

Robust Learning under Distribution Shifts

Machine Learning Summer School, Okinawa

Mar. 8th, 2024

Han Zhao

hanzhao@illinois.edu

Assistant Professor

Department of Computer Science

University of Illinois Urbana-Champaign



Robustness

Domain-Invariant Representations

Fundamental limit in domain adaptation

Invariant Risk Minimization

An efficient algorithm for IRM via post-processing

Gradual Domain Adaptation

Learning under continuous distribution shifts

Robust Multitask Learning

Understanding multi-objective optimization

Invariant Representation Learning

Information-theoretic analysis between invariance & accuracy, with applications in distributional robustness and fairness

Tradeoff between robustness and fairness

Understanding the impact of adversarial robustness on fairness

Learning Fair Representations

Tradeoff between fairness & accuracy

Fairness

Fair and Optimal Classification (I)

An optimal post-processing algorithm for demographic parity

Fair and Optimal Classification (II)

An optimal post-processing algorithm for equalized odds

Trustworthy Machine Learning



Differentially-Private and Fair Regression

Design fair & private regression algorithm

Privacy-Preserving Learning

Learning under attribute-inference attack

Machine Unlearning

Removing a subset of specified data from the learned model

Efficiency & Accuracy

Maximum Influence Subset

Understand and select out the subset of data with maximum influence to the learned model

Structured Representations

Learning representations with class hierarchical information

Privacy-Preserving Learning for Graph Neural Networks

Provide defenses strategies under attribute-inference attacks over graphs

Interpretability

Privacy

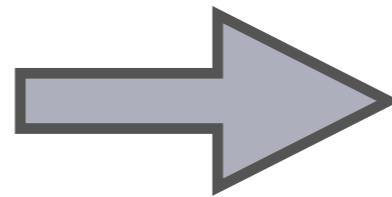
Robustness:

- Domain adaptation / generalization / Out-of-distribution generalization / transfer learning
- Invariant representation learning / invariant causal predictors

Key Factors underlying the Success

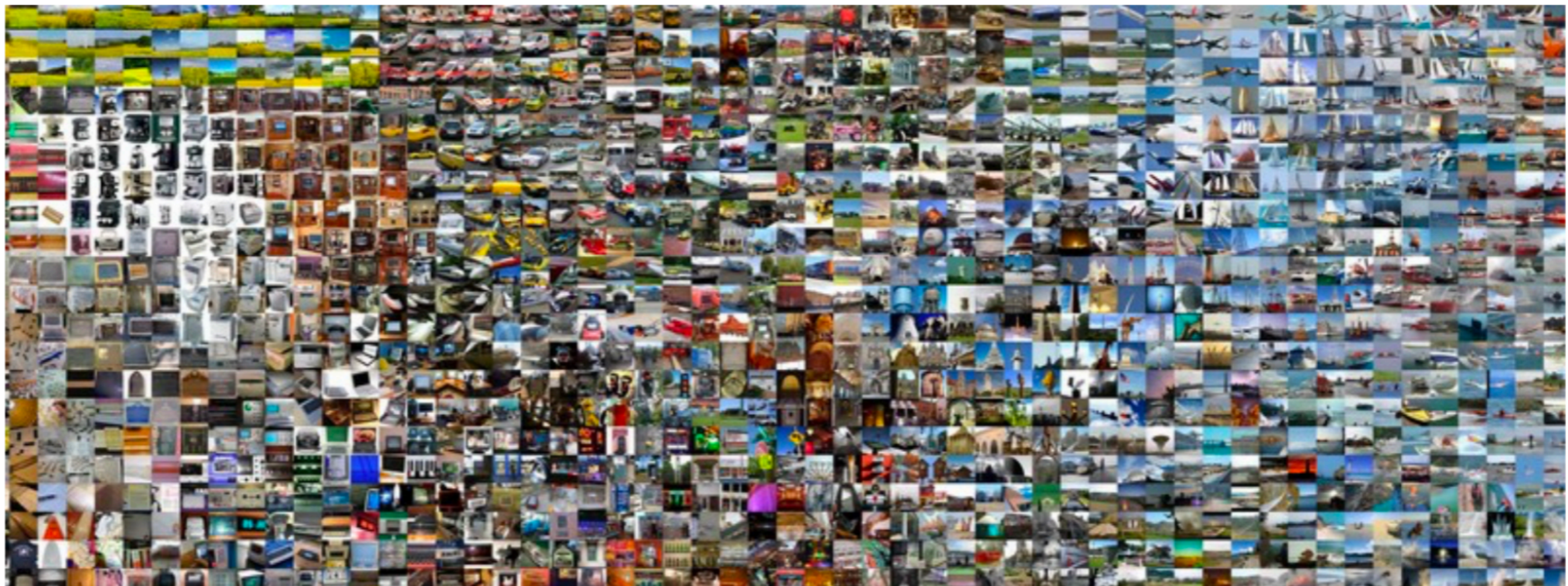
- Stationary distributions
- Large-scale labeled datasets & Powerful computation

Source (Train) = Target (Test)



Source

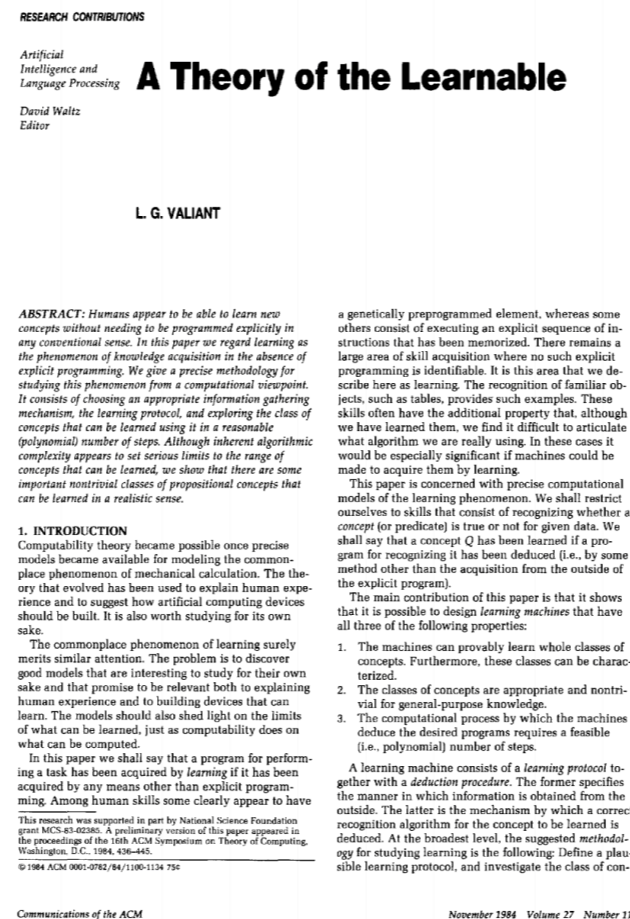
Target



Probably Approximately Correct (PAC)

Probably Approximately Correct (PAC, Valiant, CACM 1984)

- (Informal) A framework to quantify the meaning of learning a concept from samples
- With high probability (P), the learned predictor will have low generalization error (AC)
- No distributional assumption



Probably Approximately Correct (PAC)

With VC dim as the complexity measure, we have the following **uniform** generalization bound:

Given $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \sim \mu$ be a dataset of iid samples. Let \mathcal{F} be a hypothesis class of finite VC-dim, i.e., $\text{VCdim}(\mathcal{F}) < \infty$, then for $0 < \delta < 1$, with probability at least $1 - \delta$, for **all** $f \in \mathcal{F}$:

$$\varepsilon_{\mu}(f) \leq \hat{\varepsilon}_{\mathcal{D}}(f) + O\left(\sqrt{\frac{\text{VCdim}(\mathcal{F}) + \log(1/\delta)}{n}}\right)$$

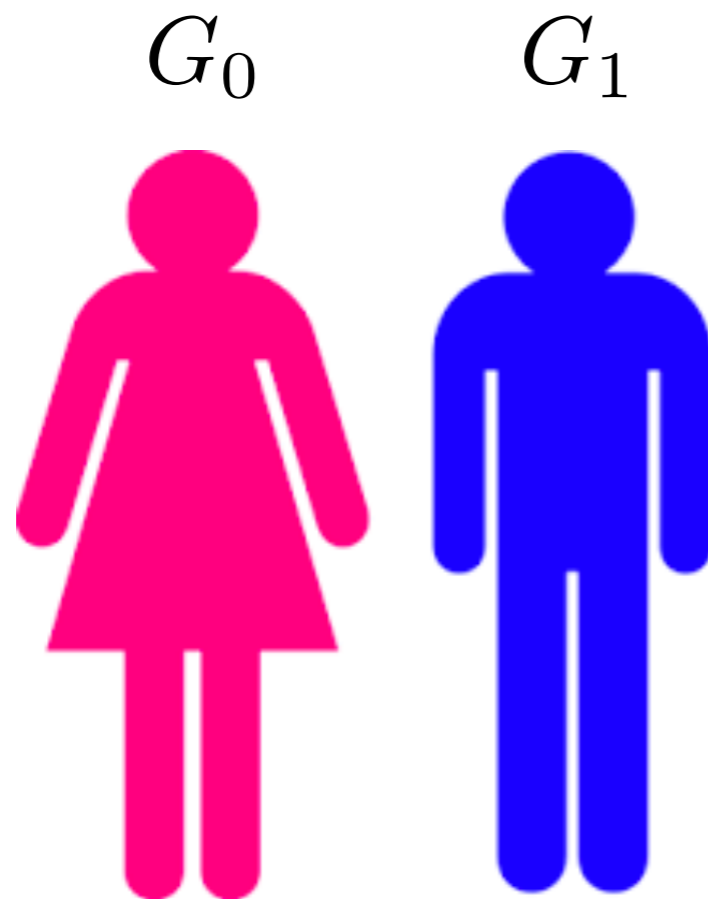
Note:

- The bound above gives the generalization error, and we can use the generalization error bound to provide an upper bound on the excess risk as well, i.e., $\varepsilon_{\mu}(f) - \inf_{f' \in \mathcal{F}} \varepsilon_{\mu}(f')$
- There are other forms of complexity measures to characterize the expressiveness/richness/powerfulness of a given hypothesis class, e.g., Rademacher complexity, covering number, algorithmic stability, etc, but the high-level form is the same
- The bound above could be loose, i.e., the generalization error could be larger than 1 for classification problems
- Other kinds of generalization bounds exist as well, i.e., instead of providing high-probability bounds, providing guarantees on expected generalization error (PAC-Bayes, mutual information, etc)

Distribution Shift in Practice

Commercial face recognition systems

Some of the key findings from the paper:



- All classifiers perform better on male faces than female faces (8.1% – 20.6% difference in error rate)
- All classifiers perform better on lighter faces than darker faces (11.8% – 19.2% difference in error rate)
- All classifiers perform worst on darker female faces (20.8% – 34.7% error rate)
- Microsoft and IBM classifiers perform best on lighter male faces (error rates of 0.0% and 0.3% respectively)
- Face++ classifiers perform best on darker male faces (0.7% error rate)
- The maximum difference in error rate between the best and worst classified groups is 34.4%

Accuracy disparity:

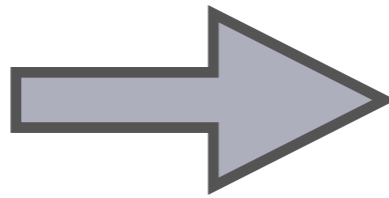
$$\Delta_{\text{Err}}(h) := |\Pr(h(X) \neq Y \mid X \sim G_0) - \Pr(h(X) \neq Y \mid X \sim G_1)|$$

Domain Adaptation / Generalization

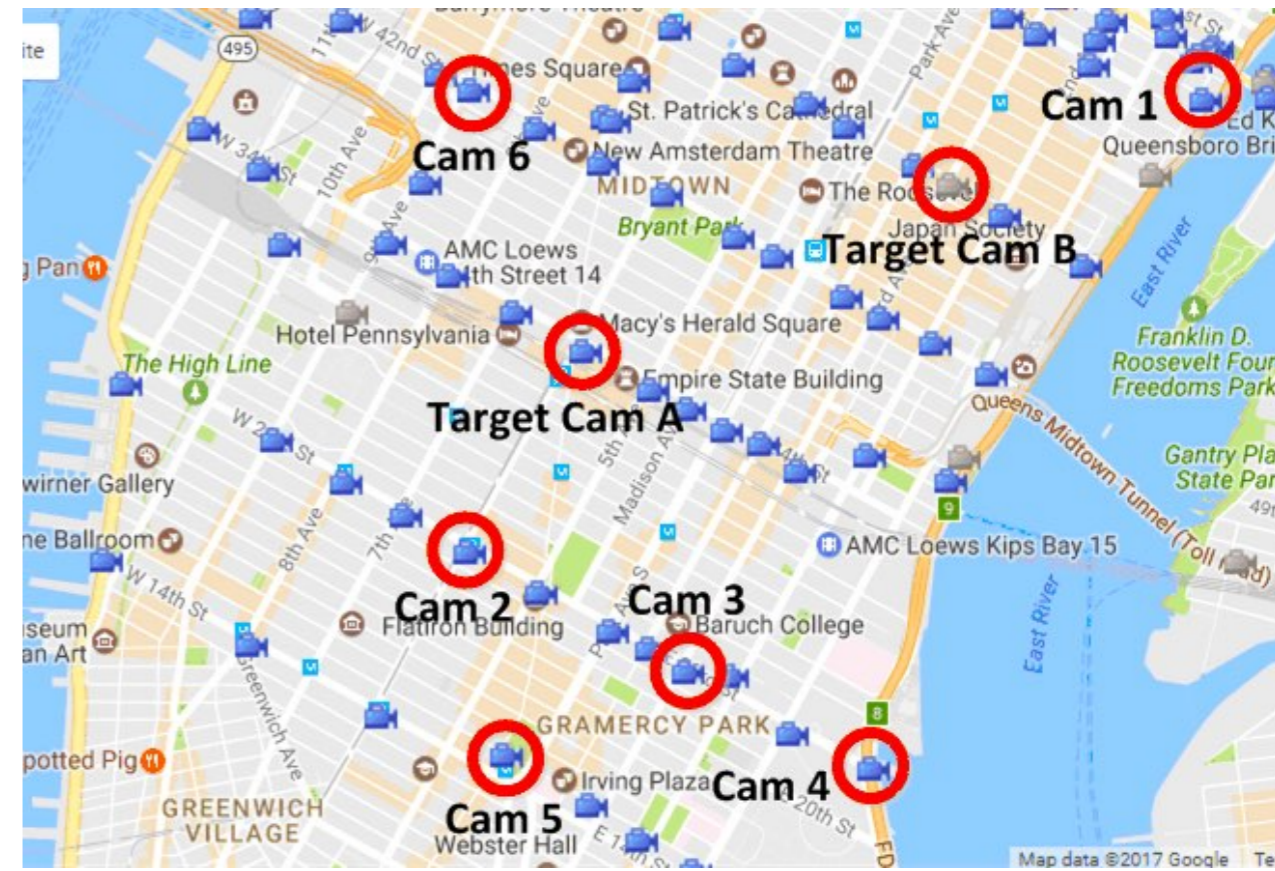
Machine learning models could be brittle



Source (with Labels)



Target (No Labels)



Domain Adaptation / Generalization

Domain adaptation: given unlabeled data from the target domain + labeled data from the source domain, can we do better?

