

# CS7150 Deep Learning

Jiaji Huang

<https://jiaji-huang.github.io>

03/09/2024

# Recap of Last Lecture

- Parameter Efficient FineTuning (PEFT)
- Pretrained LMs already solve new Tasks to some extent
  - Prompt engineering and **zero/few-shot In-context Learning**

## Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 cheese => ..... ← prompt
```

## Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← examples
3 peppermint => menthe poivrée ←
4 plush girafe => girafe peluche ←
5 cheese => ..... ← prompt
```

# Recap of Last Lecture

Yet finetuning is still necessary!

Supervised finetuning (SFT)

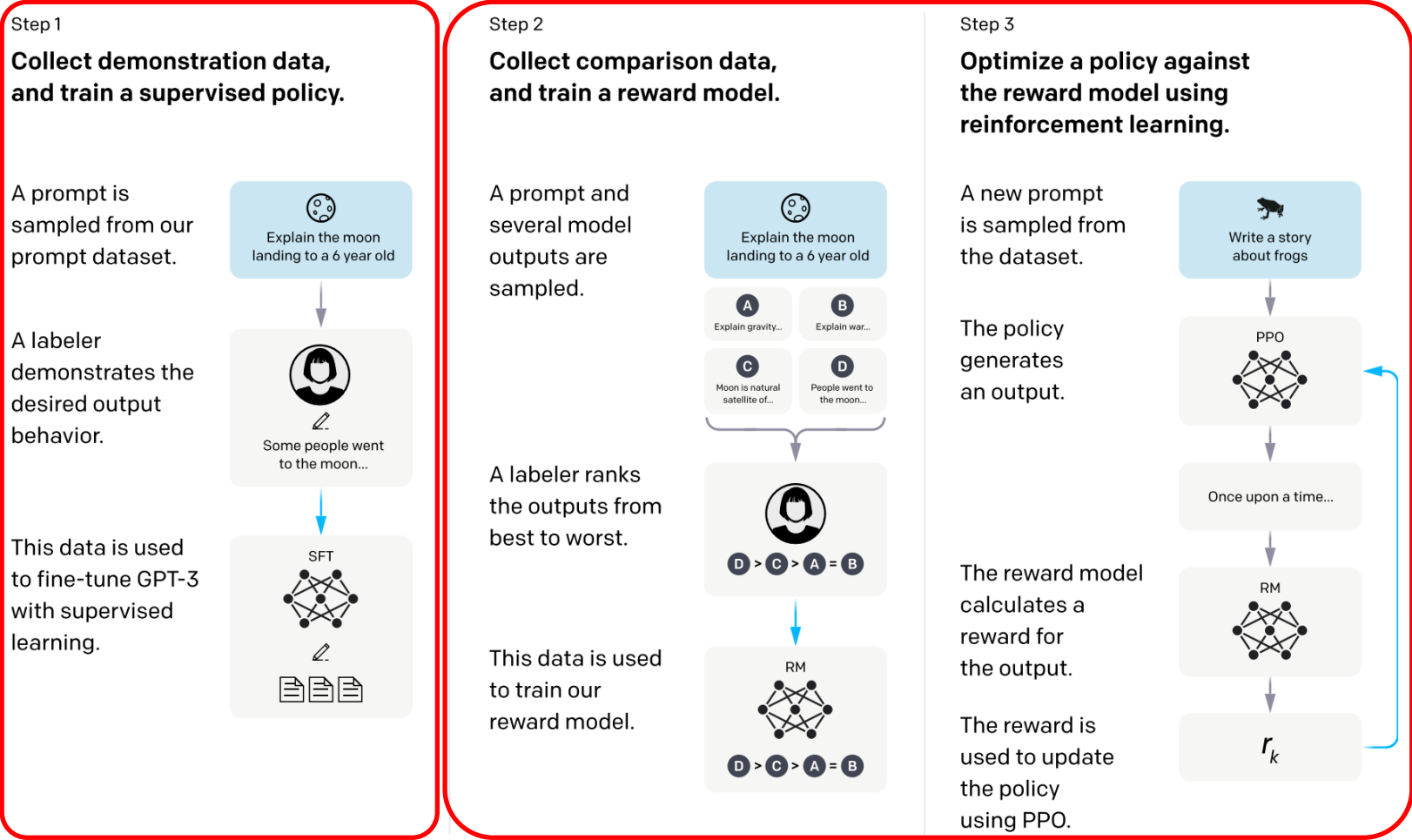
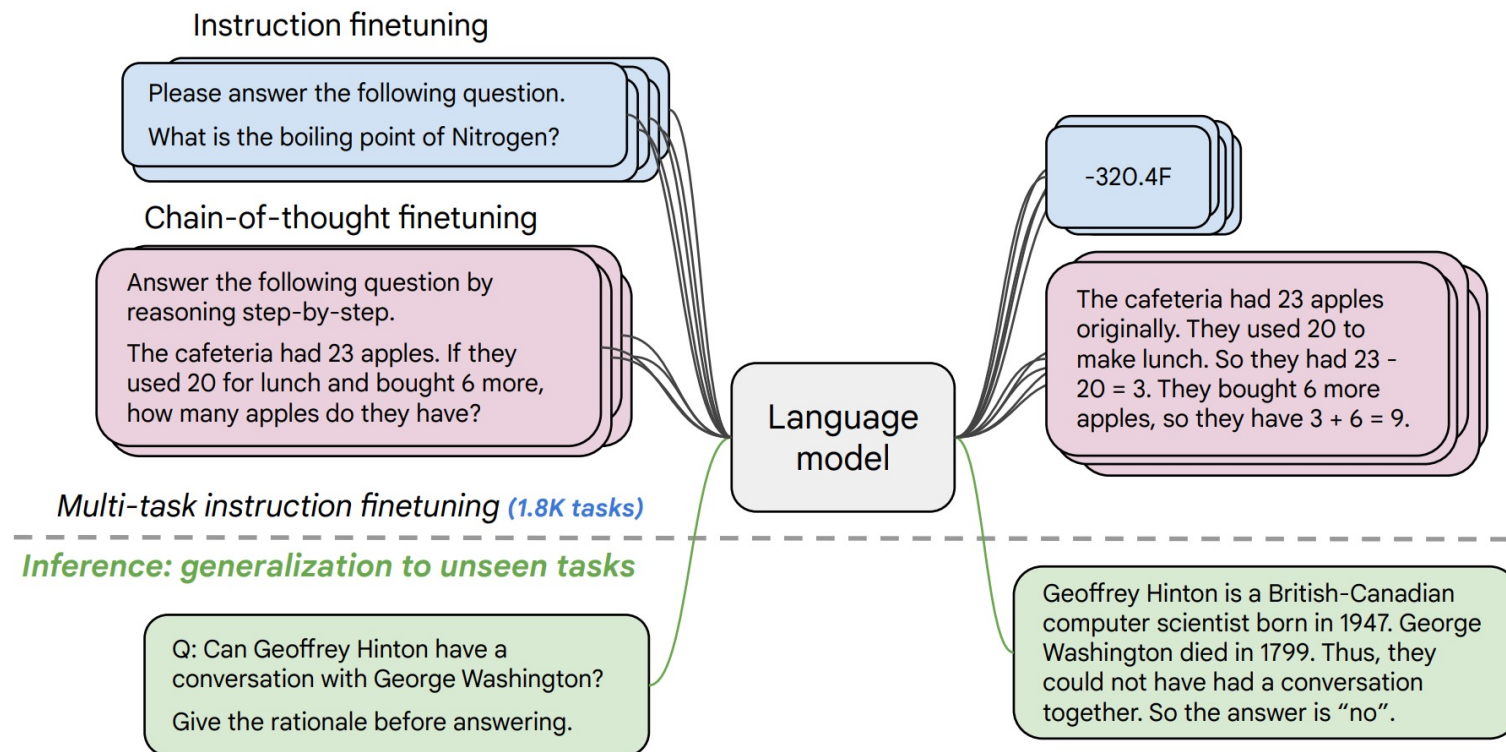


Figure from [Ouyang et. al, 2022](#)

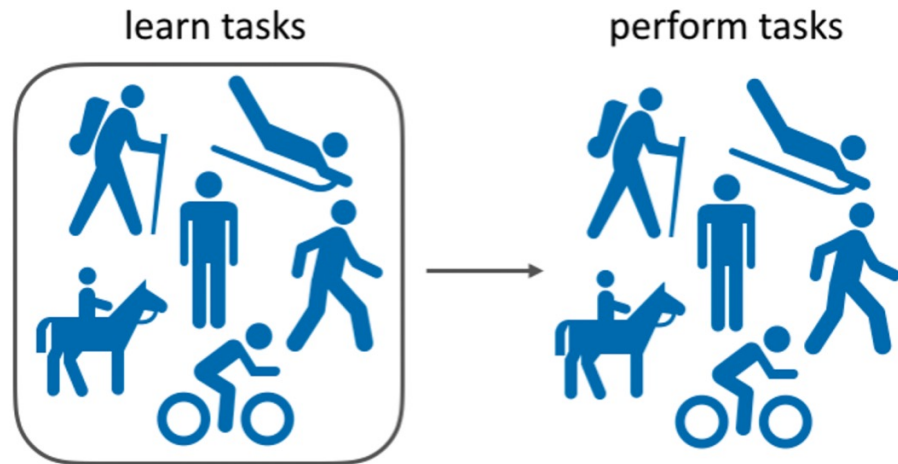
# Recap of Last Lecture: on SFT

- Instruction finetuning and FLAN (**multi-task training objective**)
- Seeing many tasks helps for solving a new task (**meta Learning**)



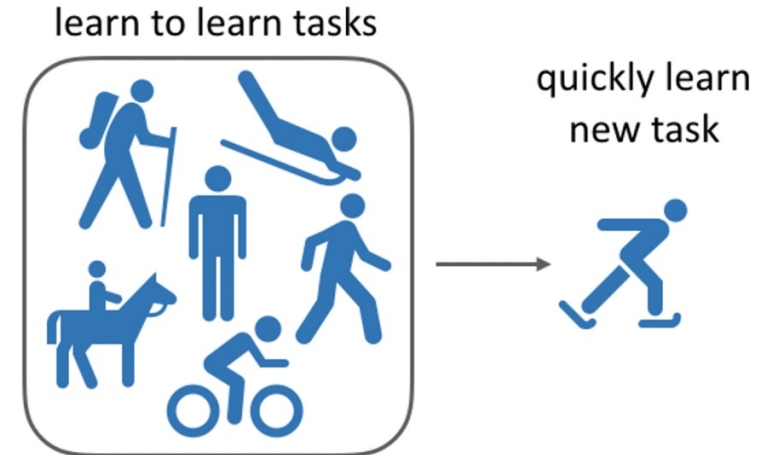
# Multi-task Learning vs Meta Learning

- Multi-task Learning



- Setting: Test tasks = Training tasks
- Goal: master this set of tasks

- Meta Learning



- Setting: Test task(s)  $\notin$  Training task
- Goal: Adapt to unseen task(s) quickly

