al. Prototypical Networks, NeurIPS '17

# Non-parametric methods

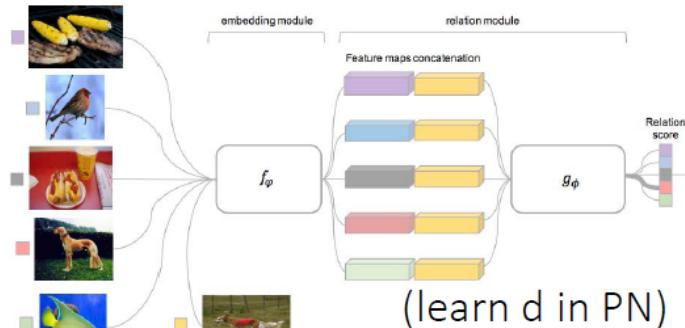
**So far:** Siamese networks, matching networks, prototypical networks

Embed, then nearest neighbors.

## Challenge

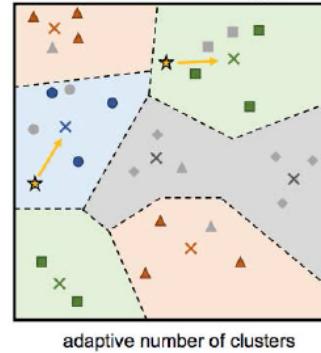
What if you need to reason about more complex relationships between datapoints?

**Idea:** Learn non-linear relation module on embeddings



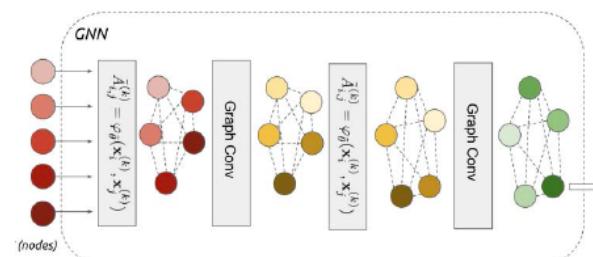
Sung et al. Relation Net

**Idea:** Learn infinite mixture of prototypes.



Allen et al. IMP, ICML '19

**Idea:** Perform message passing on embeddings



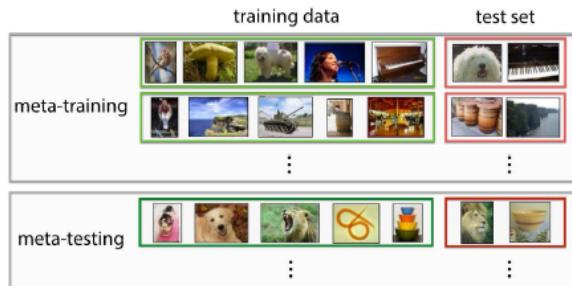
Garcia & Bruna, GNN

# Outline

- Motivation
- Problem Statement
- Meta-Learning Approaches
  - Black-Box Adaptation
  - Non-Parametric Methods
- Applications

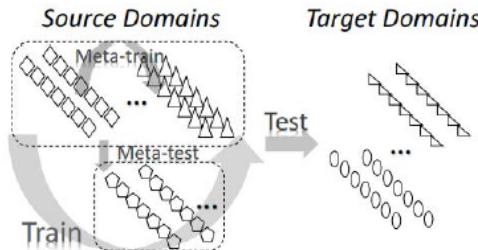
# Applications in computer vision

## few-shot image recognition



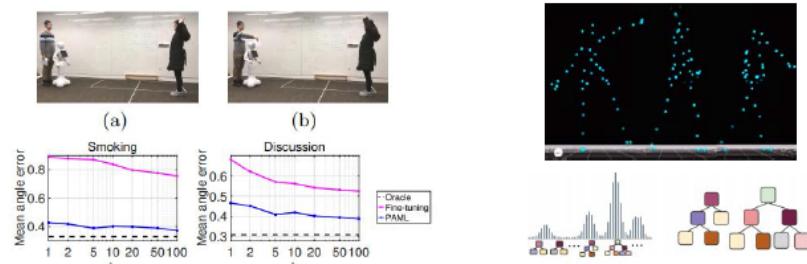
see, e.g.: Vinyals et al. **Matching Networks for One Shot Learning**, and many many others

## domain adaptation



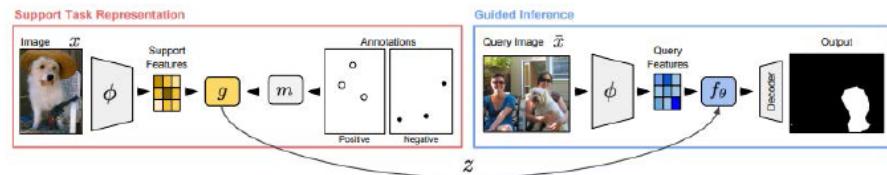
see, e.g.: Li, Yang, Song, Hospedales. **Learning to Generalize: Meta-Learning for Domain Adaptation**.