Note: closely related to the setting of semi-supervised learning, but with a key difference:

Semi-supervised learning:

Training distribution \equiv Test distribution

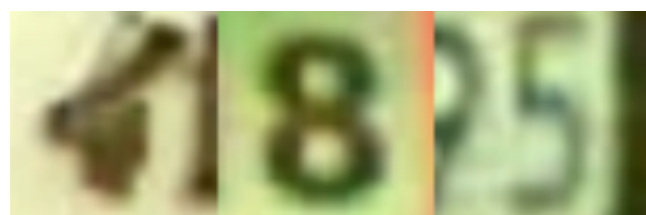
Domain adaptation:

Training distribution \neq Test distribution

Domain Adaptation / Generalization

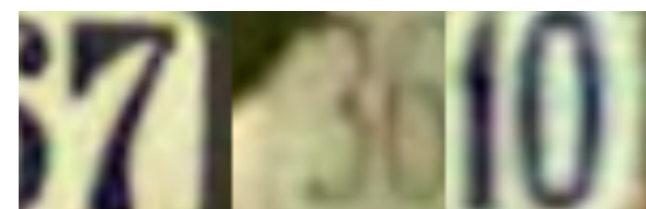
Domain adaptation: Training phase

Source domain:



(4, 8, 5)

• • •



(7, 3, 0)

+

Target domain:



• • •



Domain Adaptation / Generalization

Domain adaptation: Training phase

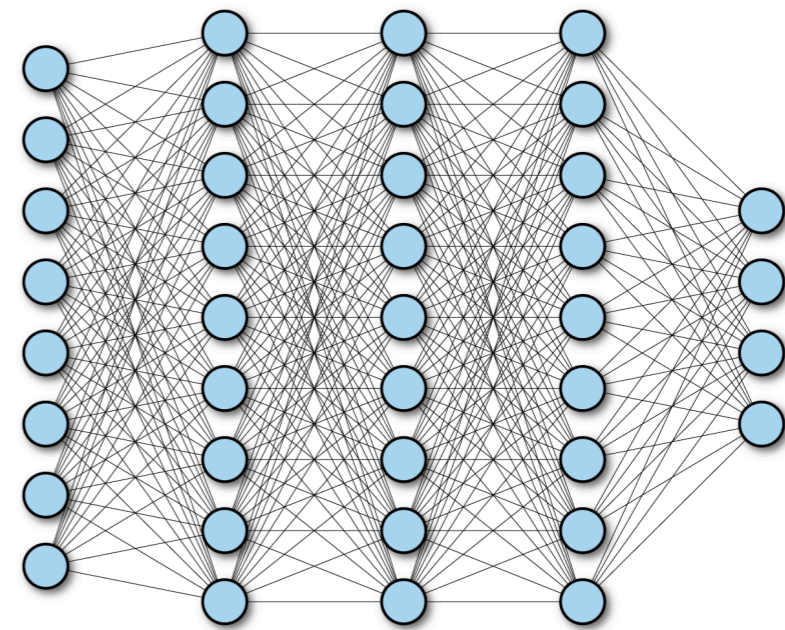
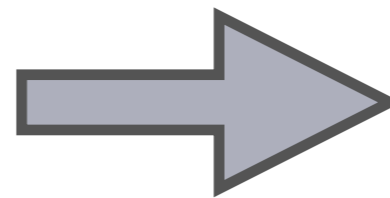
Source domain:



(4, 8, 5) (7, 3, 0)

+

Target domain:

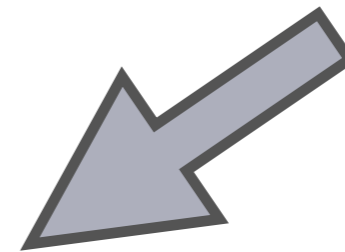
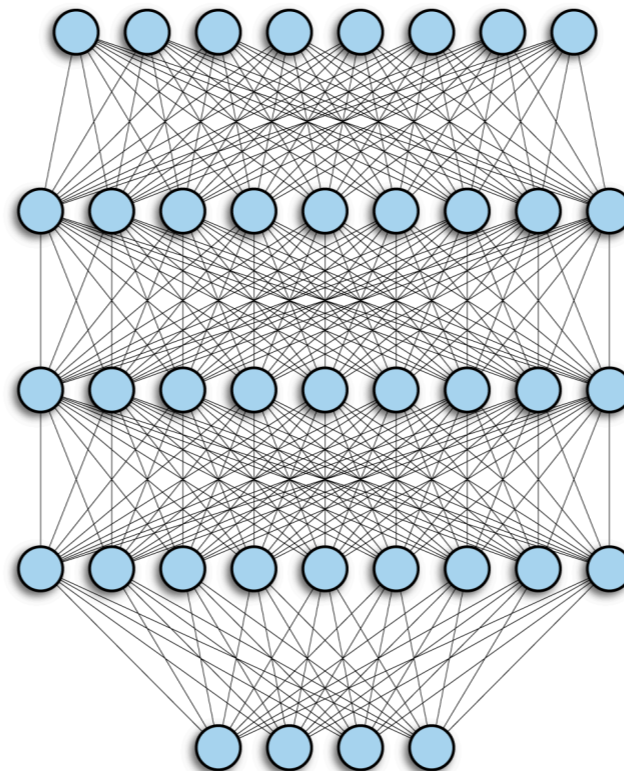
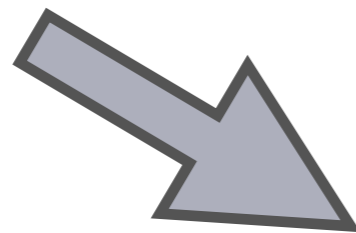


Classifier

Domain Adaptation / Generalization

Domain adaptation: Testing phase

Target domain:



(4, 0, 1)
✓ ✓ ✓

(0, 4, 2)
✗ ✓ ✓

Distribution Shift in Practice

Mismatched distributions: **Source** \neq **Target**



Is minimizing the source error sufficient for generalization on the target domain?

What we want: Generalization on the **Target** domain

What we have: Labeled data from **Source** domain + unlabeled data from **Target** domain

Distribution Shift in Practice

Mismatched distributions: **Source** \neq **Target**

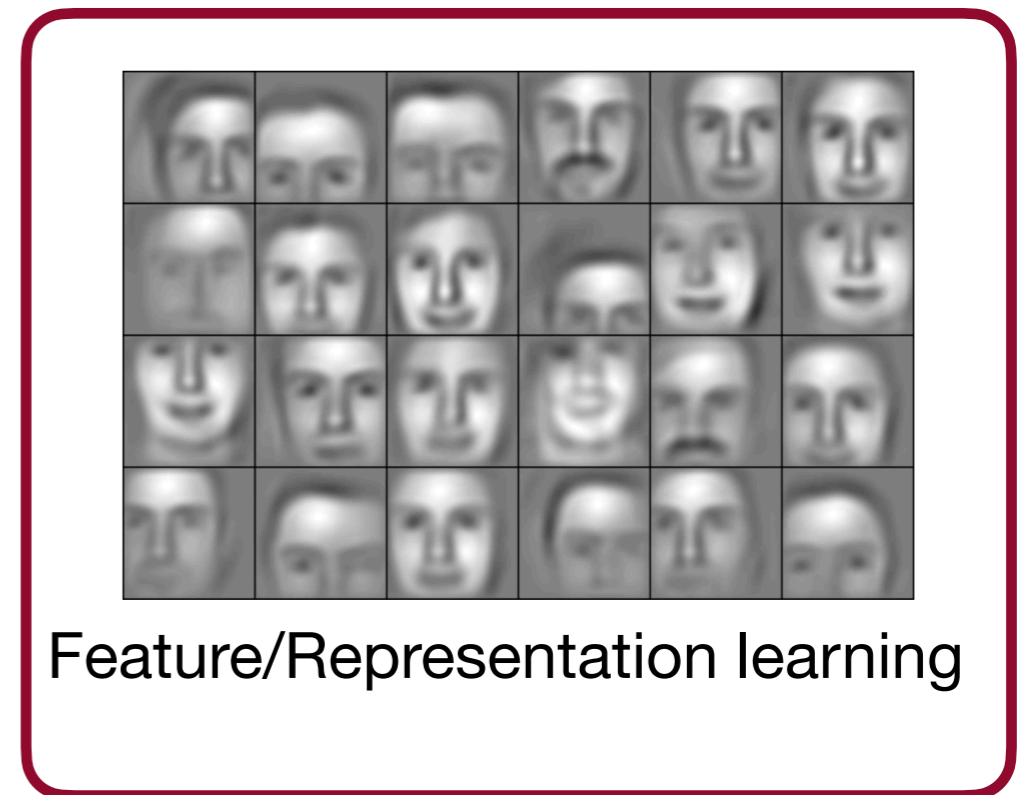
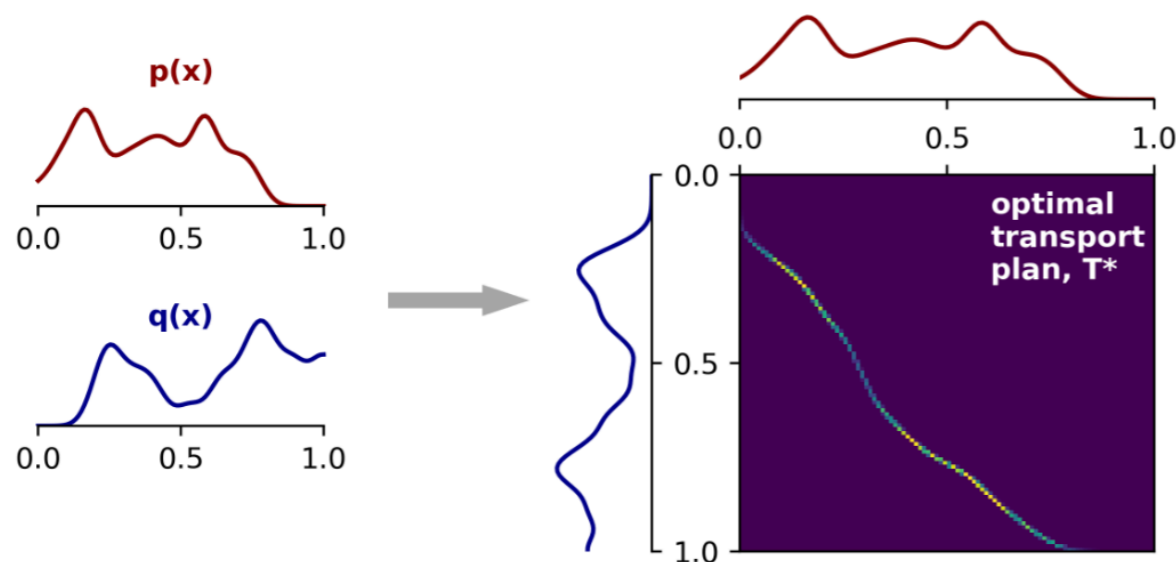
Is minimizing the source error sufficient for generalization on the target domain?

Importance Sampling

$$\mathbb{E}_q[h'(X)] = \int_{\mathbb{R}} h(x) \frac{p(x)}{q(x)} q(x) dx$$

Importance distribution q , Importance weight $\frac{p(x)}{q(x)}$

Importance sampling



Optimal transport

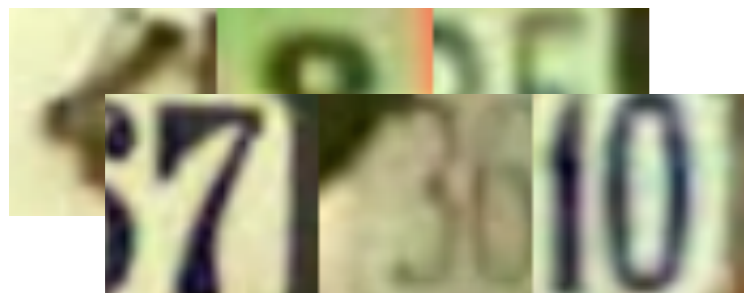
Domain-Invariant Representation Learning

How to bridge the gap?