# Agenda for Today

- Multi-task Learning (MTL)
- Meta Learning
- Zero-shot Learning

# Formalize: Defining Tasks

- A task has
  - Input  $\mathbf{x} \sim p(\mathbf{x})$
  - Target output  $\mathbf{y}$  given  $\mathbf{x}$ , draw from  $p(\mathbf{y}|\mathbf{x})$
  - $\mathcal{T} \triangleq (p(\mathbf{x}), p(\mathbf{y}|\mathbf{x}))$
- Example: Different  $p(\mathbf{x})$ 
  - Scene image classification v.s. medical image classification
- Example: Same  $p(\mathbf{x})$  but different  $p(\mathbf{y}|\mathbf{x})$ 
  - Scene classification:  $\mathbf{x}$  scene images,  $\mathbf{y}$  scene label
  - Object detection from Scene image:  $\mathbf{y}$  object bounding box

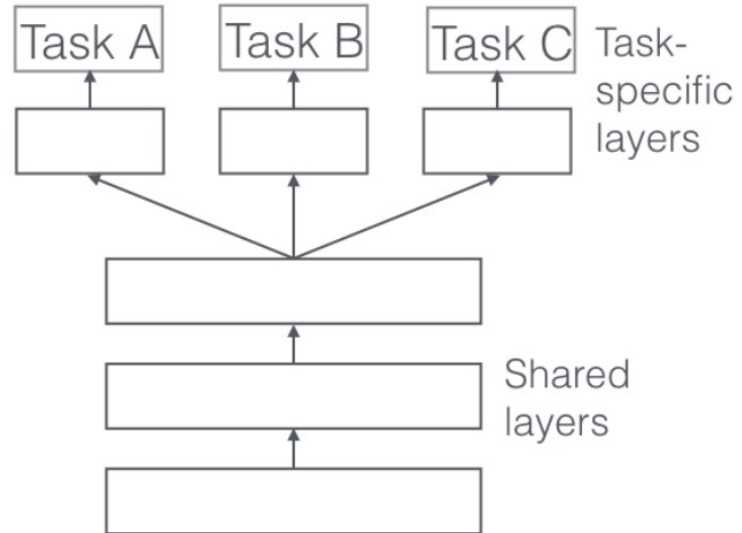
# Formalize: Multi-Task Learning (MTL)

- $\mathcal{J}_i \triangleq (p_i(\mathbf{x}), p_i(\mathbf{y}|\mathbf{x})), i = 1, \dots, T$
- Training data  $\mathcal{D}_i^{tr}$ , testing data  $\mathcal{D}_i^{te}$  draw from each  $\mathcal{J}_i$
- Train on  $\mathcal{D}_i^{tr}$  ( $i = 1, \dots, T$ ) and test on each  $\mathcal{D}_i^{te}$
- Assumption: the tasks are **relevant**
- Otherwise, we may just train a model for each task

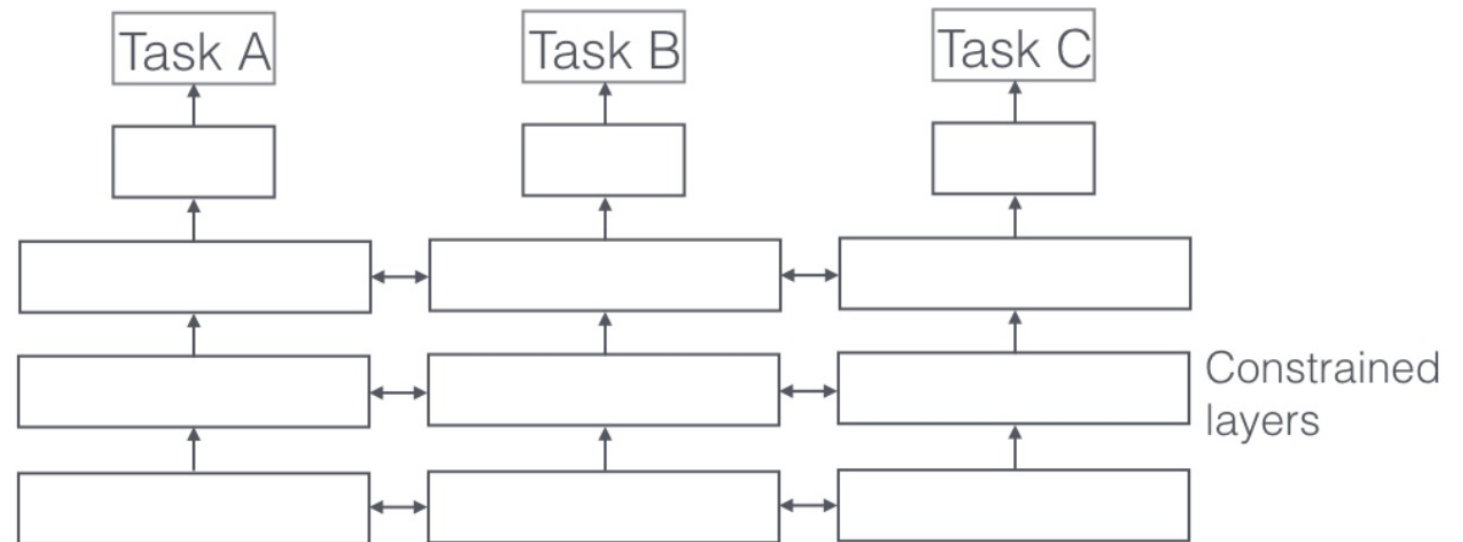


# Sharing Model Parameters for MTL

- Hard sharing

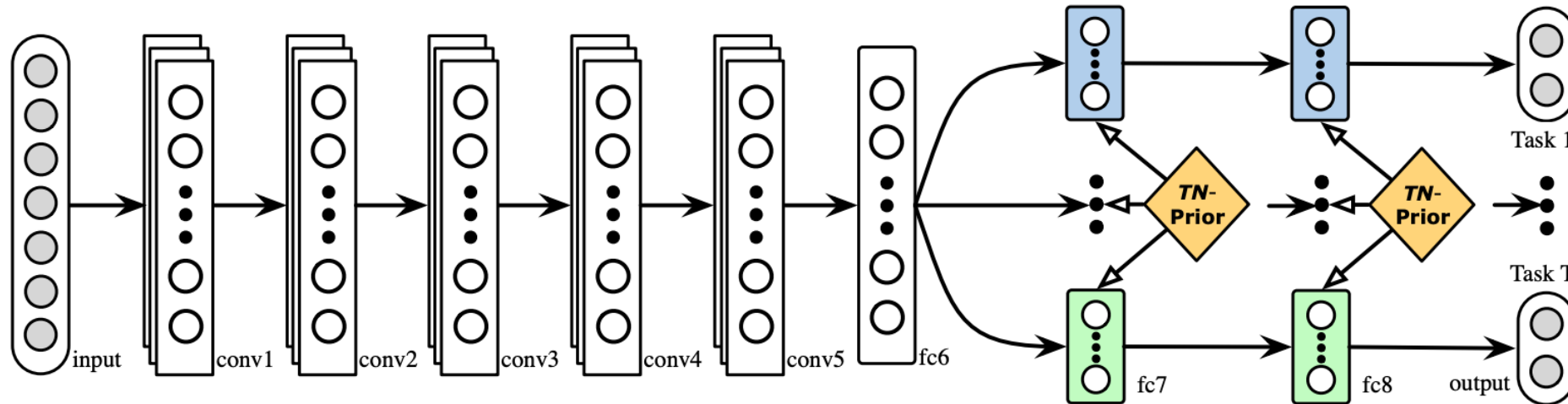


- Soft sharing



# Example of Hard Parameter Sharing

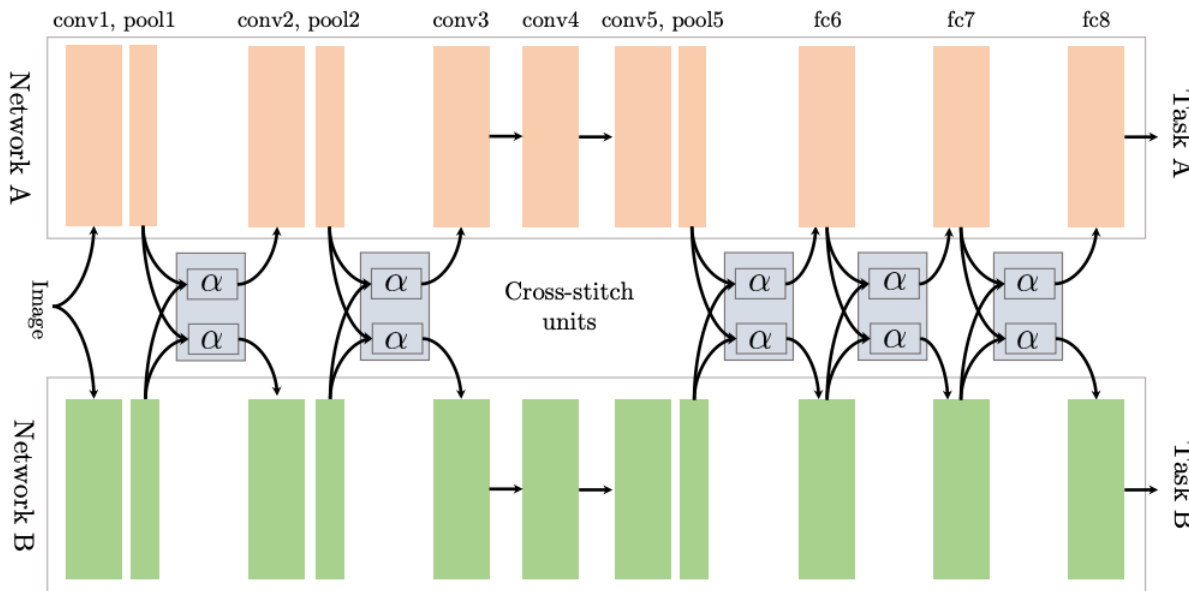
- Deep Relation Network ([Long and Wang, 2015](#))



- Share conv layers
- Prior on Fc7, fc8's weight matrices: encodes task relationship

# Example of Soft Parameter Sharing

- Cross-stitch Network ([Misra et. al, 2016](#))
- Start from two networks (same architecture) for two tasks
- Learn linear combinator  $\alpha$  for feature maps



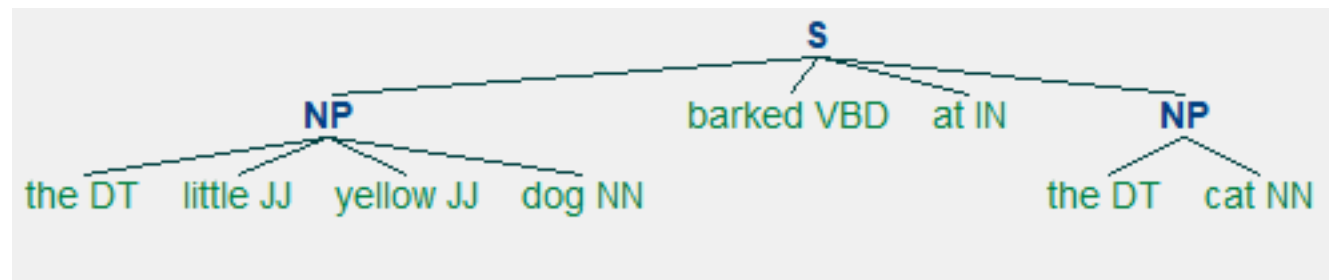
$$\begin{bmatrix} \tilde{x}_A^{ij} \\ \tilde{x}_B^{ij} \end{bmatrix} = \underbrace{\begin{bmatrix} \alpha_{AA} & \alpha_{AB} \\ \alpha_{BA} & \alpha_{BB} \end{bmatrix}}_{\text{Encodes our knowledge of task relevance}} \begin{bmatrix} x_A^{ij} \\ x_B^{ij} \end{bmatrix}$$