## human motion and pose prediction



see, e.g.: Gui et al. **Few-Shot Human Motion Prediction via Meta-Learning**.  
Alet et al. **Modular Meta-Learning**.

## few-shot segmentation



see, e.g.: Shaban, Bansal, Liu, Essa, Boots. **One-Shot Learning for Semantic Segmentation**.  
Rakelly, Shelhamer, Darrell, Efros, Levine. **Few-Shot Segmentation Propagation with Guided Networks**.  
Dong, Xing. **Few-Shot Semantic Segmentation with Prototype Learning**.

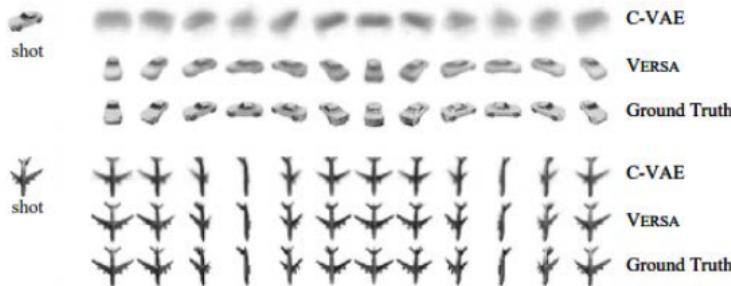
# Applications in image & video generation

## few-shot image generation



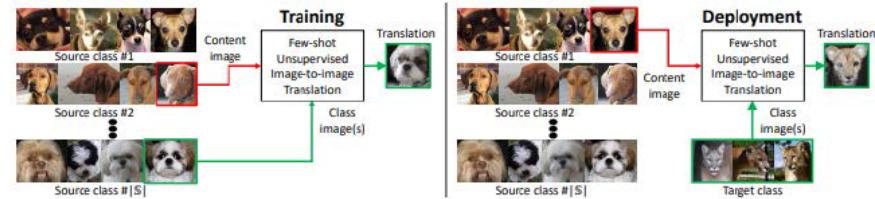
see, e.g.: Reed, Chen, Paine, van den Oord, Eslami, Rezende, Vinyals, de Freitas. **Few-Shot Autoregressive Density Estimation**. and many many others.

## generation of novel viewpoints



see, e.g.: Gordon, Bronskill, Bauer, Nowozin, Turner. **VERSA: Versatile and Efficient Few-Shot Learning**.

## few-shot image-to-image translation



see, e.g.: Liu, Huang, Mallya, Karras, Aila, Lehtinen, Kautz. **Few-Shot Unsupervised Image-to-Image Translation**.

## generating talking heads from images



see, e.g.: Zakharov, Shysheya, Burkov, Lempitsky. **Few-Shot Adversarial Learning of Realistic Neural Talking Head Models**

# Meta-Learning for *Language*

## Adapting to *new programs*

Meta Program Induction  
Learn new program from a few I/O examples.

*Devlin\*, Bunel\* et al. NeurIPS '17*

### Program Synthesis

Question:  
How many CFL teams are from York College?

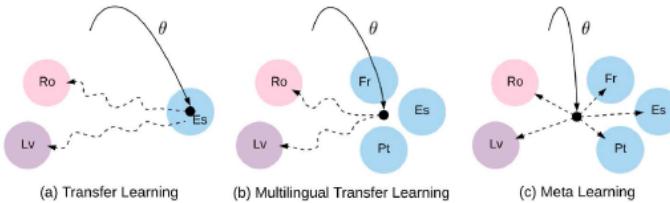
SQL:  
`SELECT COUNT CFL Team FROM CFLDraft WHERE College = "York"`

Result:  
2

Construct pseudo-tasks with relevance function

*Huang et al. NAACL '18*

## Adapting to *new languages* Low-Resource Neural Machine Translation



Learn to translate new language pair w/o a lot of paired data?  
*Gu et al. EMNLP '18*

## Learning *new words*

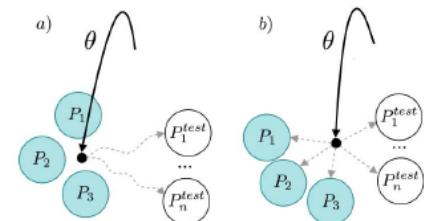
### One-Shot Language Modeling

Learn how to use a new word from one example usage.

*Vinyals et al. Matching Networks, '16*

## Adapting to *new personas*

### Personalizing Dialogue Agents



Adapt dialogue to a persona with a few examples

*Lin\*, Madotto\* et al. ACL '19*