What we want: Generalization on the **Target** domain

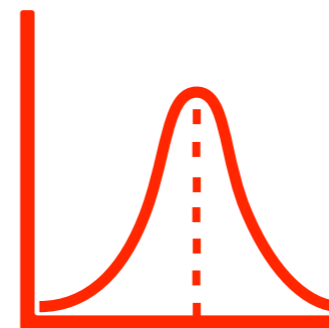
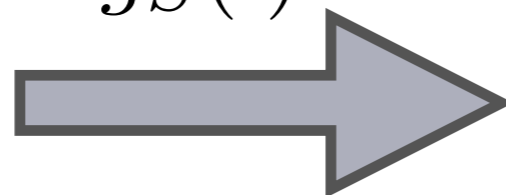
What we have: Labeled data from **Source** domain + unlabeled data from **Target** domain

Make the two domains closer

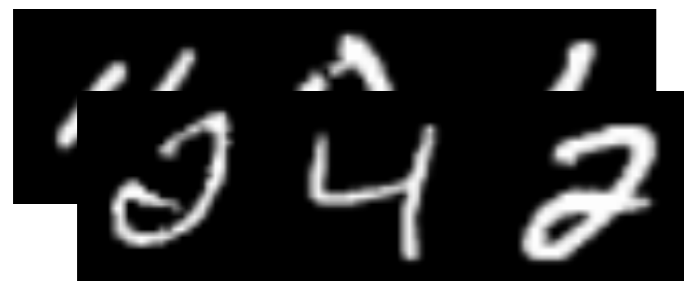


Source

$g_S(\cdot)$

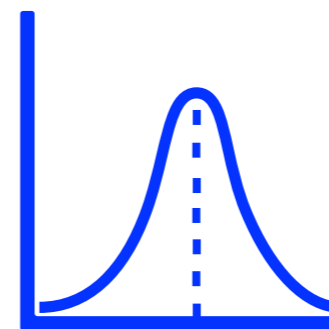
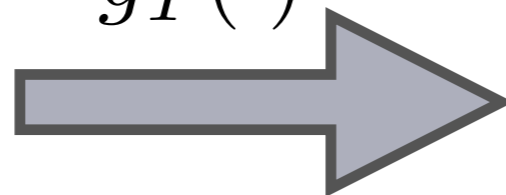


7	3	4	6	1	8	1	0	9	8
0	3	1	2	7	0	2	9	6	0
1	6	7	1	9	7	6	5	5	8
8	3	4	4	8	7	3	6	4	6
6	3	8	8	9	9	4	4	0	7
8	1	0	0	1	8	5	7	1	7
5	5	9	9	4	2	5	3	7	4
6	6	0	1	0	1	2	4	8	5
3	5	0	0	6	4	3	8	3	7
1	4	3	9	2	2	0	3	6	6



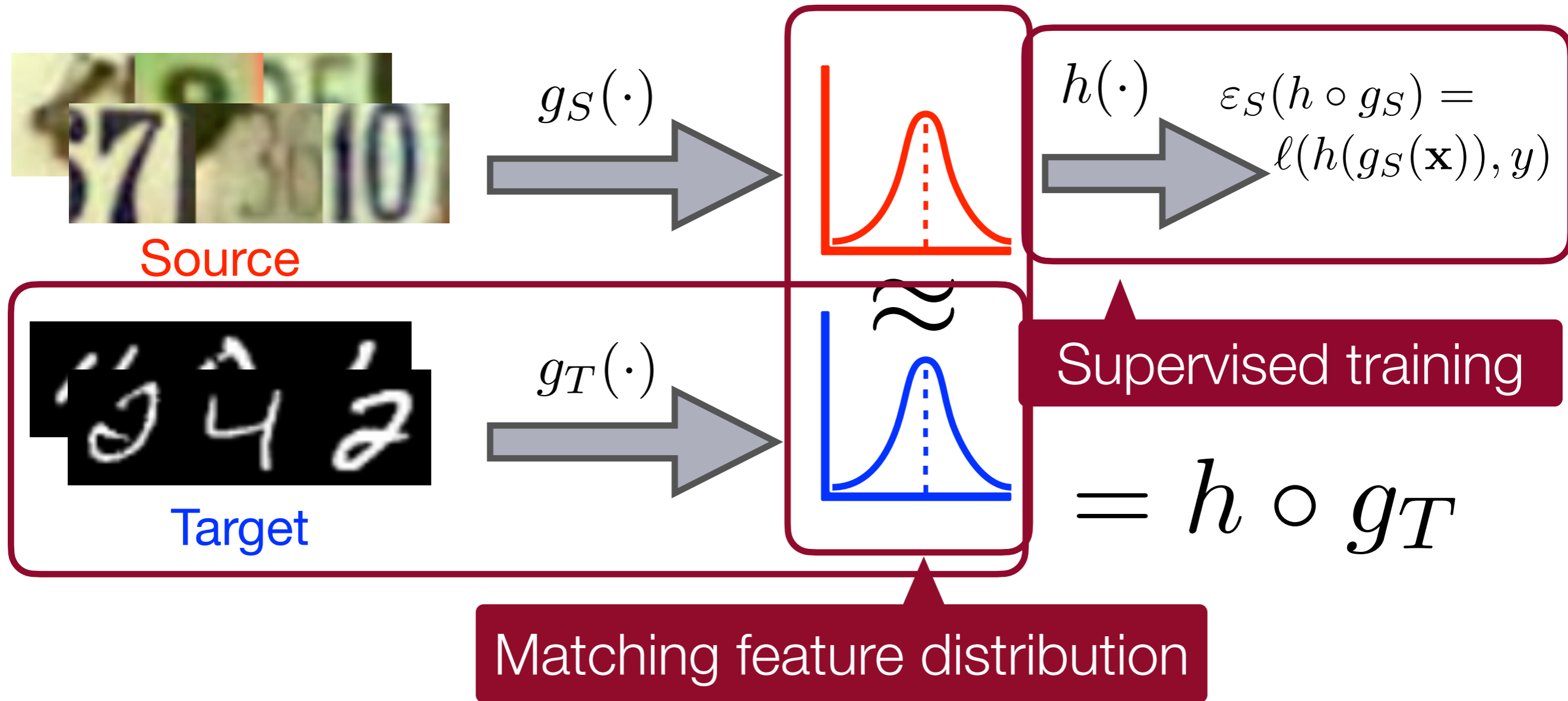
Target

$g_T(\cdot)$



7	3	4	6	1	8	1	0	9	8
0	3	1	2	7	0	2	9	6	0
1	6	7	1	9	7	6	5	5	8
8	3	4	4	8	7	3	6	4	6
6	3	8	8	9	9	4	4	0	7
8	1	0	0	1	8	5	7	1	7
5	5	9	9	4	2	5	3	7	4
6	6	0	1	0	1	2	4	8	5
3	5	0	0	6	4	3	8	3	7
1	4	3	9	2	2	0	3	6	6

Domain-Invariant Representation Learning

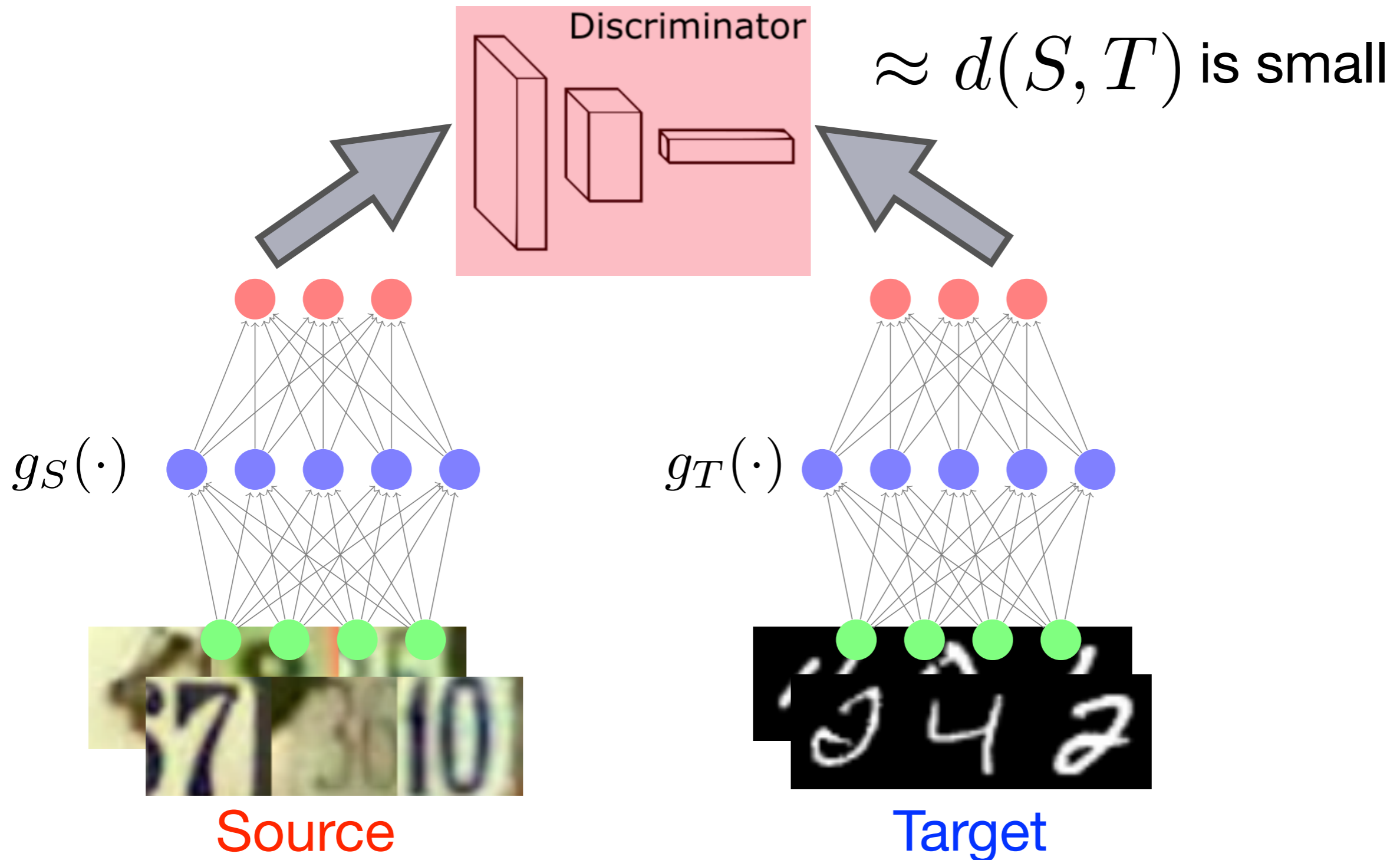


- Adversarial discriminator [Ganin et al. ICML15]
- Maximum-mean discrepancy [Long et al. ICML 15]
- Transportation distance [Shen et al. AAAI 18]
- Multiple-domain extensions [Zhao et al. NeurIPS 18]

Goal: generalization on target domain

Domain-Invariant Representation Learning

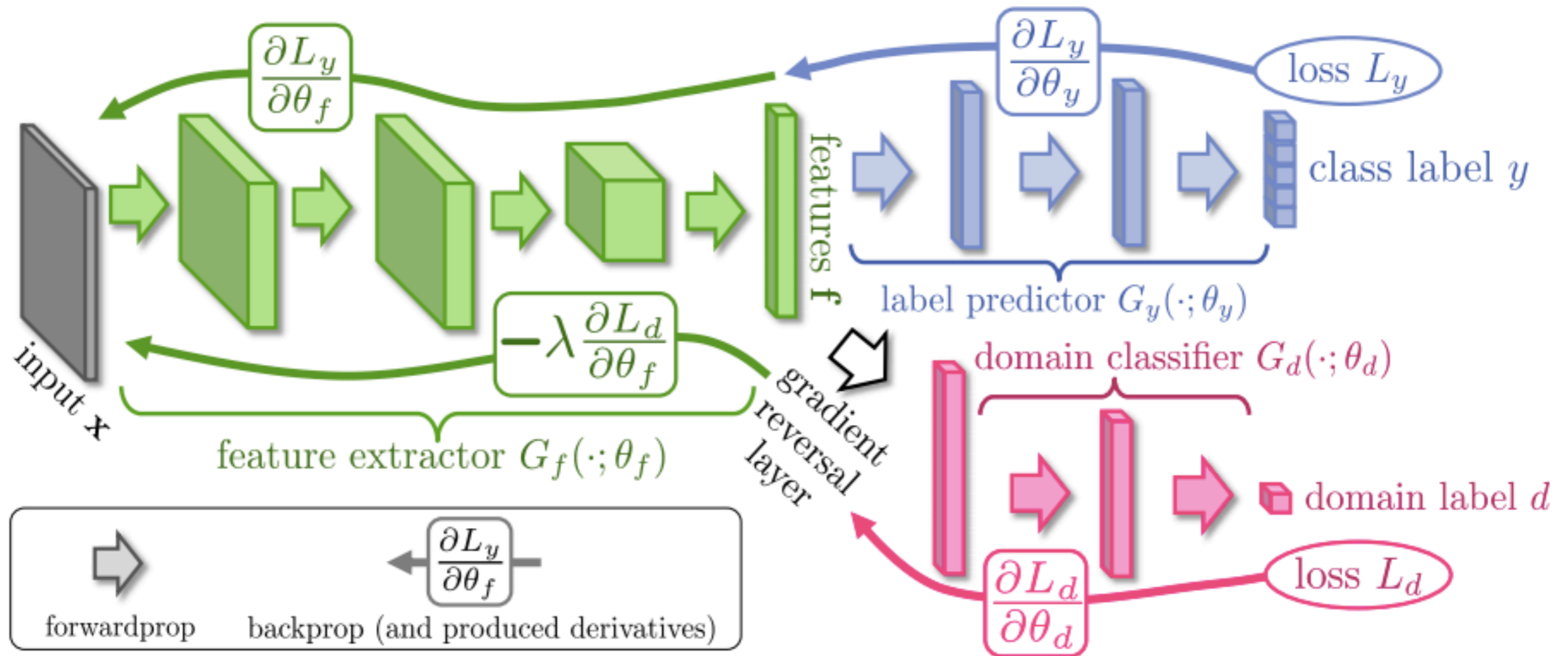
[Ganin et al. ICML15] **Source** or **Target**?