Encodes our knowledge of task relevance

# What Parameters/Layers to be Shared

- Common to share the bottom layers, with task-specific “head”
- Sometimes a task is more fundamental than the others ([Søgaard and Goldberg, 2016](#))
- E.g., Chunking works on top of POS (part of speech) tags

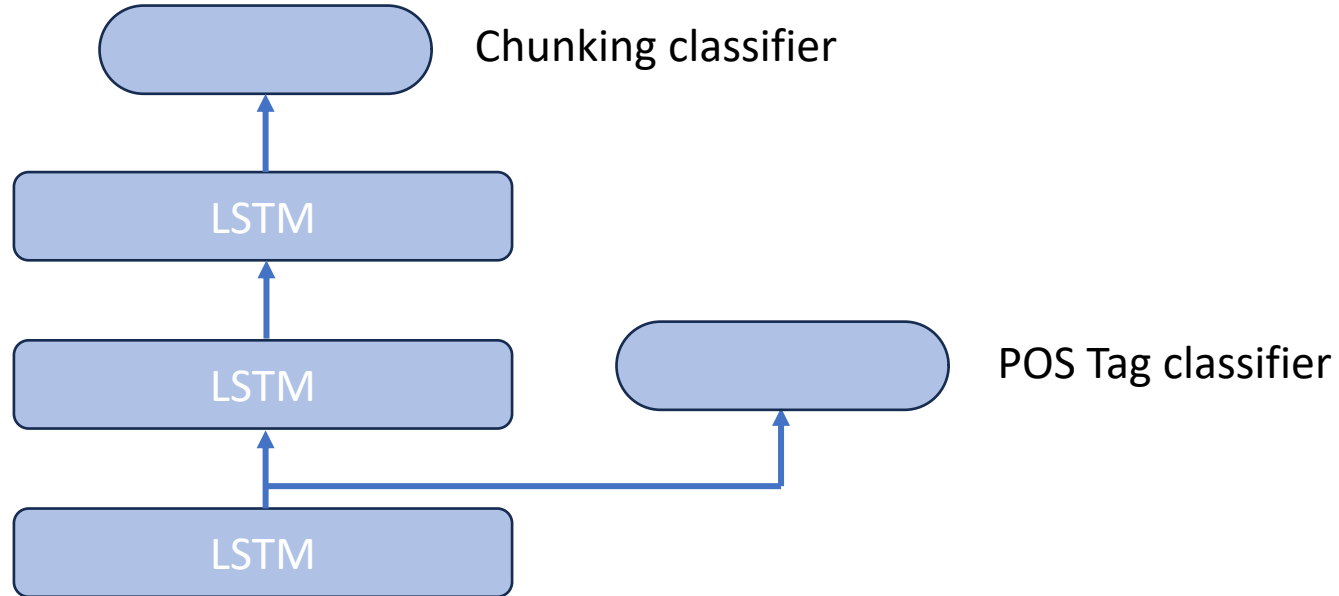
POS Tag	Abbr.	words
Determiner	DT	a, an, the, this
Adjective	JJ	big, kind, cool, ...
Noun	NN	dog, cat
preposition	IN	at, into, over, ...
Verb (past tense)	VBD	walked, talked, ...



Example from [medium](#)

# Jointly learn chunking with POS tagging

- Discussion: how do we share the parameters/layers?



# Jointly learn chunking with POS tagging

- Results

	LAYERS		DOMAINS			
	CHUNKS	POS	BROADCAST (6)	BC-NEWS (8)	MAGAZINES (1)	WEBLOGS (6)
BI-LSTM	3	-	88.98	91.84	90.09	90.36
	3	3	88.91	91.84	90.95	90.43
	3	1	<b>89.48</b>	<b>92.03</b>	<b>91.53</b>	<b>90.78</b>

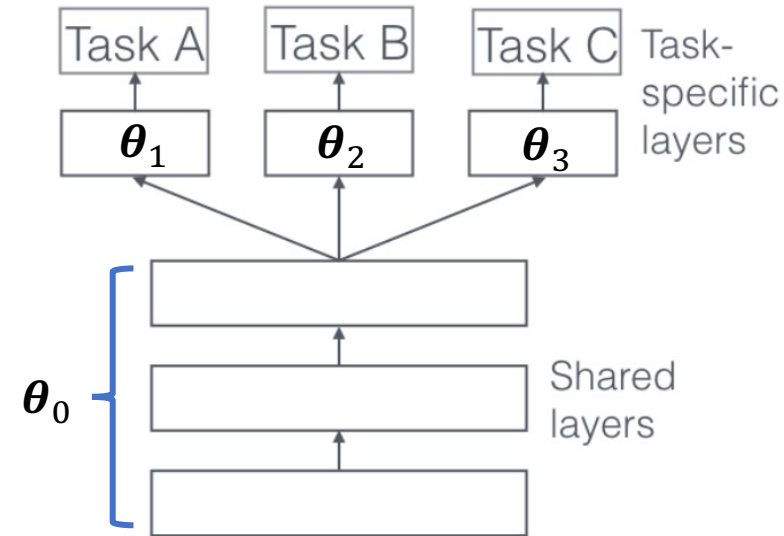
- More helpful to use low-level task at lower layer

# MTL Objective functions

- $\theta_0$ : shared parameter,  $\theta_i$ : task-specific parameter
- Commonly seen, additive

$$\min_{\theta_0, \dots, \theta_T} \sum_{i=1}^T w_i \left\{ \mathcal{L}_i \triangleq \sum_{(x,y) \in \mathcal{D}_i^{tr}} \ell_i(\theta_0, \theta_i; x, y) \right\}$$

- $w_i$ : importance of the  $i$ -th task
- $w_i$  such that tasks with similar gradient magnitude ([Chen et. al, 2018](#))



# Optimize the Objective

$$\min_{\boldsymbol{\theta}_0, \dots, \boldsymbol{\theta}_T} \sum_{i=1}^T w_i \left\{ \mathcal{L}_i \triangleq \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_i^{tr}} \ell_i(\boldsymbol{\theta}_0, \boldsymbol{\theta}_i; \mathbf{x}, \mathbf{y}) \right\}$$

- Sample a minibatch of tasks, indices  $\mathcal{J} \subseteq \{1, \dots, T\}$
- For each task  $i \in \mathcal{J}$ , sample a batch of  $(\mathbf{x}, \mathbf{y})$ 's, denoted as  $\mathcal{X}_i \subseteq \mathcal{D}_i^{tr}$
- Compute (stochastic) loss

$$\hat{\mathcal{L}} = \sum_{i \in \mathcal{J}} w_i \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{X}_i} \ell_i(\boldsymbol{\theta}_0, \boldsymbol{\theta}_i; \mathbf{x}, \mathbf{y})$$

- Back-prop to compute gradients,  $\frac{\partial \hat{\mathcal{L}}}{\partial \boldsymbol{\theta}_0}$  and  $\frac{\partial \hat{\mathcal{L}}}{\partial \boldsymbol{\theta}_i}$  ( $i \in \mathcal{J}$ )
- Update the  $\boldsymbol{\theta}_0$  and  $\boldsymbol{\theta}_i$  's with Adam, etc.



# Potential Issues

- Choice of  $w_i$  can be tricky
- Tasks may compete (negative transfer), i.e.,  
$$\mathcal{L}_1(\boldsymbol{\theta}_0, \boldsymbol{\theta}_1) < \mathcal{L}_1(\overline{\boldsymbol{\theta}}_0, \overline{\boldsymbol{\theta}}_1)$$

but

$$\mathcal{L}_2(\boldsymbol{\theta}_0, \boldsymbol{\theta}_2) > \mathcal{L}_2(\overline{\boldsymbol{\theta}}_0, \overline{\boldsymbol{\theta}}_2)$$

Improve for task 1, but harm task 2

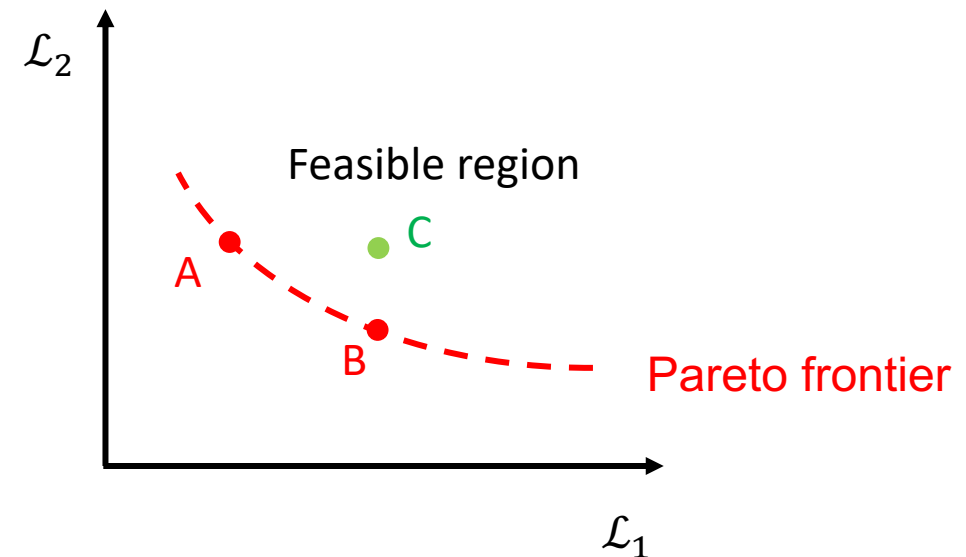
# Multi-objective MTL ([Sener and Koltun, 2018](#))

$$\min_{\theta_0, \theta_1, \dots, \theta_T} \{\mathcal{L}_1(\theta_0, \theta_1), \dots, \mathcal{L}_T(\theta_0, \theta_T)\}$$

- Pareto optimality:

$(\theta_0^*, \theta_1^*, \dots, \theta_T^*)$  is Pareto optimal if any other  $(\theta_0, \theta_1, \dots, \theta_T)$  harms at least one task

- Optimize so we arrive onto the **Pareto frontier**



# Optimizer for the Multi-objective MTL

**Pareto stationary point**  $(\boldsymbol{\theta}_0, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_T)$  satisfies (**KKT condition**):

- For task-specific parameters:

$$\nabla_{\boldsymbol{\theta}_i} \mathcal{L}_i(\boldsymbol{\theta}_0, \boldsymbol{\theta}_i) = 0 \text{ for all task } i = 1, \dots, m$$

- For shared parameters:

Exist  $w_1, \dots, w_T \geq 0$  where  $\sum_{i=1}^T w_i = 1$  such that

$$\sum_{i=1}^T w_i \nabla_{\boldsymbol{\theta}_0} \mathcal{L}_i(\boldsymbol{\theta}_0, \boldsymbol{\theta}_i) = 0$$

# Optimizer for the Multi-objective MTL

While not converged:

Update task specific  $\theta_i \leftarrow \theta_i - \eta_i \nabla_{\theta_i} \mathcal{L}_i(\theta_0, \theta_i)$

Solve for  $w_1, \dots, w_T \geq 0$  where  $\sum_{i=1}^T w_i = 1$  such that

$$\min_{w_1, \dots, w_T} \left\| \sum_{i=1}^T w_i \nabla_{\theta_0} \mathcal{L}_i(\theta_0, \theta_i) \right\|^2$$

Update  $\theta_0 \leftarrow \theta_0 - \eta \sum_{i=1}^T w_i \nabla_{\theta_0} \mathcal{L}_i(\theta_0, \theta_i)$

- Note: we can convert the above to stochastic gradient descent

# Results

- Experiment on [multiMNIST dataset](#)

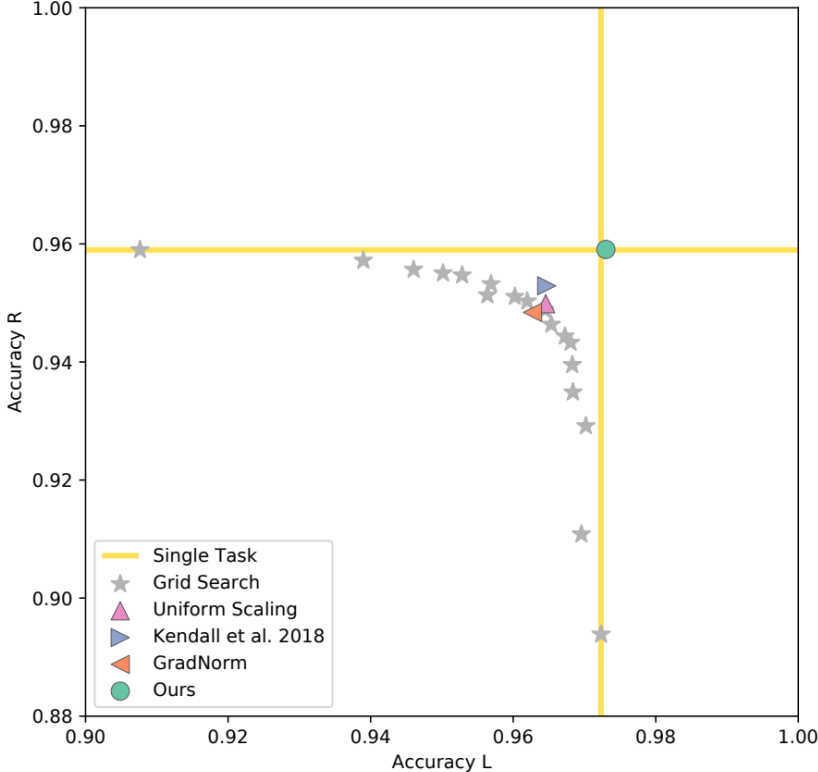
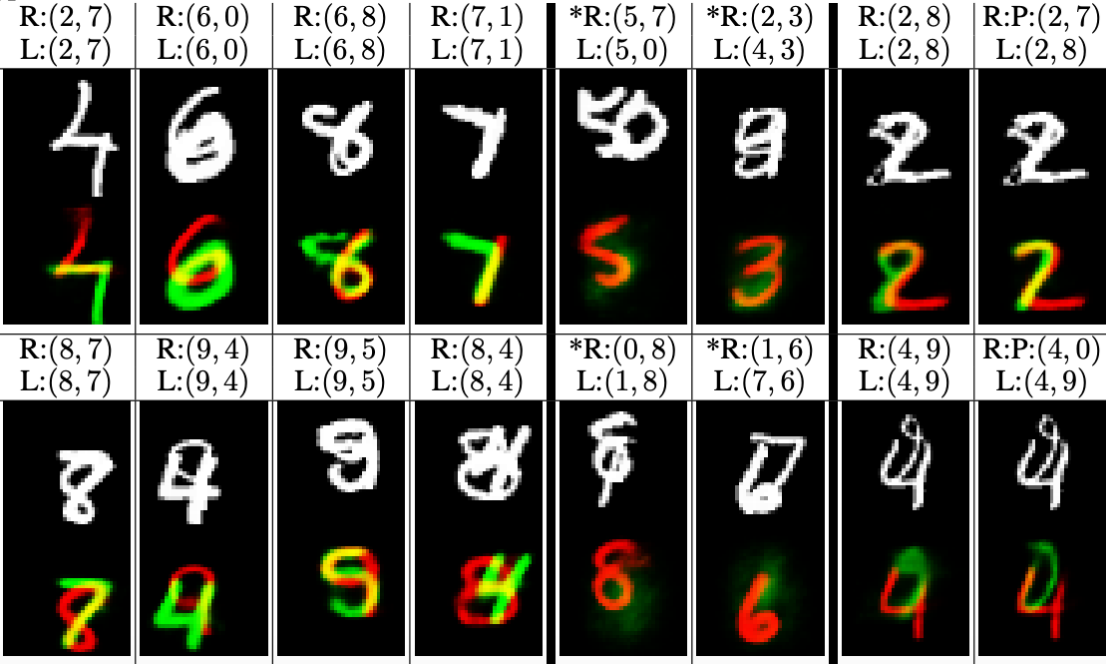


Figure 3: **MultiMNIST accuracy profile.** We plot the obtained accuracy in detecting the left and right digits for all baselines. The grid-search results suggest that the tasks compete for model capacity. Our method is the only one that finds a solution that is as good as training a dedicated model for each task. Top-right is better.

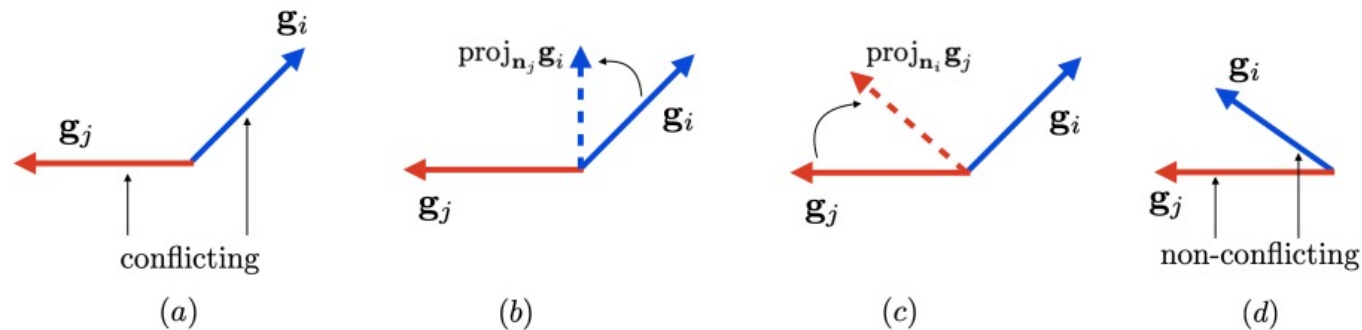
Fig. 3 from ([Sener and Koltun, 2018](#))

# Diagnose Negative Transfer via Gradients

- Again Consider additive objective

$$\min_{\theta_0, \dots, \theta_T} \sum_{i=1}^T w_i \mathcal{L}_i$$

- Remove conflicting components ([Yu, et. al, 2022](#))



# Results

	% accuracy	
task specific, 1-fc [46]	42	} Naïve MTL inferior to independently trained
task specific, all-fc [46]	49	
cross stitch, all-fc [40]	53	
independent	67.7	
PCGrad (proposed)	71	

- But can we predict if two tasks are relevant?

# On Task Relevance

- [Taskonomy](#) by Stanford
- Measured as performance of transfer learning

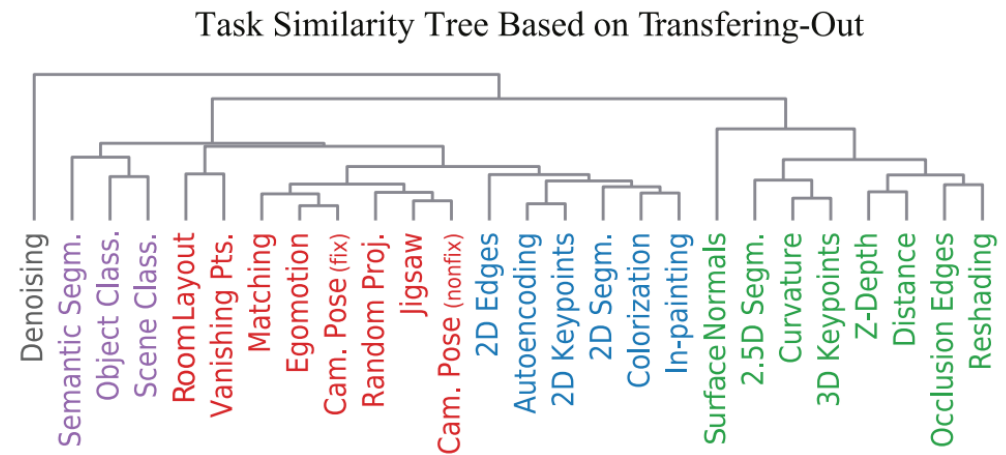
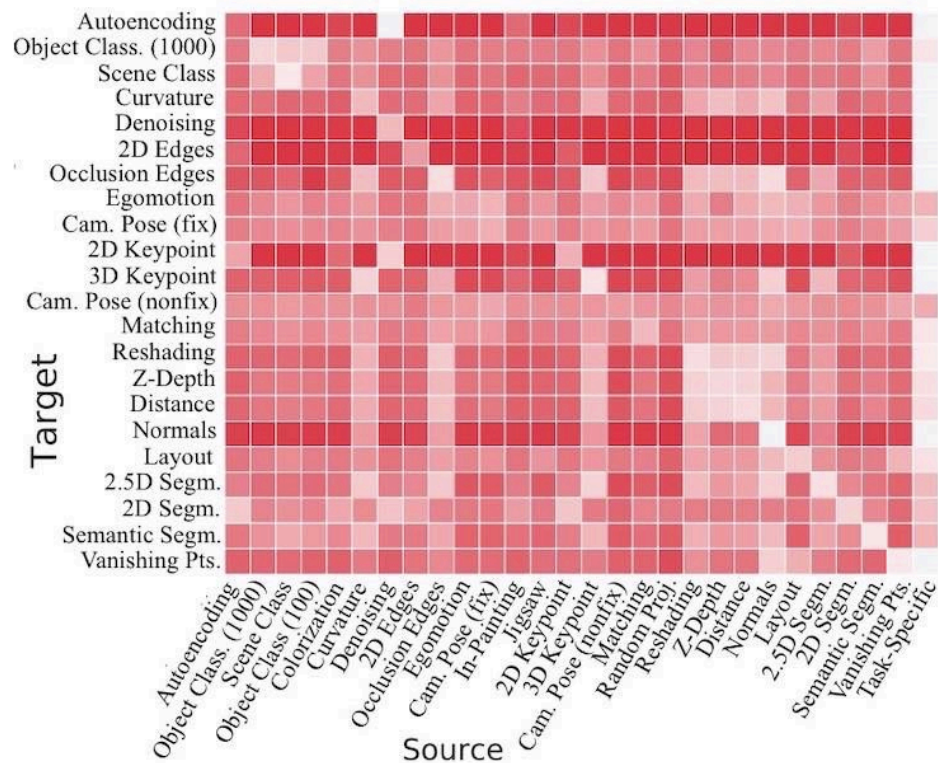
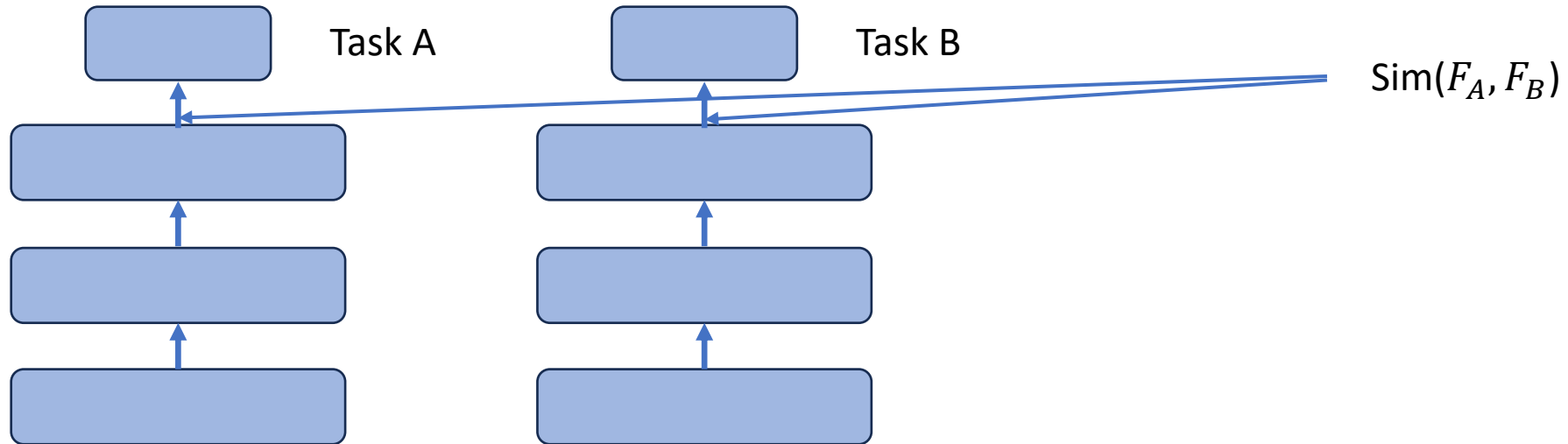