Domain-Invariant Representation Learning

Domain-Adversarial Neural Networks (DANN):

- Goal 1: Learn discriminative features for the target task of interest
- Goal 2: Learn domain-invariant features to confuse the domain classifier

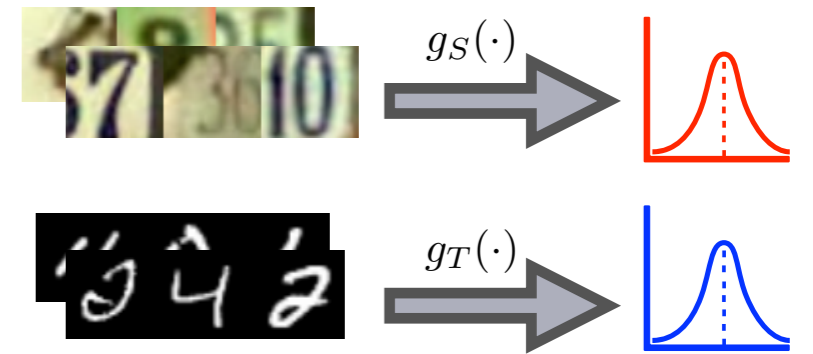


[Figure credit: Ganin et al.' 16]

A Theory of Learning from Different Domains

Theorem (Ben-David et al.' 07):

$$\varepsilon_T(h) \leq \varepsilon_S(h) + \frac{1}{2}d(S; T) + \lambda^*$$



- $\varepsilon_T(h)/\varepsilon_S(h)$: true target/source errors
- $d(S; T)$: divergence between target/source input distributions
- $\lambda^* := \min_{h' \in \mathcal{H}} \varepsilon_S(h') + \varepsilon_T(h')$: optimal joint error

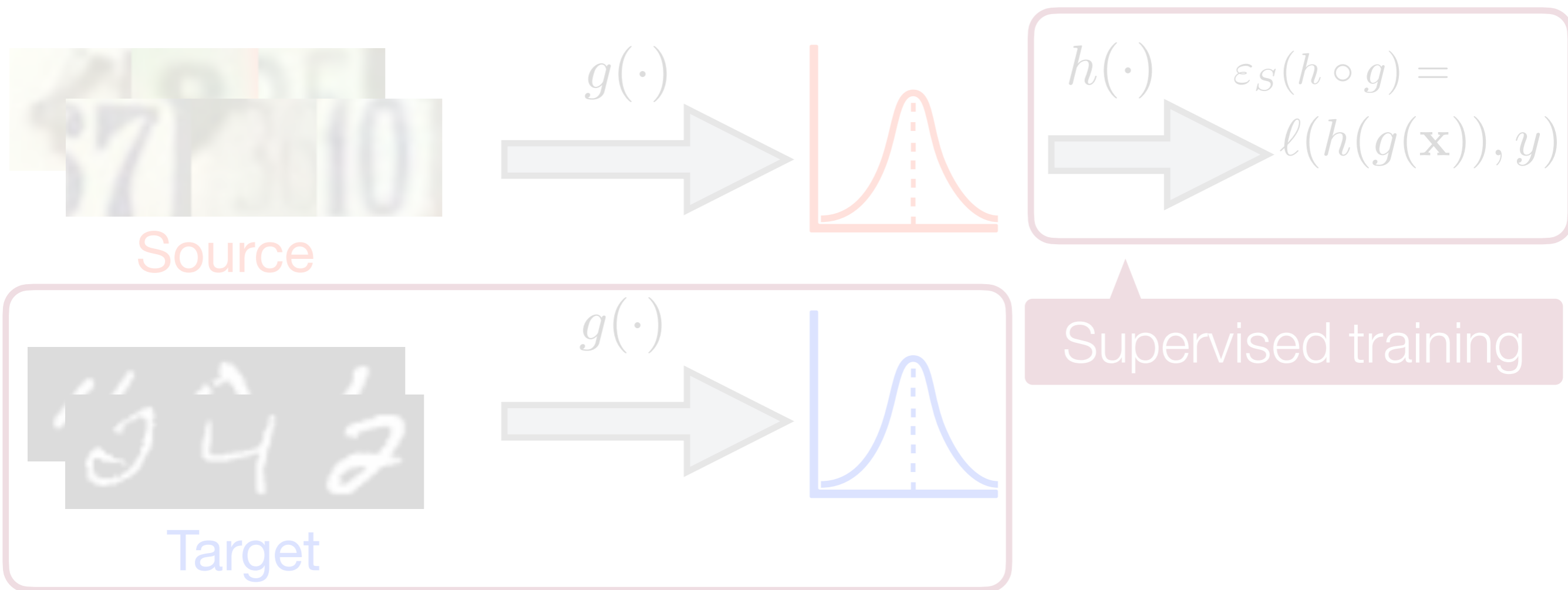
Bound-minimizing algorithm:

$$\min \varepsilon_S(h) + \frac{1}{2}d(S; T)$$

error minimization on the **source** domain

distribution matching between **source** and **target** domains

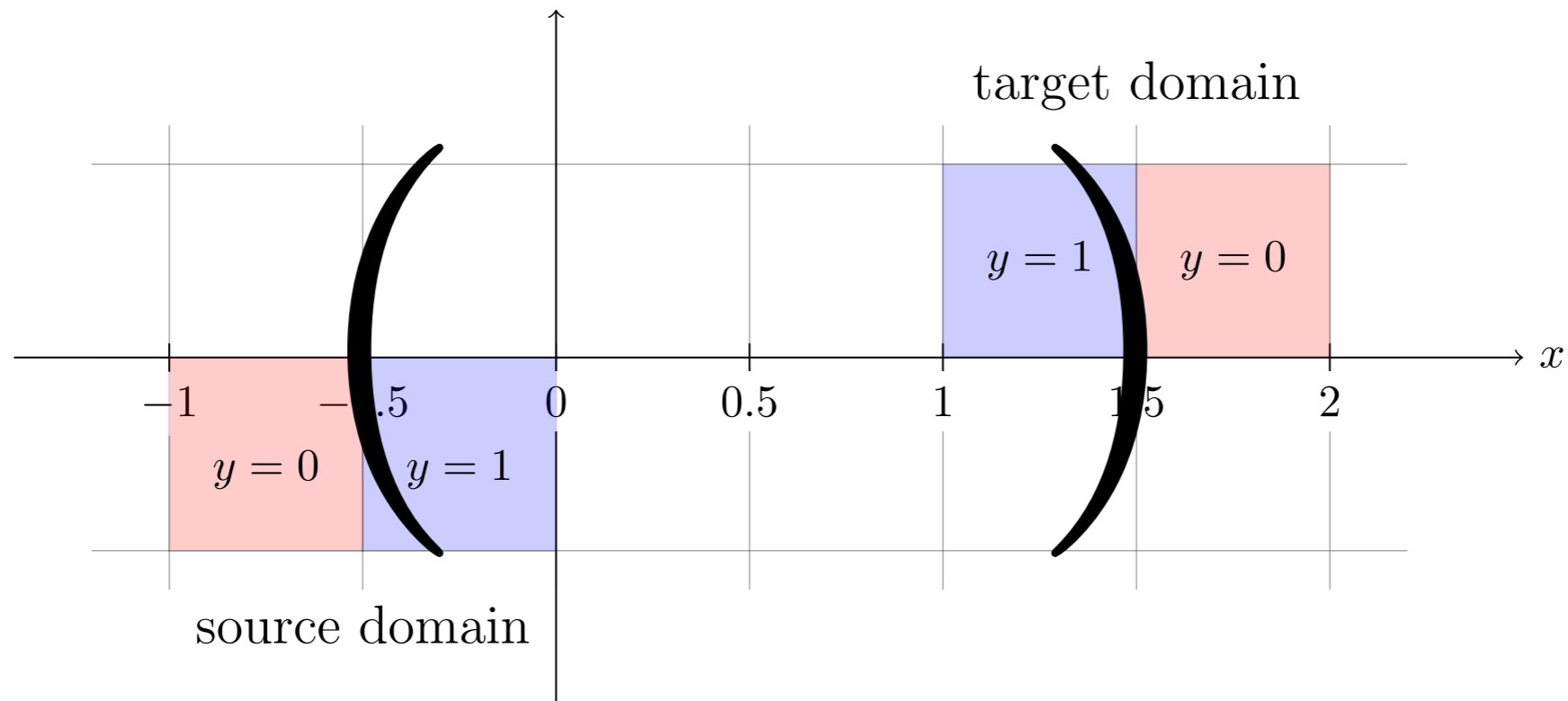
Domain-Invariant Representation Learning



Question: Is it guaranteed to generalize on target domain?

A simple 1D adaptation example

Before adaptation:



Source: $\mathcal{D}_S = U(-1, 0)$

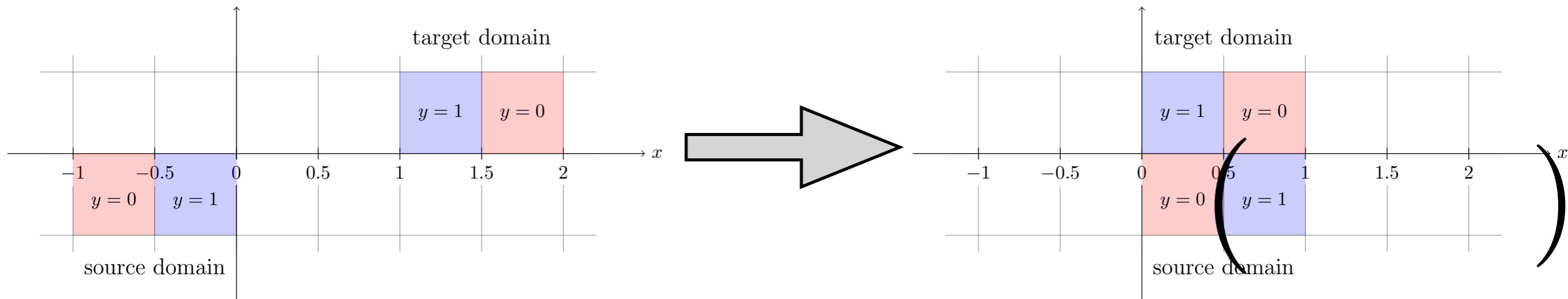
Target: $\mathcal{D}_T = U(1, 2)$

$$h^*(x) = 1 \text{ iff } x \in (-1/2, 3/2)$$

$$\lambda^* = \min_{h'} \varepsilon_S(h') + \varepsilon_T(h') = 0$$

A simple 1D adaptation example

After adaptation:



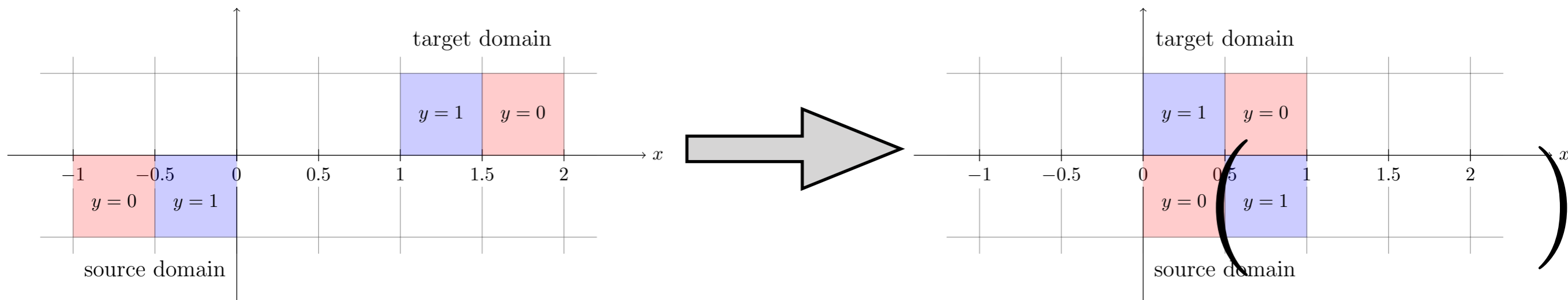
Bound-minimizing algorithm: $\min \varepsilon_S(h) + \frac{1}{2}d(S; T)$

Source: $\mathcal{D}'_S = U(0, 1)$
Target: $\mathcal{D}'_T = U(0, 1)$ $\Rightarrow d(S, T) = 0$

If $\varepsilon_S(h) = 0$, then $\varepsilon_T(h) = 1$!

A simple 1D adaptation example

After adaptation:



Bound-minimizing algorithm: $\min \epsilon_S(h) + \frac{1}{2}d(S; T)$

What's wrong with this example?

$$\epsilon_T(h) \leq \epsilon_S(h) + \frac{1}{2}d(S; T) + \lambda^*$$

$$\begin{array}{ccc} \parallel & & \parallel \\ 0 & & 0 \end{array}$$

In fact, $\forall h \in 2^{[0,1]}$, $\epsilon_S(h) + \epsilon_T(h) = 1 \implies \lambda^* = 1$

Key Message: Matching features is not sufficient!

An information-theoretic lower bound

Well, that's just a worst-case example, in general it works...

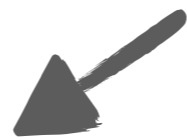
Consider the general adaptation scenario:

$$X \xrightarrow{g} Z \xrightarrow{h} \hat{Y}$$

- $g(\cdot)$ (nonlinear) feature transformation
- $h(\cdot) \in \{0, 1\}$ (randomized) classification function

Theorem [Zhao et al., ICML 19]: suppose the Markov chain holds and $d(\mathcal{D}_S^Y, \mathcal{D}_T^Y) \geq d(\mathcal{D}_S^Z, \mathcal{D}_T^Z)$, then:

$$\varepsilon_S(h \circ g) + \varepsilon_T(h \circ g) \geq \frac{1}{2} \left(d(\mathcal{D}_S^Y, \mathcal{D}_T^Y) - d(\mathcal{D}_S^Z, \mathcal{D}_T^Z) \right)^2$$



Distance between marginal label distributions

Distance between feature distributions

Key Message: the better the distribution alignment & the source risk, the worse the target risk

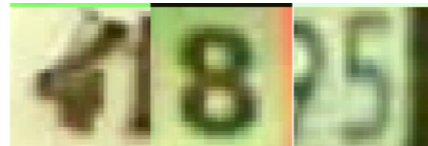
Empirical results

Digit classification:

MNIST



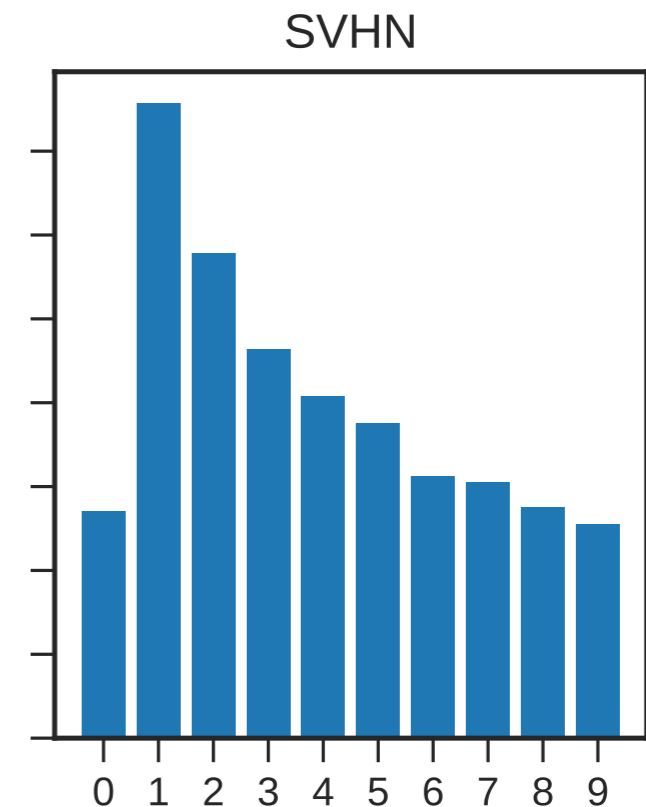
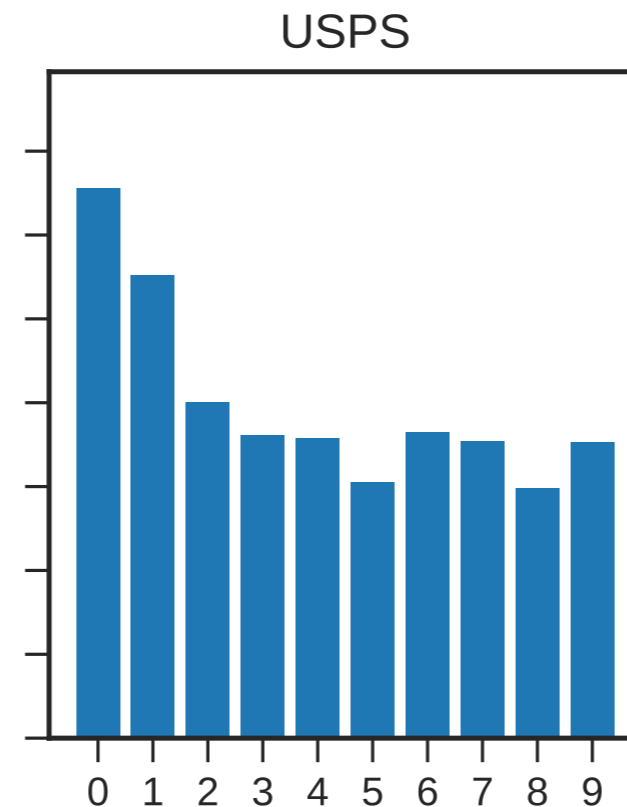
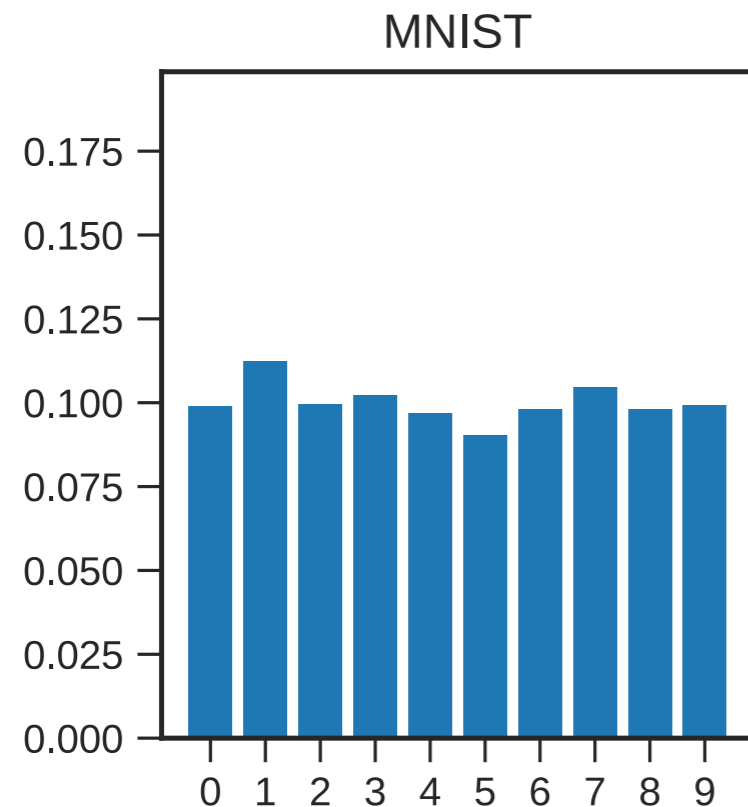
SVHN



USPS



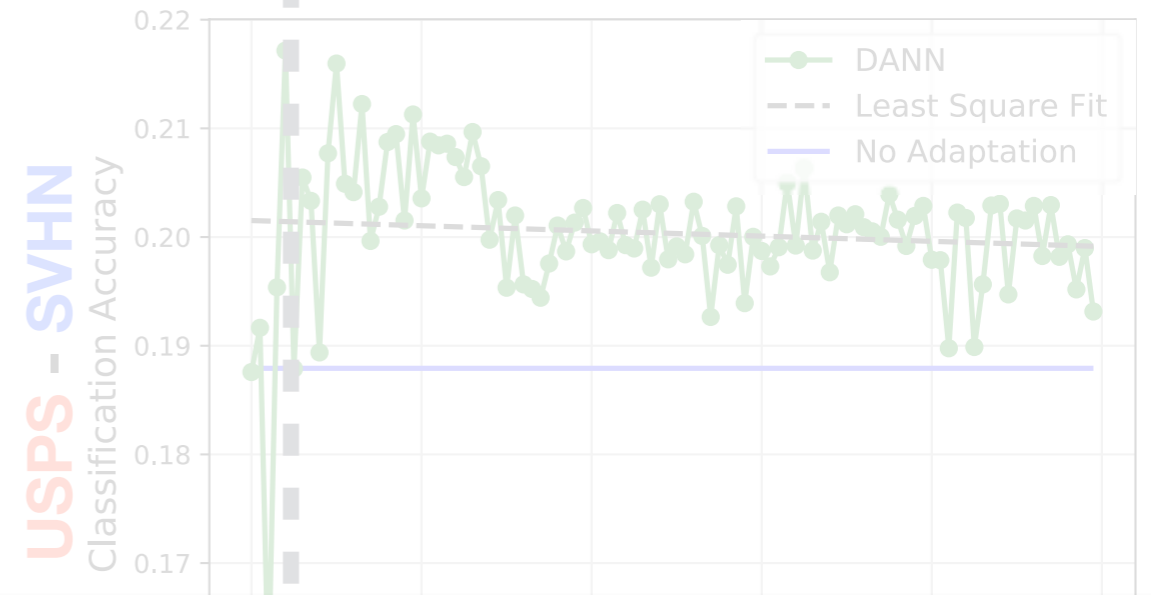
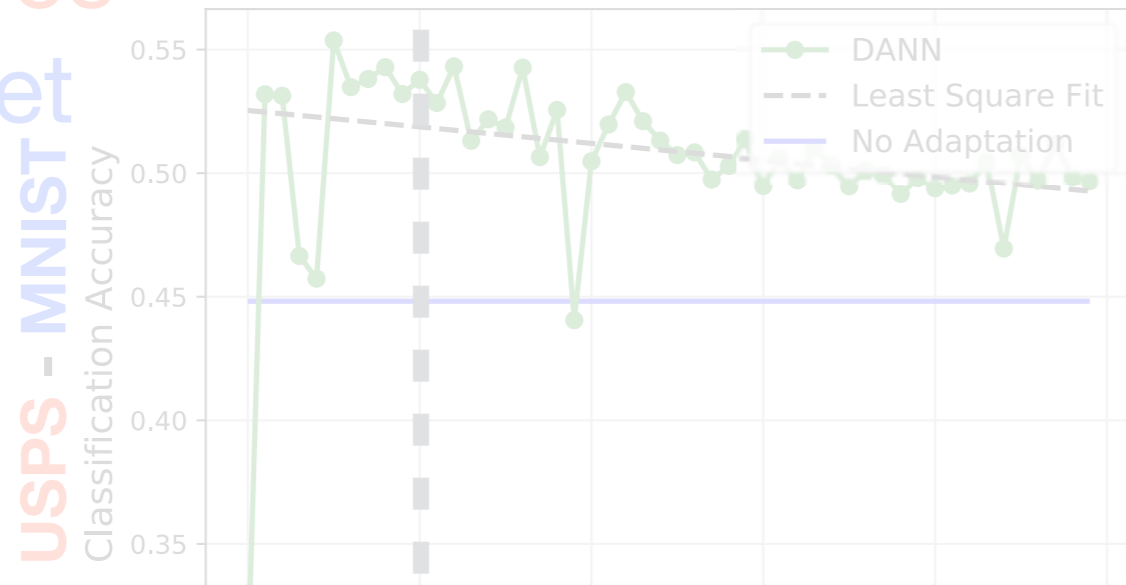
Marginal label distribution:



Negative transfer between source and target

Source

Target



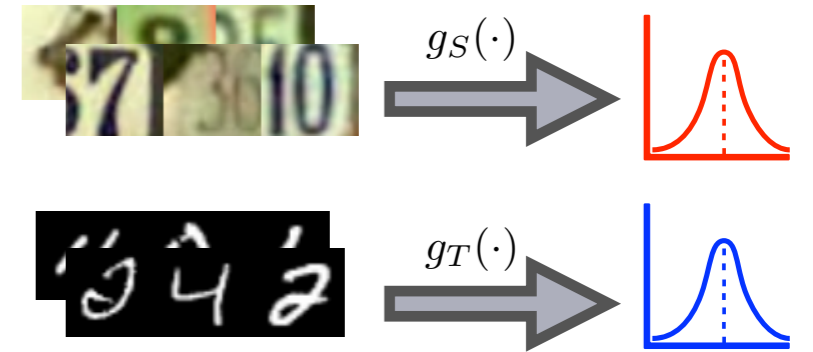
Negative Transfer: When marginal label distributions differ, invariant representations & small training error implies large target error



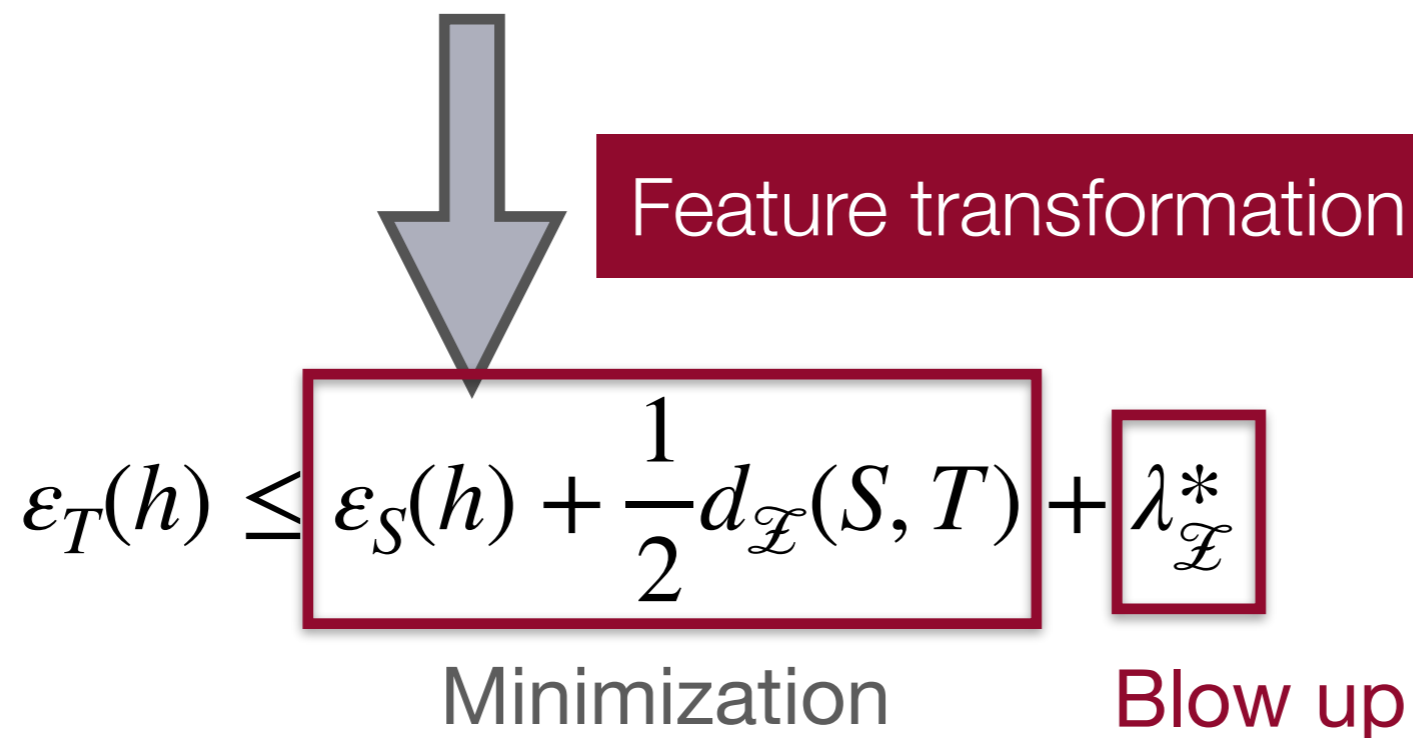
A Theory of Learning from Different Domains

Theorem (Ben-David et al.' 07):

$$\varepsilon_T(h) \leq \varepsilon_S(h) + \frac{1}{2}d_{\mathcal{X}}(S, T) + \lambda_{\mathcal{X}}^*$$



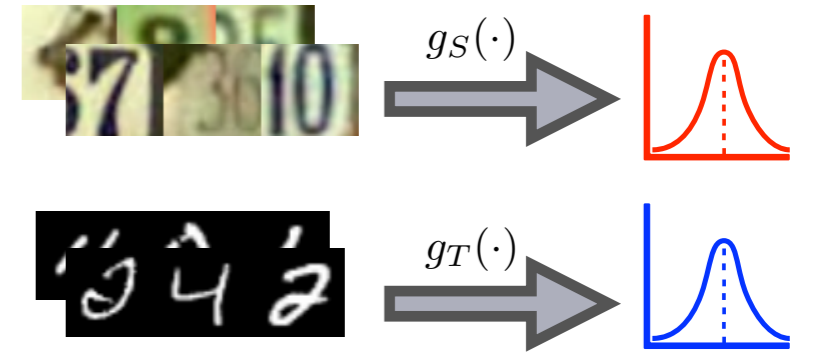
- $\varepsilon_T(h)/\varepsilon_S(h)$: true target/source errors
- $d_{\mathcal{X}}(S, T)$: divergence between target/source input distributions **over \mathcal{X}**
- $\lambda_{\mathcal{X}}^* := \min_{h': \mathcal{X} \rightarrow \{0,1\}} \varepsilon_S(h') + \varepsilon_T(h')$: optimal joint error **obtainable from \mathcal{X}**



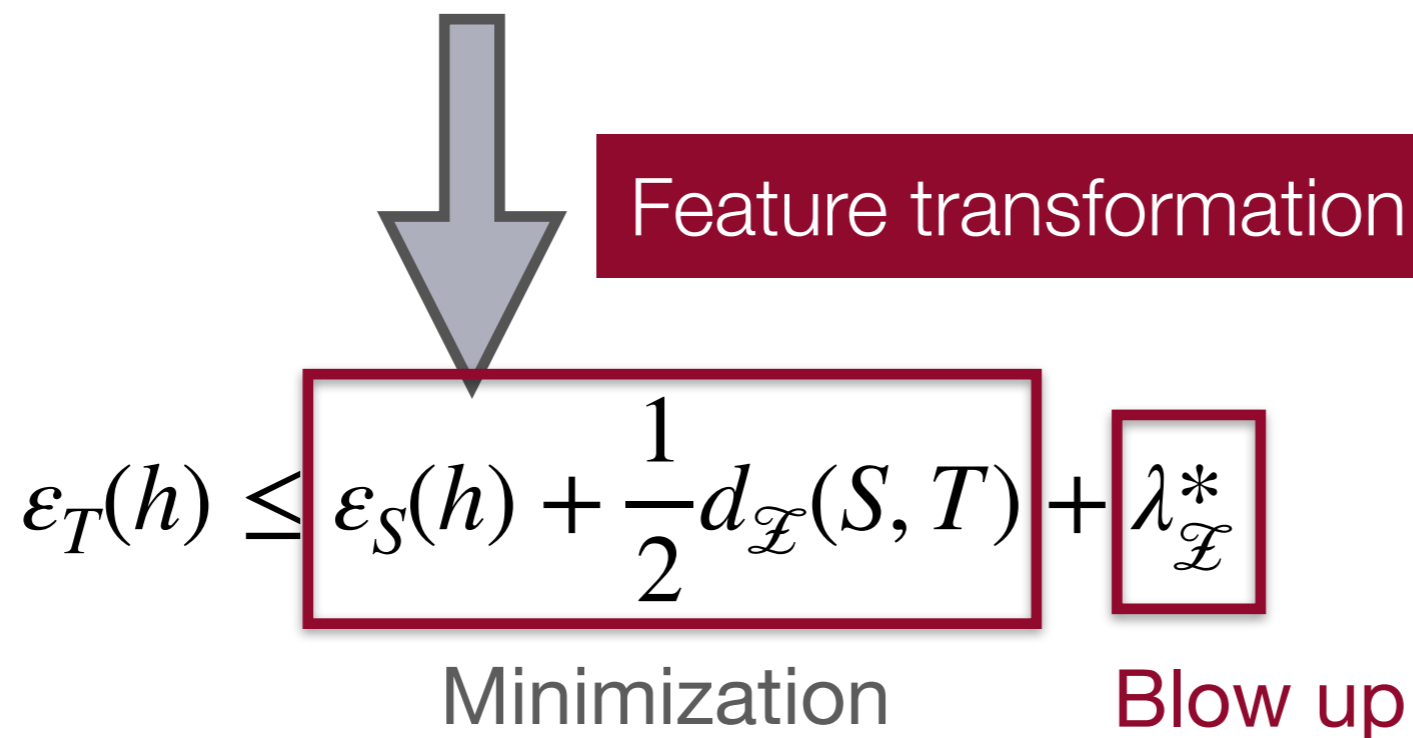
A Theory of Learning from Different Domains

Theorem (Ben-David et al.' 07):

$$\varepsilon_T(h) \leq \varepsilon_S(h) + \frac{1}{2}d_{\mathcal{X}}(S, T) + \lambda_{\mathcal{X}}^*$$



- $\varepsilon_T(h)/\varepsilon_S(h)$: true target/source errors
- $d_{\mathcal{X}}(S, T)$: divergence between target/source input distributions **over \mathcal{X}**
- $\lambda_{\mathcal{X}}^* := \min_{h': \mathcal{X} \rightarrow \{0,1\}} \varepsilon_S(h') + \varepsilon_T(h')$: optimal joint error **obtainable from \mathcal{X}**



A Decomposition of the Error Difference

Theorem: Let $\langle \mathcal{D}_S, f_S \rangle$ and $\langle \mathcal{D}_T, f_T \rangle$ be the source and target domains, respectively. For any function class $\mathcal{H} \subseteq [0, 1]^{\mathcal{X}}$, and $\forall h \in \mathcal{H}$, the following inequality holds:

$$\varepsilon_T(h) \leq \varepsilon_S(h) + d_{\tilde{\mathcal{H}}}(\mathcal{D}_S, \mathcal{D}_T) + \min \{ \mathbb{E}_{\mathcal{D}_S} [|f_S - f_T|], \mathbb{E}_{\mathcal{D}_T} [|f_S - f_T|] \}$$

where $\tilde{\mathcal{H}} := \{ \text{sgn}(|h(x) - h'(x)| - t) : h, h' \in \mathcal{H}, t \in [0, 1] \}$


- $\mathcal{D}_S(\mathcal{D}_T)$ marginal distributions over features in source (target) domains
- $f_S(f_T)$ optimal predictors from features to labels in source (target) domains
- The output of the predictor could be probabilistic, i.e., $[0, 1]$


A Decomposition of the Error Difference

Theorem: Let $\langle \mathcal{D}_S, f_S \rangle$ and $\langle \mathcal{D}_T, f_T \rangle$ be the source and target domains, respectively. For any function class $\mathcal{H} \subseteq [0, 1]^{\mathcal{X}}$, and $\forall h \in \mathcal{H}$, the following inequality holds:

$$\varepsilon_T(h) \leq \varepsilon_S(h) + d_{\tilde{\mathcal{H}}}(\mathcal{D}_S, \mathcal{D}_T) + \min \{ \mathbb{E}_{\mathcal{D}_S} [|f_S - f_T|], \mathbb{E}_{\mathcal{D}_T} [|f_S - f_T|] \}$$

where $\tilde{\mathcal{H}} := \{ \text{sgn}(|h(x) - h'(x)| - t) : h, h' \in \mathcal{H}, t \in [0, 1] \}$


$$d(\phi(X_S), \phi(X_T))$$


$$d(\mathbb{E}[Y_S | \phi(X_S)], \mathbb{E}[Y_T | \phi(X_T)])$$

- An error decomposition theorem
- The first term \sim Invariant representations
- The second term \sim Invariant predictors
- Free of the original λ^* term

Invariant Risk Minimization

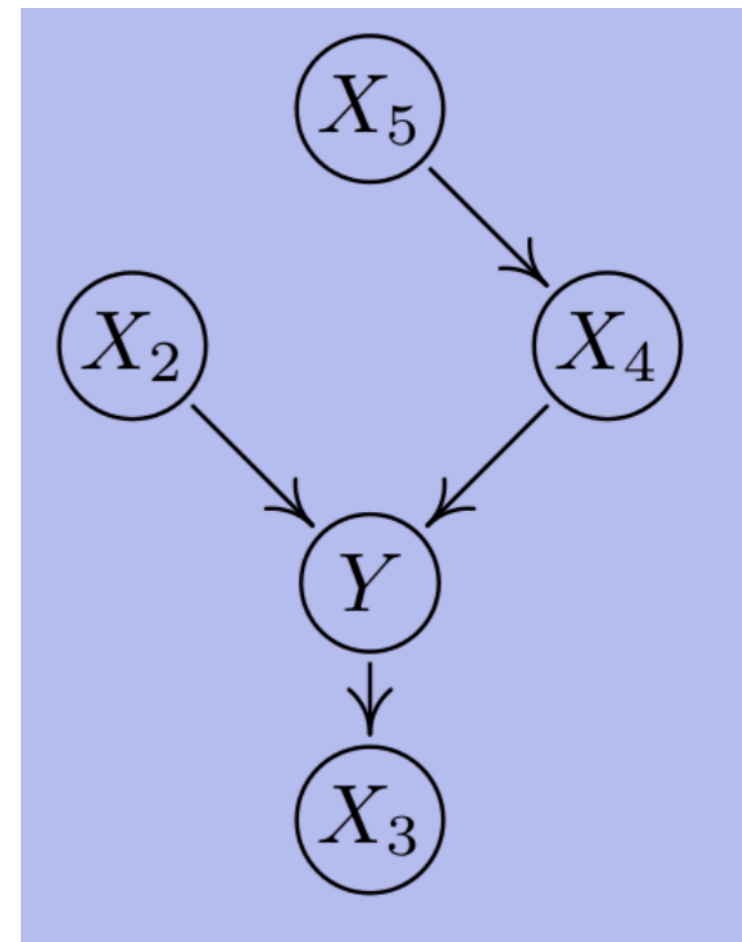
Bayesian Network: a graph used to encode the conditional independence in a joint distribution

- A directed acyclic graph \mathcal{G} , where each node $X \in \mathcal{G}$ corresponds to an RV
- Let $\text{Pa}(X)$ denote the parents of node X , then each node is associated with a conditional probability distribution (CPD):
 $\Pr(X | \text{Pa}(X))$

Let X_1, \dots, X_n be a topological ordering of all the nodes in \mathcal{G} , then the following decomposition of the joint probability holds:

$$\begin{aligned}\Pr(X_1, \dots, X_n) &= \prod_{i=1}^n \Pr(X_i | X_1, \dots, X_{i-1}) \\ &= \prod_{i=1}^n \Pr(X_i | \text{Pa}(X_i))\end{aligned}$$

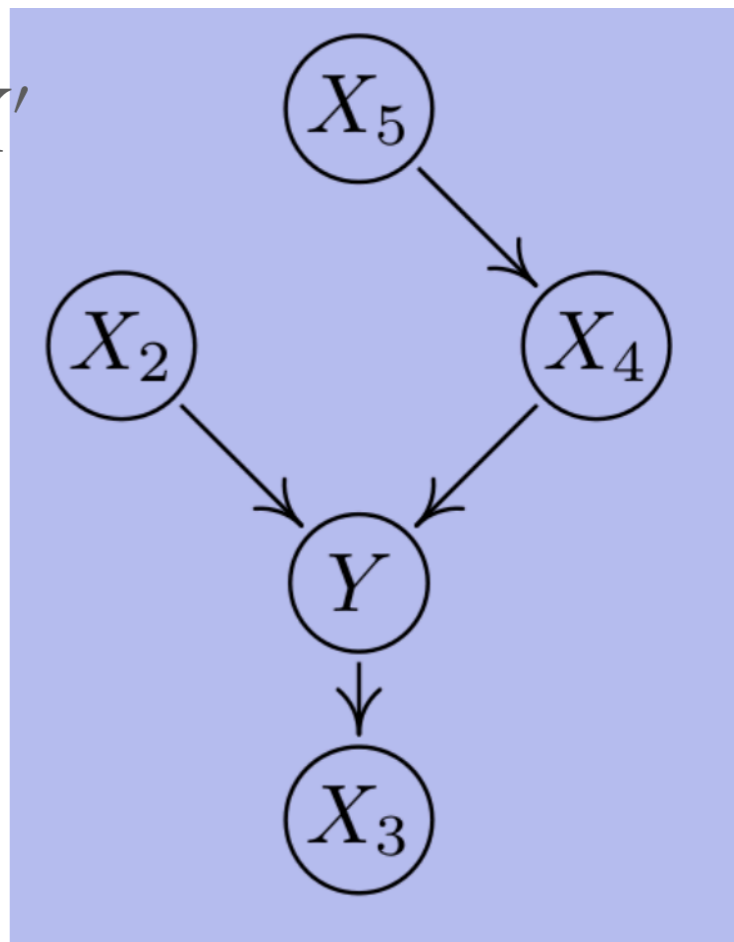
In short: every RV is conditionally independent of non-descendants given its parents.



Invariant Risk Minimization

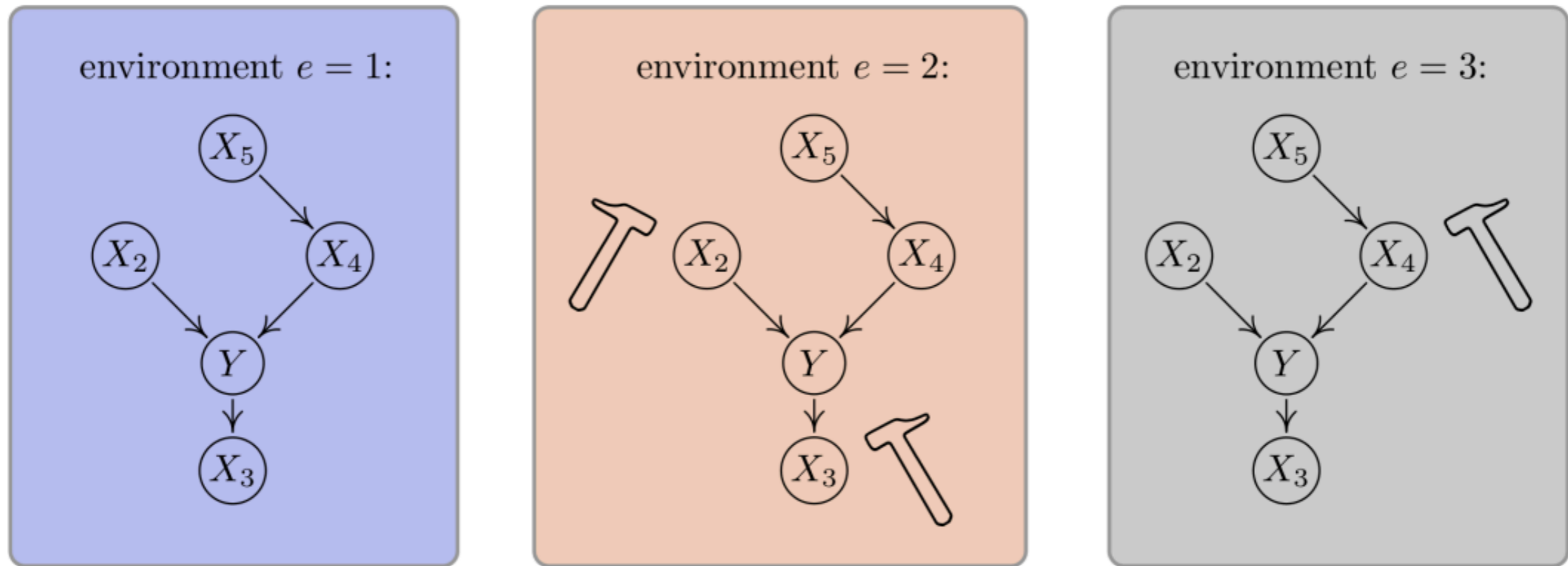
Structural Causal Model (SCM):

- A directed acyclic graph \mathcal{G} , where each node $X \in \mathcal{G}$ corresponds to an RV (called **Endogenous variables**)
- At each node X , there is also a noise RV U_X (called **Exogenous variables**)
- If $X \rightarrow X'$ in \mathcal{G} , then $\exists f : X \times U_X \rightarrow X'$ where $X' = f(X, U_X)$ (we call X is the cause of X')
- If $X \rightarrow X'$ then we say X has a causal effect on X'



Invariant Risk Minimization

Structural Causal Model (SCM):



$$Y = f(X_2, X_4)$$

Suppose from $e = 1$ we learned:

- What if in $e = 2$, we fix the value of X_2 ?
- What if in $e = 3$, we fix the value of X_4 ?
- More generally, what if we arbitrarily change the marginal distributions of X_2 , X_4 and others (except Y)?

Invariant Risk Minimization

Given the causal graph, how to find the subset of causal features?