Figure 13: **Task Similarity Tree.** Agglomerative clustering of tasks based on their transferring-out patterns (i.e. using columns of normalized affinity matrix as task features). **3D**, **2D**, **low dimensional geometric**, and **semantic** tasks clustered together using a fully computational approach.



# Task Relevance: More Analytical way

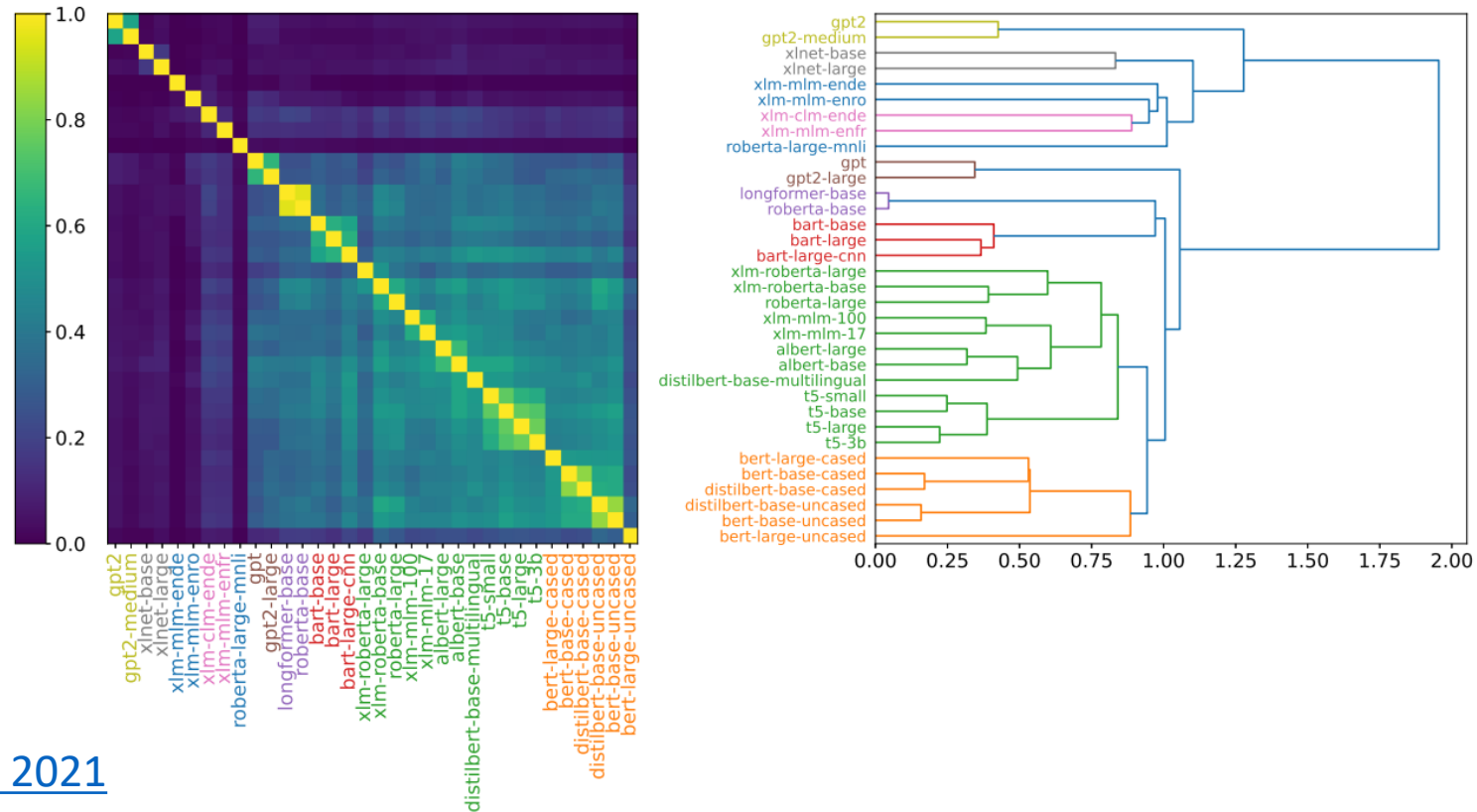
- Consider two tasks with same  $p(\mathbf{x})$ , but different  $p(\mathbf{y}|\mathbf{x})$ 's
- Assume we have trained a model for each of the two tasks



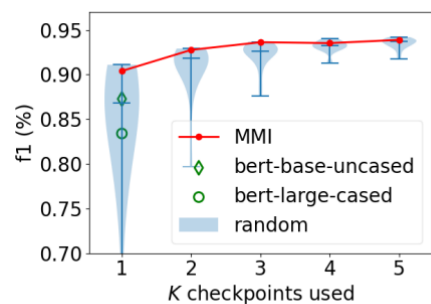
- Measure task relevance using features' similarity ([Huang et. al, 2021](#))

# Task Relevance: More Analytical way

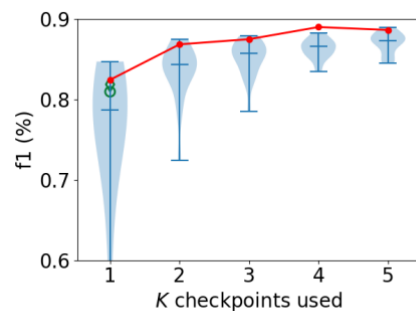
- On  $\text{Sim}(F_A, F_B)$ : **invariance** w.r.t linear transform (revisit in next lecture)



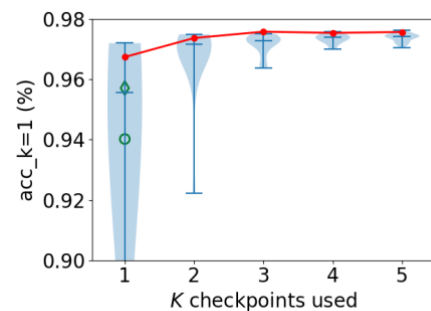
# Pick checkpoints for new tasks



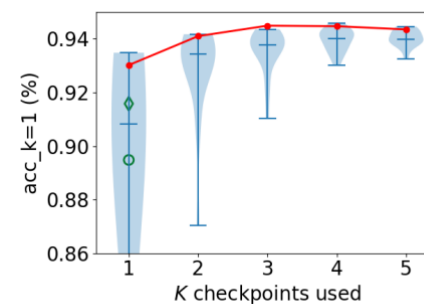
(a) Chunking



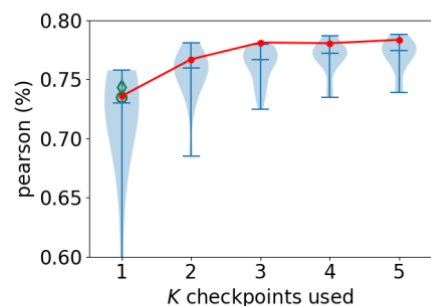
(b) NER



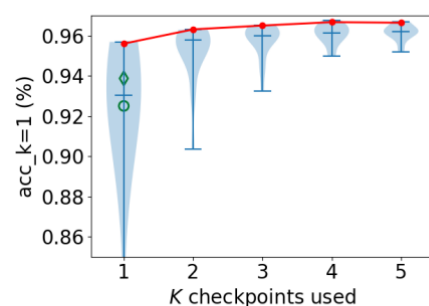
(c) POS Tagging



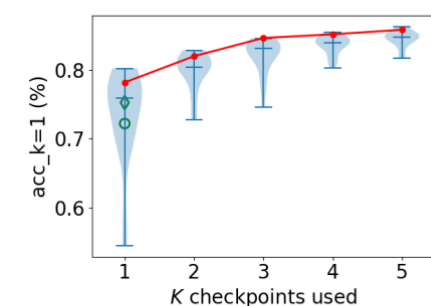
(d) Semantic tagging



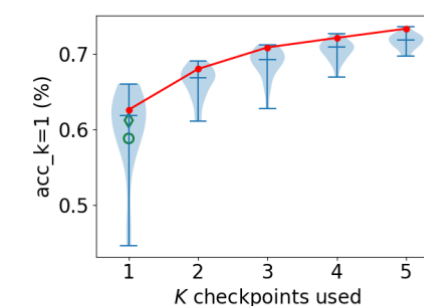
(e) Event factuality



(f) Tagging parent in phrase-structure tree



(g) Tagging grandparent in phrase-structure tree



(h) Tagging great grandparent in phrase-structure tree

Red: selected ckpt ; Blue: randomly picked ckpt

# Agenda for Today

- Multi-task Learning (MTL)
- Meta Learning
- Zero-shot Learning

# Motivating Meta Learning

- Sometimes, we may have to **learn a model from very few samples**
- i.e., few-shot learning
- e.g., 5-way, 1-shot classification

Given 1 example of 5 classes:

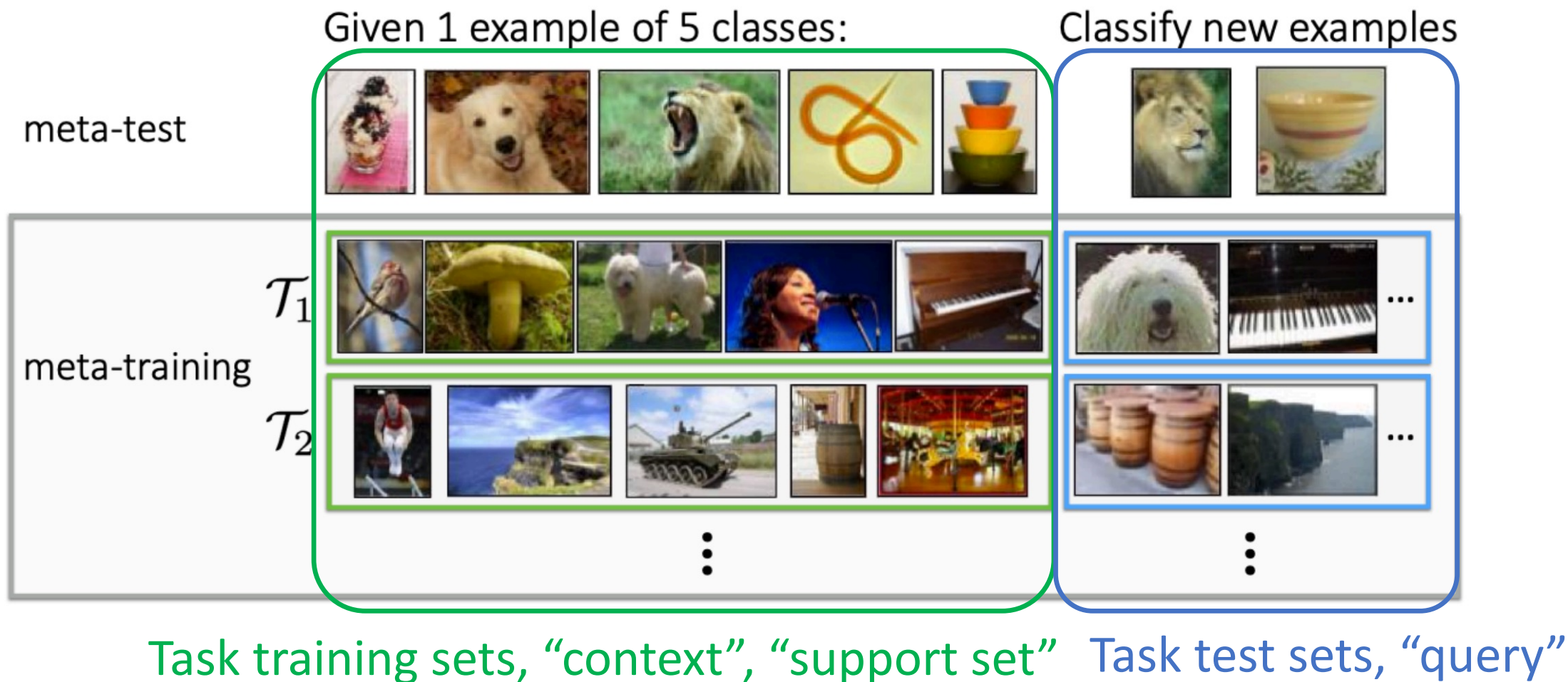


Classify new examples



- Seems very hard if we train a randomly initialized network!
- Can we start from a network that is good at few-shot learning?

# Motivating Meta Learning



# Formalize: Meta Learning

- Meta Training Set

- Tasks  $\mathcal{T}_1, \dots, \mathcal{T}_T$ , datasets  $\mathcal{D}_1, \dots, \mathcal{D}_T$ ;
- Each  $\mathcal{D}_i = \mathcal{D}_i^{tr} \cup \mathcal{D}_i^{te}$  (task training and test sets)

- Meta Test Set

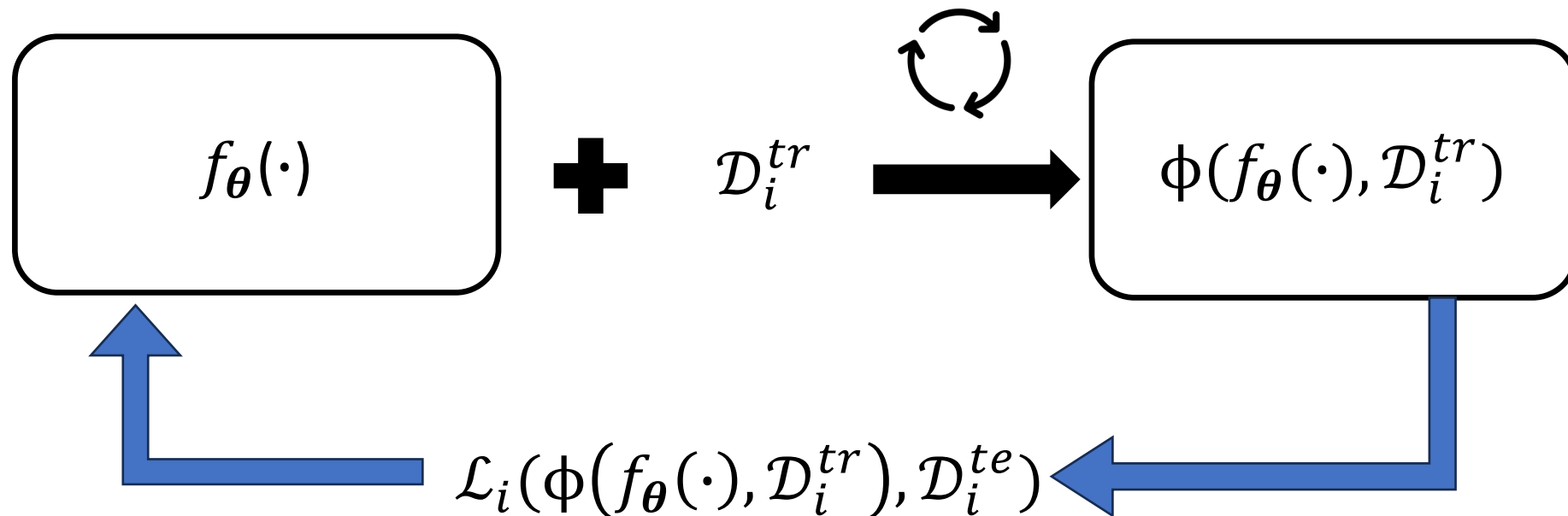
- New task  $\mathcal{T}_{T+1}$ , training samples  $\mathcal{D}_{T+1}^{tr}$ , test samples  $\mathcal{D}_{T+1}^{te}$

- Objective

Find a network  $f_\theta(\cdot)$ , so that if we few-shot train it on  $\mathcal{D}_i^{tr}$ , test result on  $\mathcal{D}_i^{te}$  is good

# Meta Learning: General Framework

1. Few shot Training
2. Get loss on task's test set
3. Back-prop loss to update  $\theta$





# Meta Learning: General Algorithm

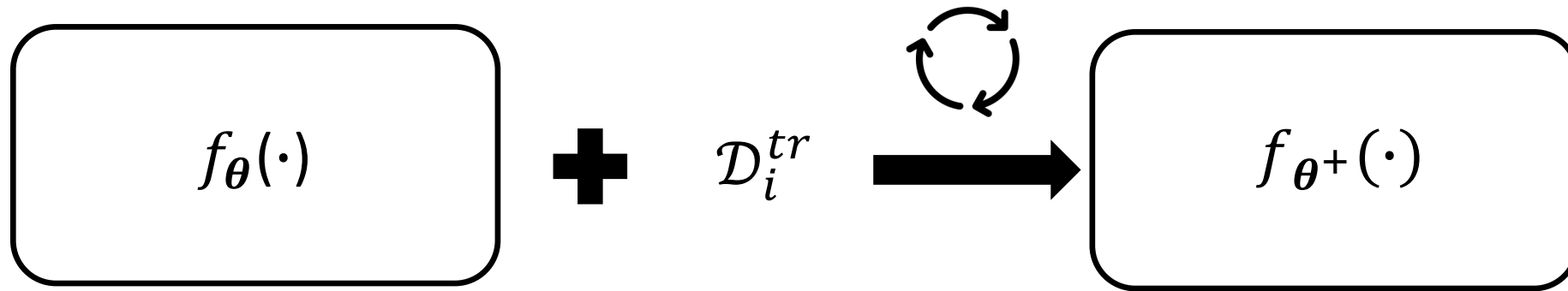
While not converged:

1. Sample task  $\mathcal{T}_i$
2. Starting from current network  $f_{\theta}(\cdot)$ , few-shot train on  $\mathcal{D}_i^{tr}$ ,  
Denote the task-specific model as  $\phi(f_{\theta}(\cdot), \mathcal{D}_i^{tr})$
3. Get test loss of  $\phi(f_{\theta}(\cdot), \mathcal{D}_i^{tr})$  on  $\mathcal{D}_i^{te}$ , denoted as  
$$\mathcal{L}_i(\phi(f_{\theta}(\cdot), \mathcal{D}_i^{tr}), \mathcal{D}_i^{te})$$
4. Update  $\theta$  via gradient descent

**Question: How is it different from transfer learning?**

# Model Agnostic Meta Learning (MAML)

- $\phi(f_{\theta}(\cdot), \mathcal{D}_i^{tr})$  is simply a gradient step on  $\theta$ :  
$$\theta^+ = \theta - \eta \nabla \mathcal{L}_i(\theta; \mathcal{D}_i^{tr})$$



- Evaluate test loss by

$$\mathcal{L}_i(\theta^+; \mathcal{D}_i^{te})$$

# Optimization for MAML Training Loss

- Overall training loss

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^m \{ \tilde{\mathcal{L}}_i(\boldsymbol{\theta}; \mathcal{D}_i^{te}) \triangleq \mathcal{L}_i(\boldsymbol{\theta}_i^+; \mathcal{D}_i^{te}) \}$$

with  $\boldsymbol{\theta}_i^+ = \boldsymbol{\theta} - \eta \nabla \mathcal{L}_i(\boldsymbol{\theta}; \mathcal{D}_i^{tr})$

- $\nabla \tilde{\mathcal{L}}_i(\boldsymbol{\theta}; \mathcal{D}_i^{te}) = [1 - \eta \nabla \otimes \nabla \mathcal{L}_i(\boldsymbol{\theta}; \mathcal{D}_i^{tr})] \nabla \mathcal{L}_i(\boldsymbol{\theta}_i^+; \mathcal{D}_i^{te})$

Hessian big, let's ignore it

# Results on mini-ImageNet

MiniImagenet (Ravi & Larochelle, 2017)	5-way Accuracy	
	1-shot	5-shot
fine-tuning baseline	$28.86 \pm 0.54\%$	$49.79 \pm 0.79\%$
nearest neighbor baseline	$41.08 \pm 0.70\%$	$51.04 \pm 0.65\%$
matching nets (Vinyals et al., 2016)	$43.56 \pm 0.84\%$	$55.31 \pm 0.73\%$
meta-learner LSTM (Ravi & Larochelle, 2017)	$43.44 \pm 0.77\%$	$60.60 \pm 0.71\%$
<b>MAML, first order approx.</b>	<b><math>48.07 \pm 1.75\%</math></b>	<b><math>63.15 \pm 0.91\%</math></b>

# Agenda for Today

- Multi-task Learning (MTL)
- Meta Learning
- Zero-shot Learning

# Setup: Zero-Shot Learning (ZSL)

- Training: input  $x_i$ , label  $y_i \in \mathcal{V}$
- Test: input  $x$ , predict label  $y \notin \mathcal{V}$
- Impossible if the labels are just categorical
- What if labels have semantics?