- Find a subset of the input data, such that the optimal predictor of the target given the subset does not change across environments

Invariant Risk Minimization [Arjovsky et al. 2019]

- Going beyond linear transformations (since subset selection = projection, which is linear)

IRM objective:

$$\min_{\substack{\Phi: \mathcal{X} \rightarrow \mathcal{H} \\ w: \mathcal{H} \rightarrow \mathcal{Y}}} \sum_{e \in \mathcal{E}_{\text{tr}}} R^e(w \circ \Phi)$$

subject to $w \in \arg \min_{\bar{w}: \mathcal{H} \rightarrow \mathcal{Y}} R^e(\bar{w} \circ \Phi)$, for all $e \in \mathcal{E}_{\text{tr}}$.

- The same optimal prediction rule applies to all the environments
- Find data representation $\Phi(\cdot)$ that elicits such representations
- Bi-level optimization formulation

Invariant Risk Minimization

Follow-up works on relaxing/simplifying the bi-level optimization framework of IRM:

- Minimization of conditional mutual information: “Invariant Rationalization”, Chang, Zhang, Jaakkola, ICML’ 20
- Minimizing the variance of empirical risks from multiple domains: “Out-of-Distribution Generalization via Risk Extrapolation”, Krueger et al., ICLR’ 20
- Combination of information bottleneck principle:
 - “Invariant Information Bottleneck for Domain Generalization”, Li et al., AAAI’ 22,
 - “Invariance Principle Meets Information Bottleneck for Out-of-Distribution Generalization”, NeurIPS’ 21
-

But,

- Again, it does not guarantee small error difference
- Fails when different distributions are not from the same SCM
- Enforcing the invariant optimal predictor is only a necessary condition, but not sufficient

A Unified Perspective from Representation Learning

Can we combine invariant representations and invariant predictors to provide a unified theory for OOD generalization?

Key observation: Let $\phi : \mathcal{X} \rightarrow \mathcal{Z}$ be the feature map, E be the RV corresponding to domain index, then

Invariant Representations: $\phi(X) \perp E$

Invariant Predictors: $Y \perp E \mid \phi(X)$

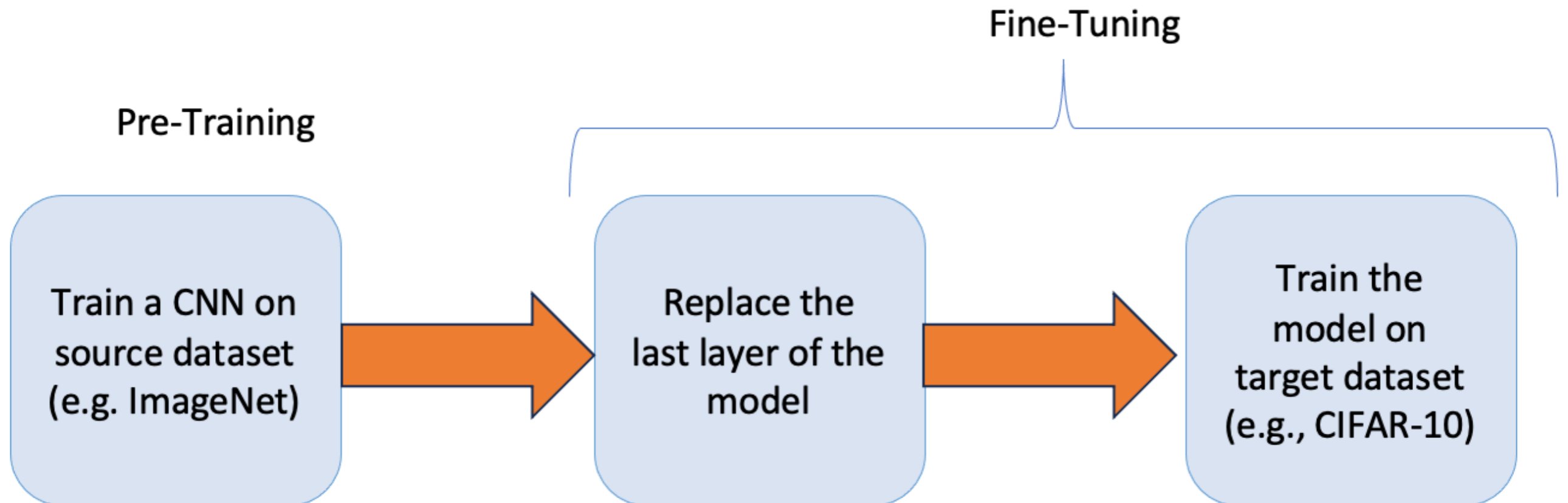
But $(\phi(X), Y)$ determines the joint distribution over features and labels, so if $(\phi(X), Y) \perp E$ then there would be no distributional shift.

In order to explain OOD generalization, we need to consider shifts in both $\phi(X)$ and $Y \mid \phi(X)$.

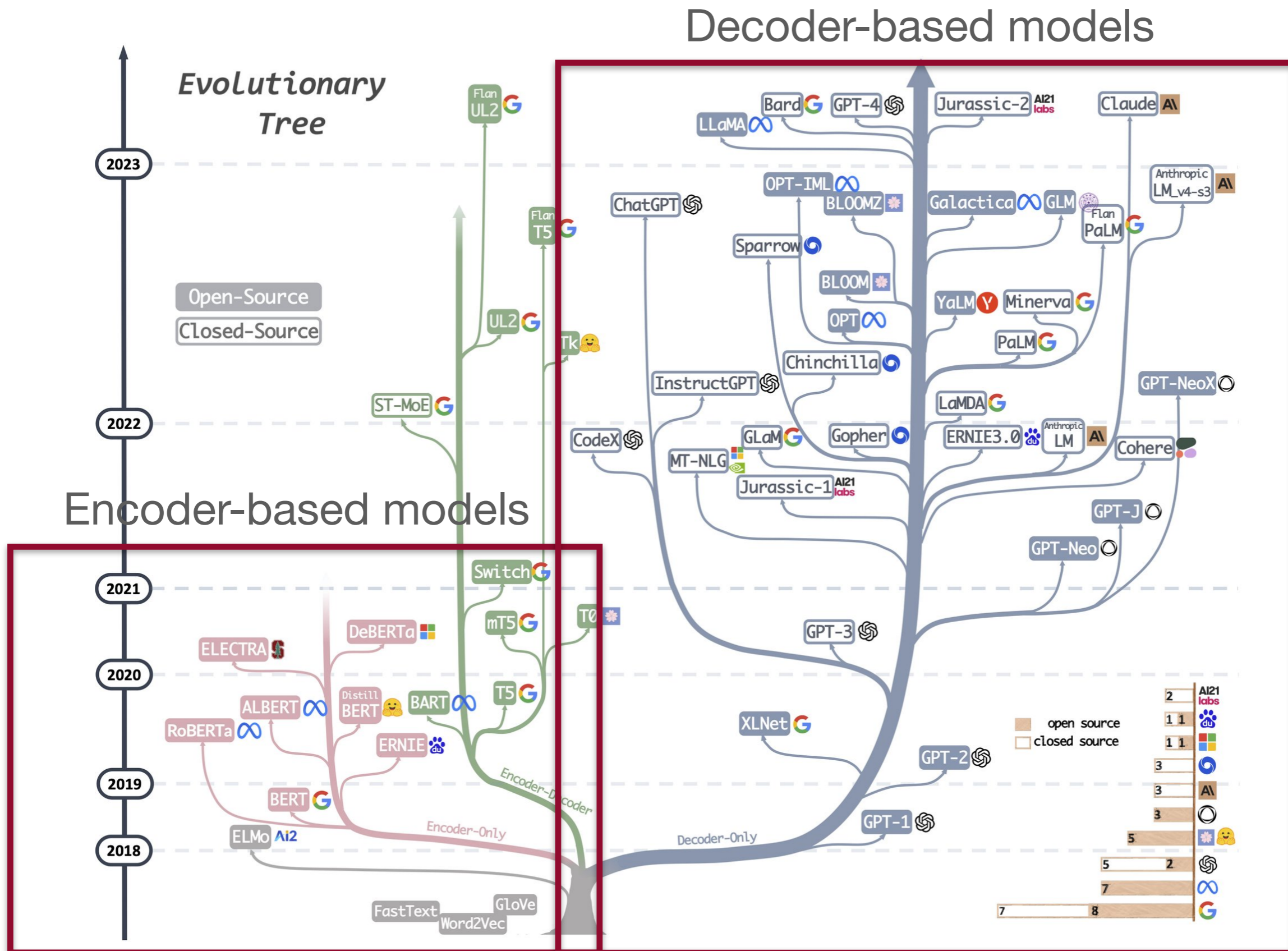
New Paradigms in Transfer Learning

Typical paradigm: pre-training + fine-tuning

1. Pre-training a feature encoder
2. Fine-tune task-specific header for prediction



New Paradigms in Transfer Learning

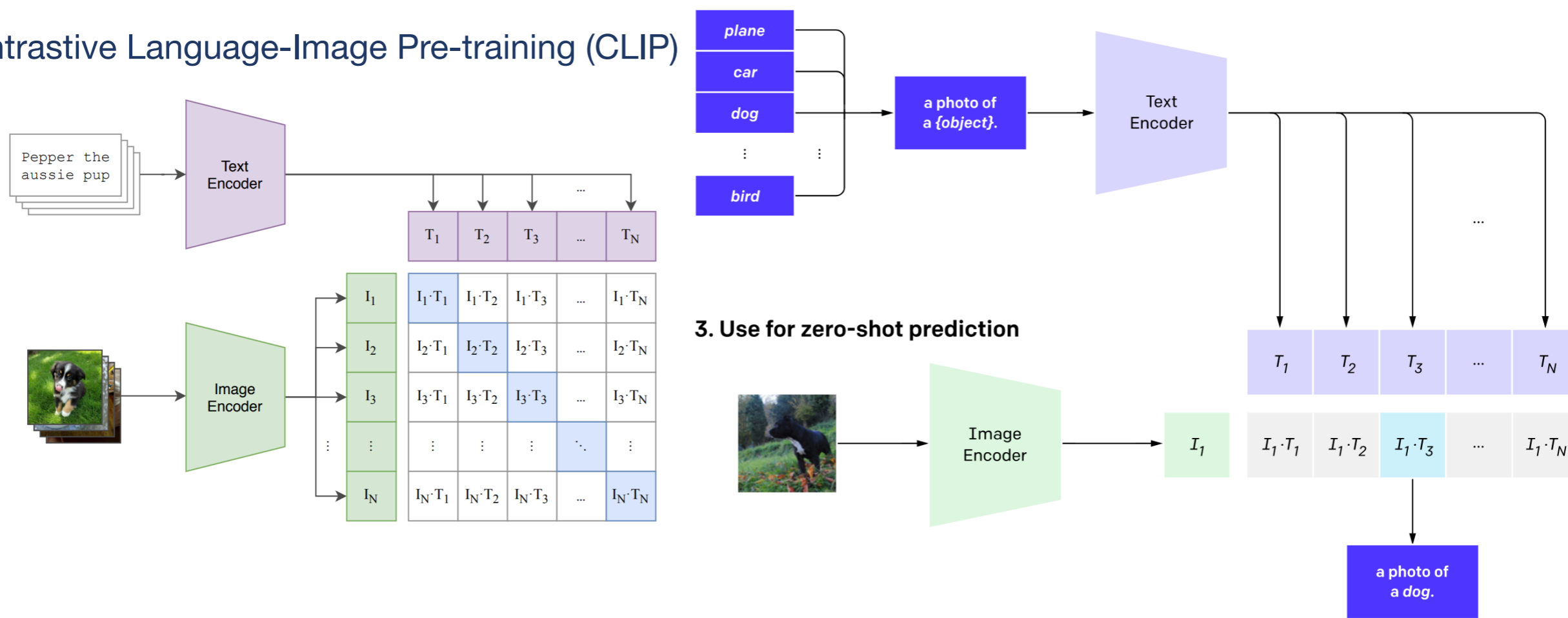


New Paradigms in Transfer Learning

New paradigm: pre-training + task mapping through prompt

1. Pre-training a feature encoder
2. Mapping the class label to a text token, and ask the model to complete the prompt

Contrastive Language-Image Pre-training (CLIP)



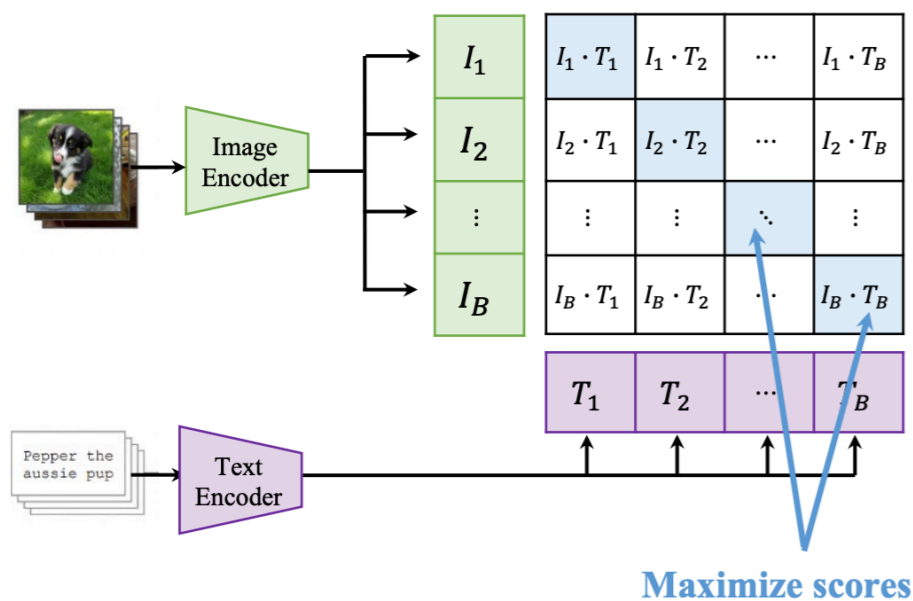
New Paradigms in Transfer Learning

Weight interpolation of model parameters

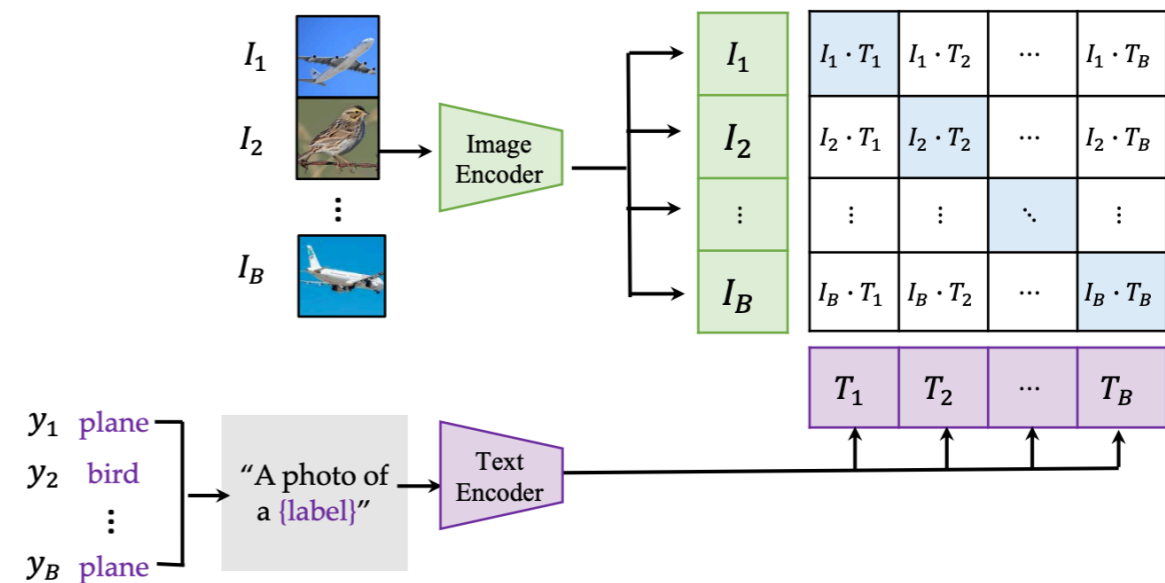
Finetune like you pretrain: Improved finetuning of zero-shot vision models