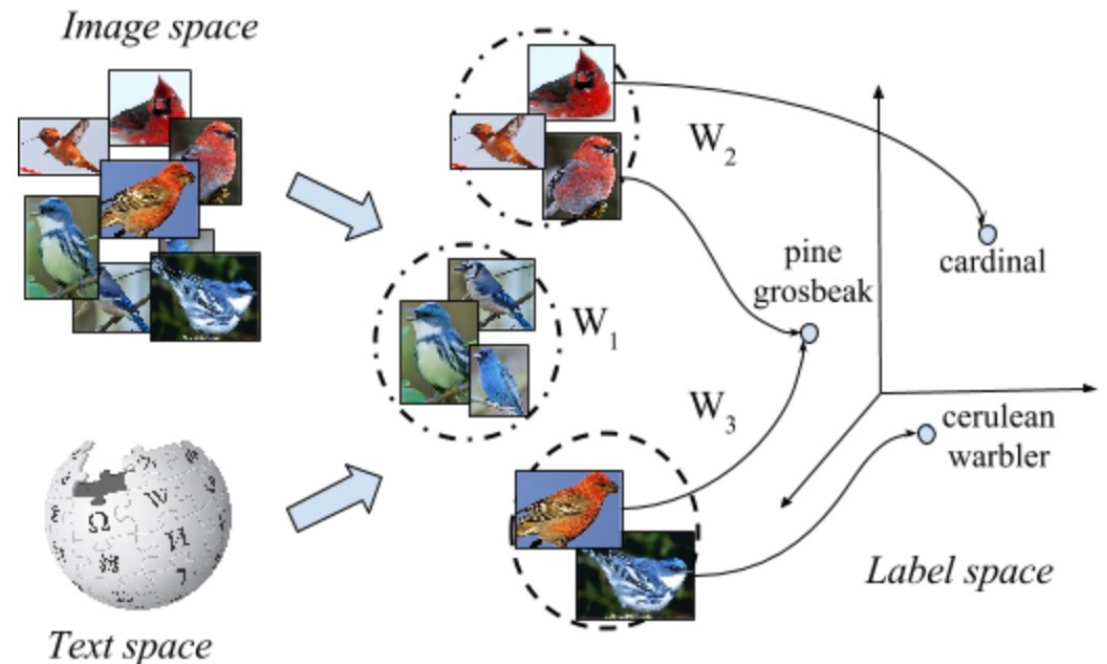
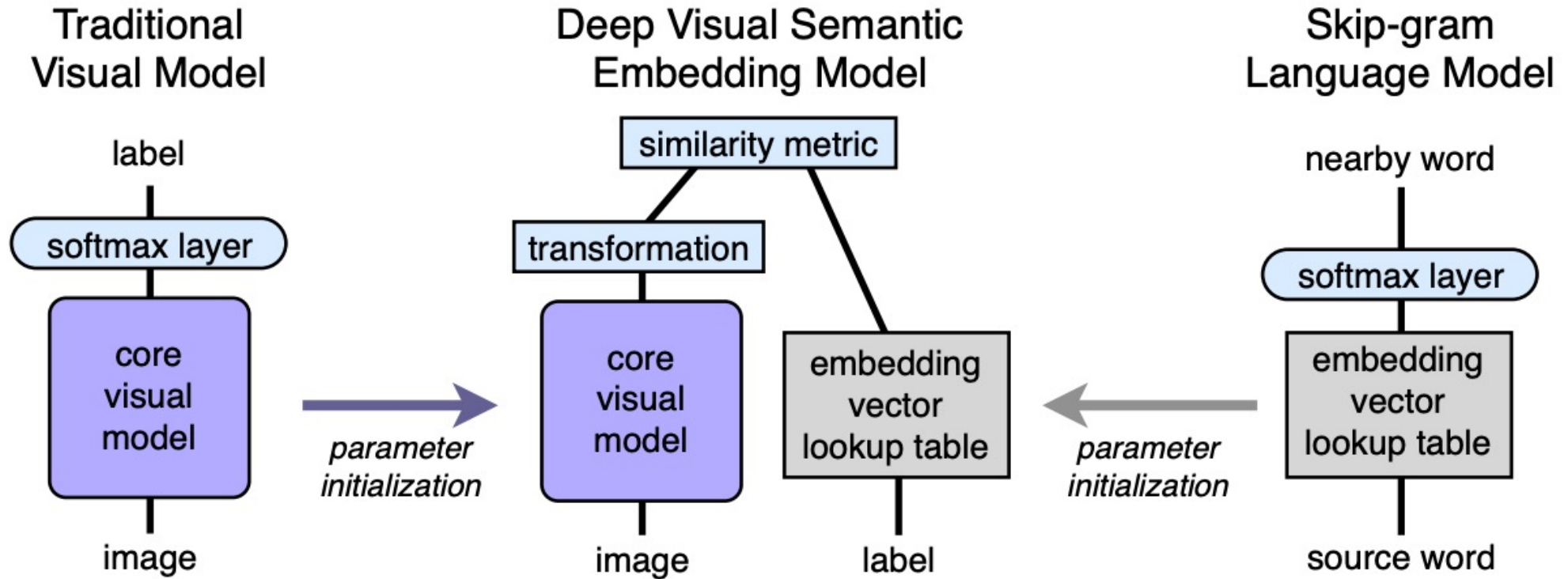


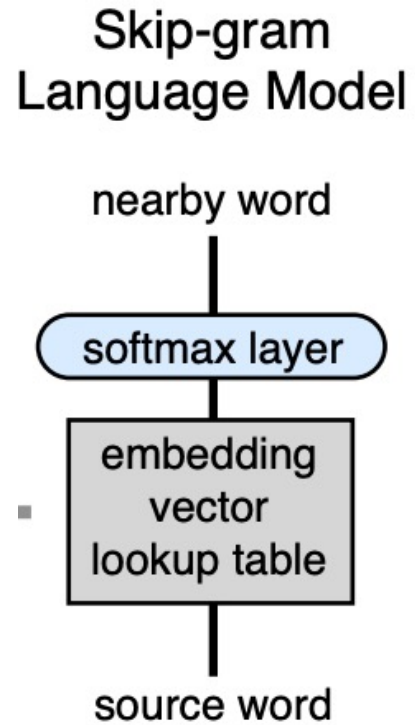
Illustration from [here](#)

# ZSL Image Classification



# Recap

- How does skip gram work?





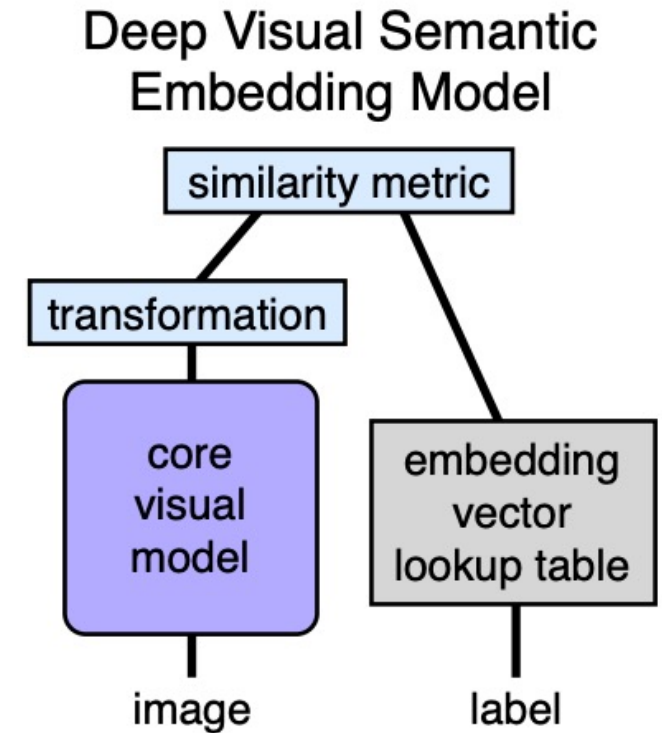
# Metric Learning Objective

- The usual way to measure two vector's similarity

$$x^T y = \sum_i x_i y_i$$

- More generally, we may want to
  - weigh the dimensions
  - Consider cross dimensions

- That's  $\sum_{i,j} m_{i,j} x_i y_j = x^T M y$

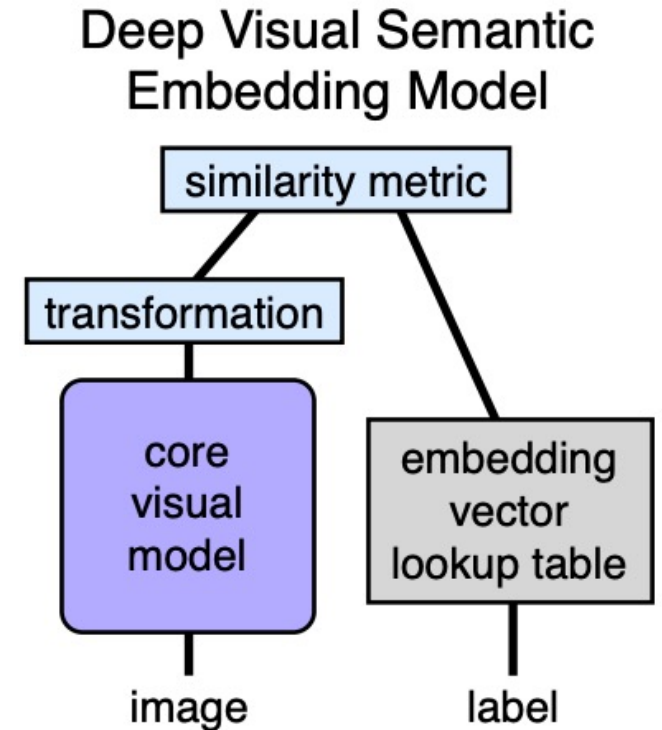
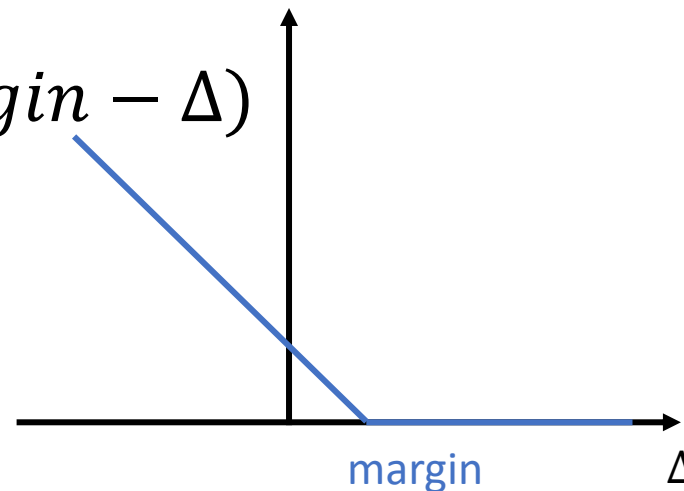


# Metric Learning Objective

- Image vector should be close to its text label
- But far away from a wrong text label
- Require the distance differ by some “margin”

$$\text{loss}(\text{image}, \text{label}) = \sum_{j \neq \text{label}} \max[0, \text{margin} - \vec{t}_{\text{label}} M \vec{v}(\text{image}) + \vec{t}_j M \vec{v}(\text{image})]$$







- Hinge loss =  $\max(0, \text{margin} - \Delta)$



# Testing Phase

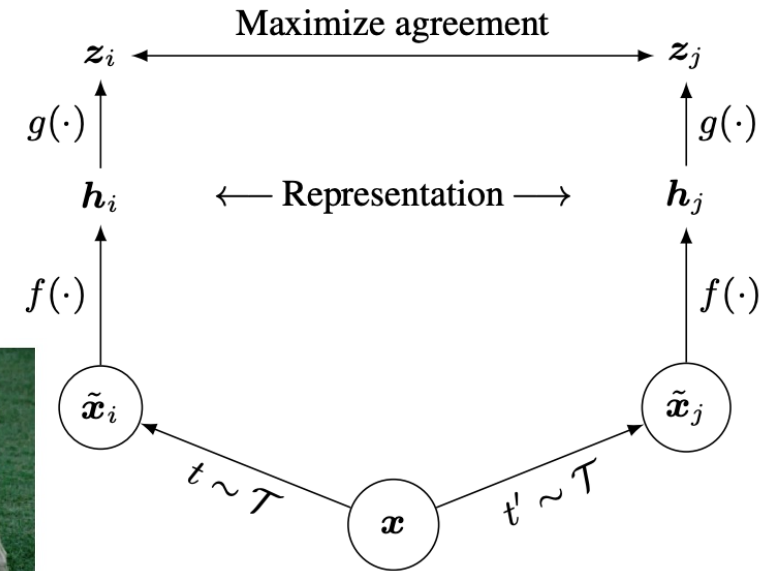
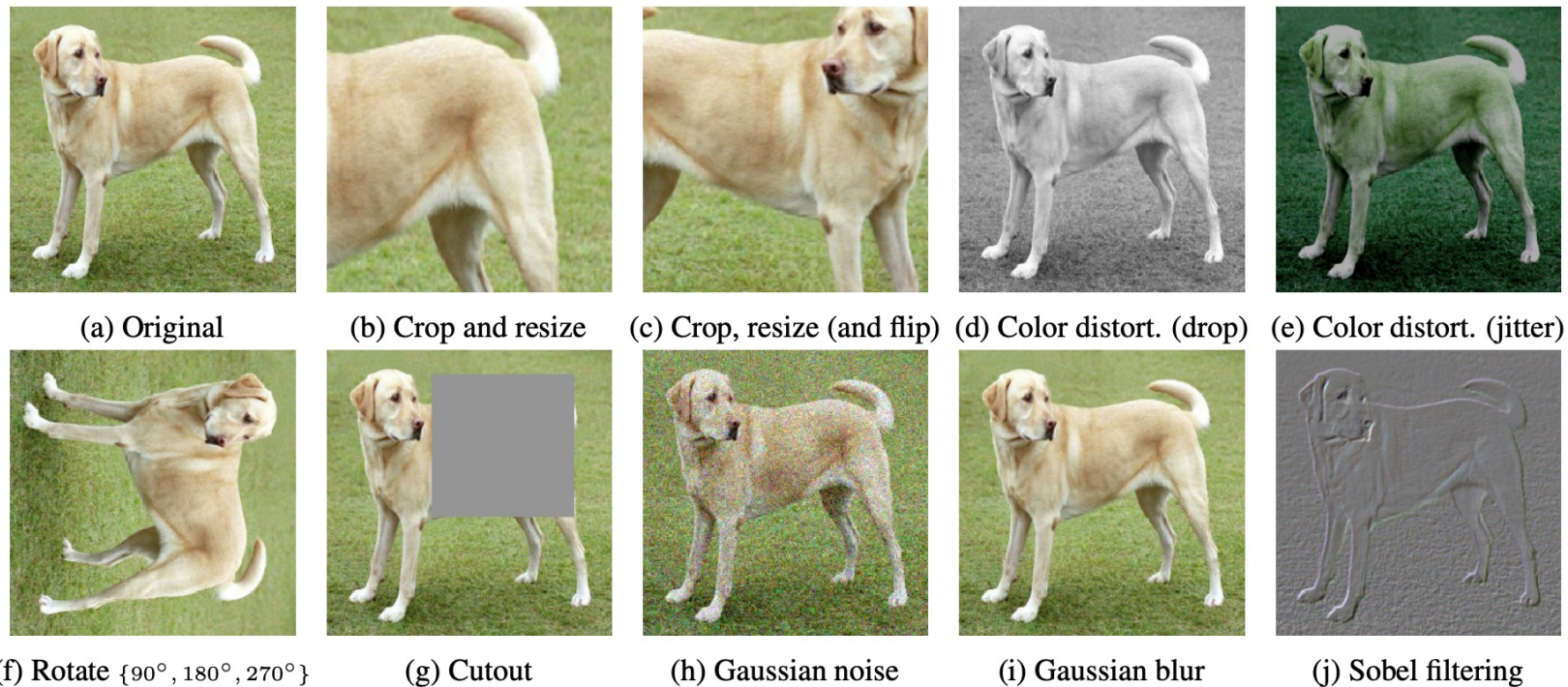
- Classify a new image by nearest neighbor search
- But with distance metric  $M$

$$\min_{text} \vec{v}(text) M \vec{v}(image)$$

	ZSL	Softmax over ImageNet 1K		ZSL	Softmax over ImageNet 1K
	<p><b>A</b></p> eyepiece, ocular Polaroid compound lens <b>telephoto lens, zoom lens</b> rangefinder, range finder	typewriter keyboard tape player reflex camera CD player space bar		<p><b>D</b></p> fruit pineapple <b>pineapple plant, Ananas</b> sweet orange sweet orange tree, ...	pineapple, ananas coral fungus .. artichoke, globe artichoke sea anemone, anemone cardoon
	<p><b>B</b></p> oboe, hautboy, hautbois bassoon <b>English horn, cor anglais</b> hook and eye hand	reel punching bag, punch bag, ... whistle bassoon letter opener, paper knife, ...		<p><b>E</b></p> comestible, edible, ... dressing, salad dressing Sicilian pizza vegetable, veggie, veg fruit	pot, flowerpot cauliflower guacamole cucumber, cuke broccoli
	<p><b>C</b></p> barbet patas, hussar monkey, ... <b>babblers, cackler</b> titmouse, tit bowerbird, catbird	patas, hussar monkey, ... proboscis monkey, Nasalis ... macaque titi, titi monkey guenon, guenon monkey		<p><b>F</b></p> dune buggy, beach buggy searcher beetle, ... seeker, searcher, quester Tragelaphus eurycerus, ... bongo, bongo drum	warplane, military plane missile projectile, missile sports car, sport car submarine, pigboat, sub, ...

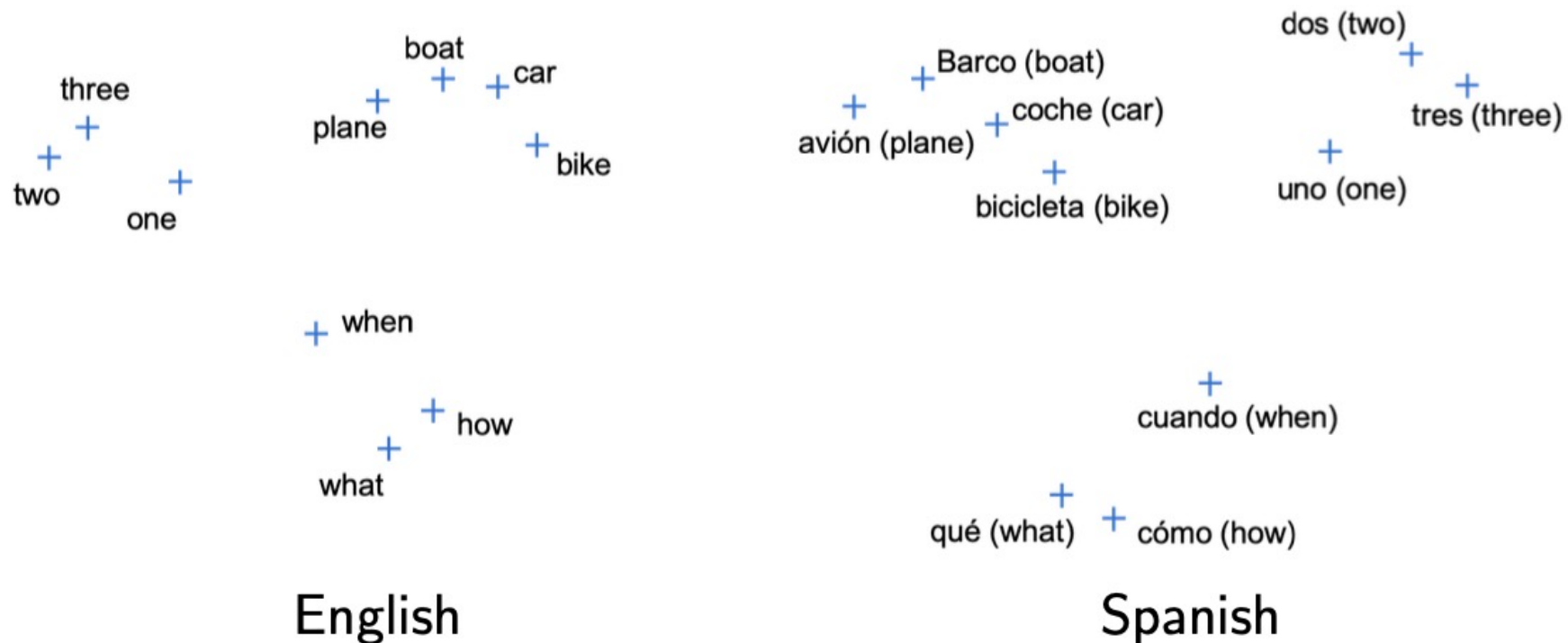
# More on Learning by Contrast

- Self-supervised pretraining of vision models
- $\mathcal{T}$  : set of augmentations



# Another ZSL: Bilingual Lexicon Induction (BLI)

- Generate word-to-word translation from very few “seeding” pairs
- Again take advantage of word embeddings



# BLI

- Source embedding  $X \in \mathbb{R}^{n \times d}$ , target embedding  $Y \in \mathbb{R}^{n \times d}$
- Learn a rotation matrix  $R \in \mathbb{R}^{d \times d}$ ,  $RR^T = I$
- Procrustes problem

$$\min_{R:RR^T=I} \|XR - Y\|^2$$



Procrustes



Theseus

# Solving the Procrustes Problem

$$\min_{R:RR^T=I} \|XR - Y\|^2$$

- We can show it's equivalent to solving

$$\max_{R:RR^T=I} \langle R, Y^T X \rangle$$

- Let the SVD of  $Y^T X = U\Lambda V^T$ , then optimum  $R^* = UV^T$
- Why impose  $RR^T = I$ ?
  - Prior knowledge: languages should share some fundamentals
  - A regularization