Sachin Goyal¹, Ananya Kumar², Sankalp Garg¹, Zico Kolter^{1,3}, and Aditi Raghunathan¹

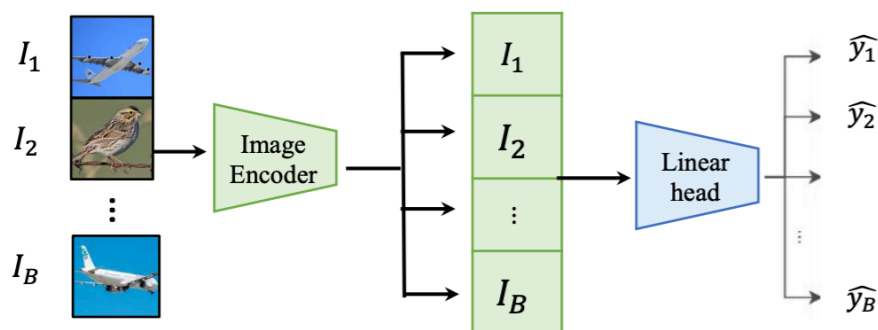
Contrastive pretraining



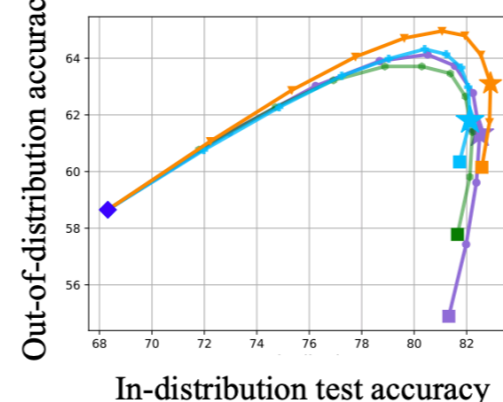
Finetune like you pretrain (FLYP)



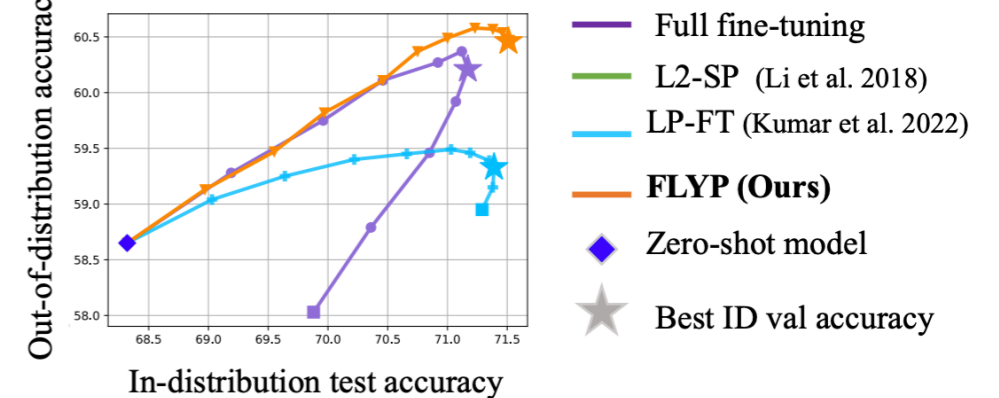
Standard finetuning



ImageNet (full)



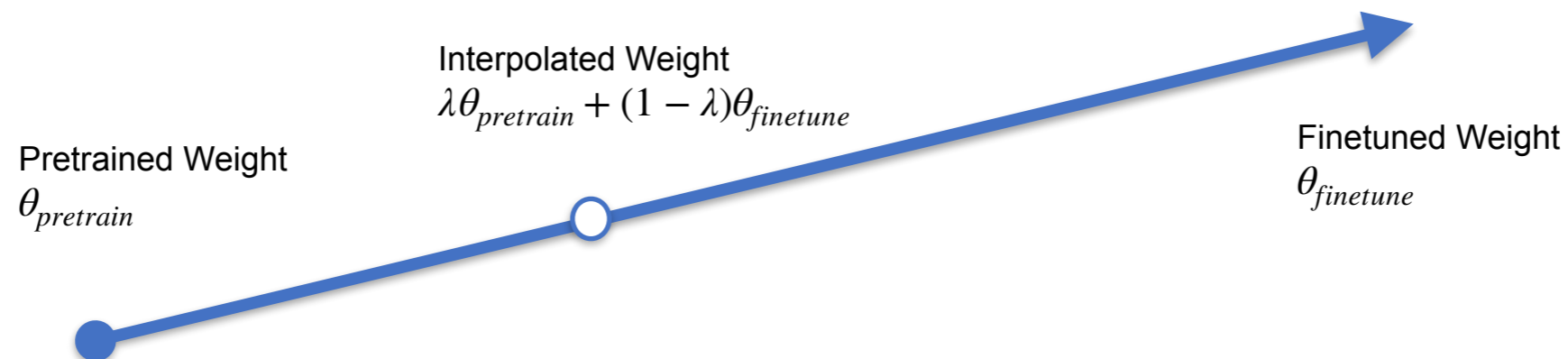
ImageNet (4-shot)



- Full fine-tuning
- L2-SP (Li et al. 2018)
- LP-FT (Kumar et al. 2022)
- **FLYP (Ours)**
- ◆ Zero-shot model
- ★ Best ID val accuracy

New Paradigms in Transfer Learning

Weight interpolation of model parameters



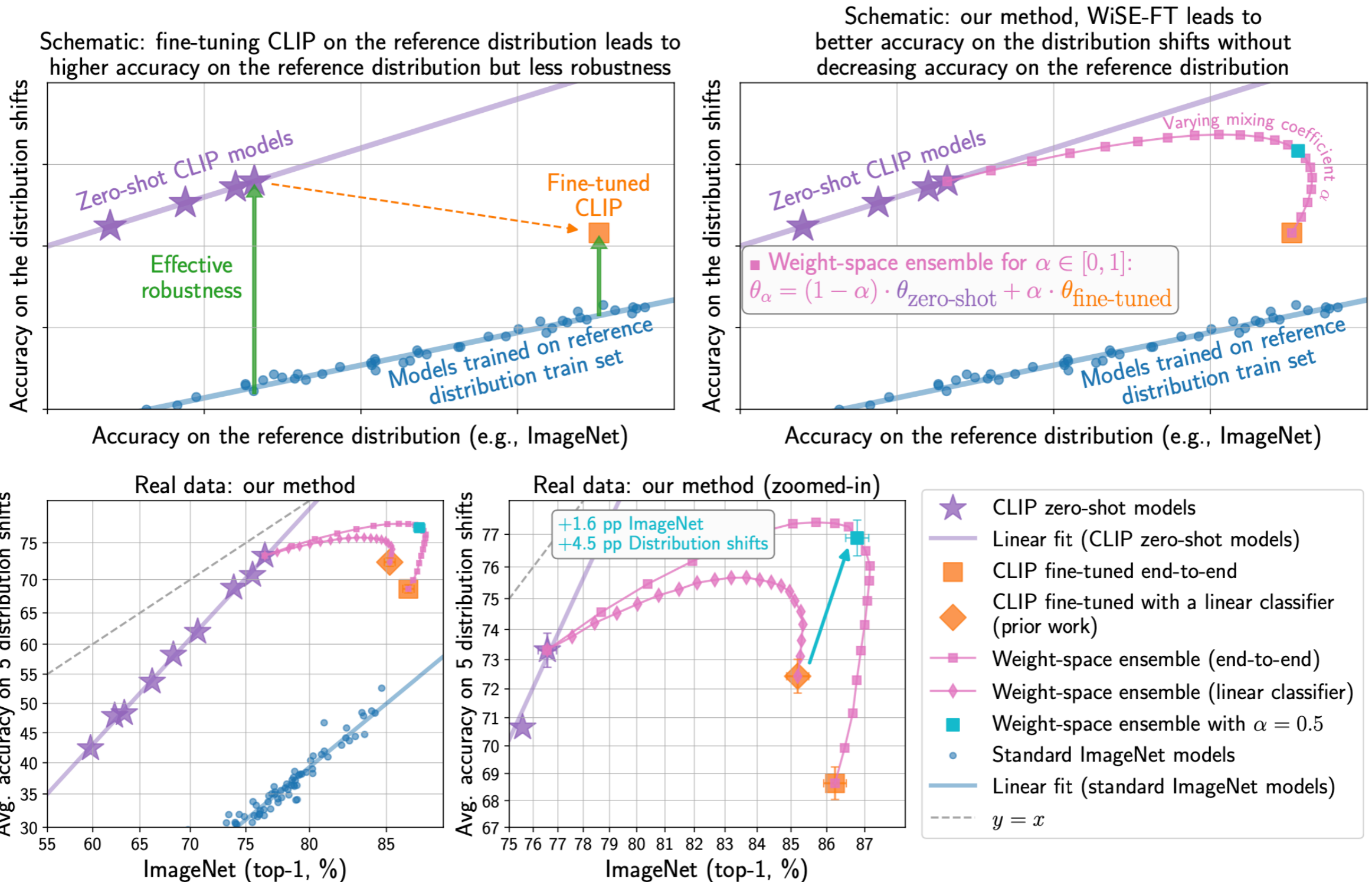
Weight Interpolation:

Simply take a linear combination of the pretrained weight and finetuned weight in the model parameter space.



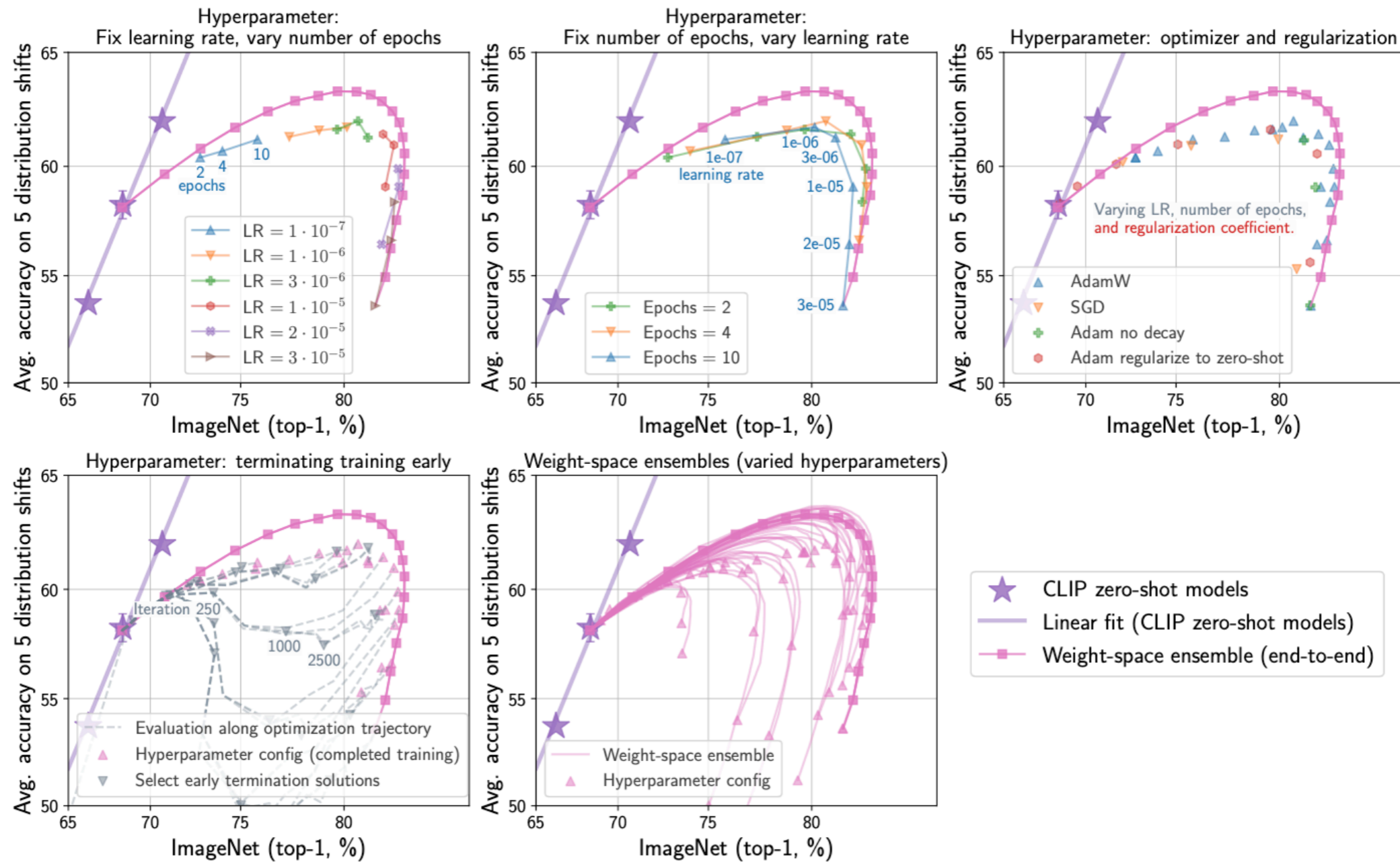
New Paradigms in Transfer Learning

Weight interpolation of model parameters



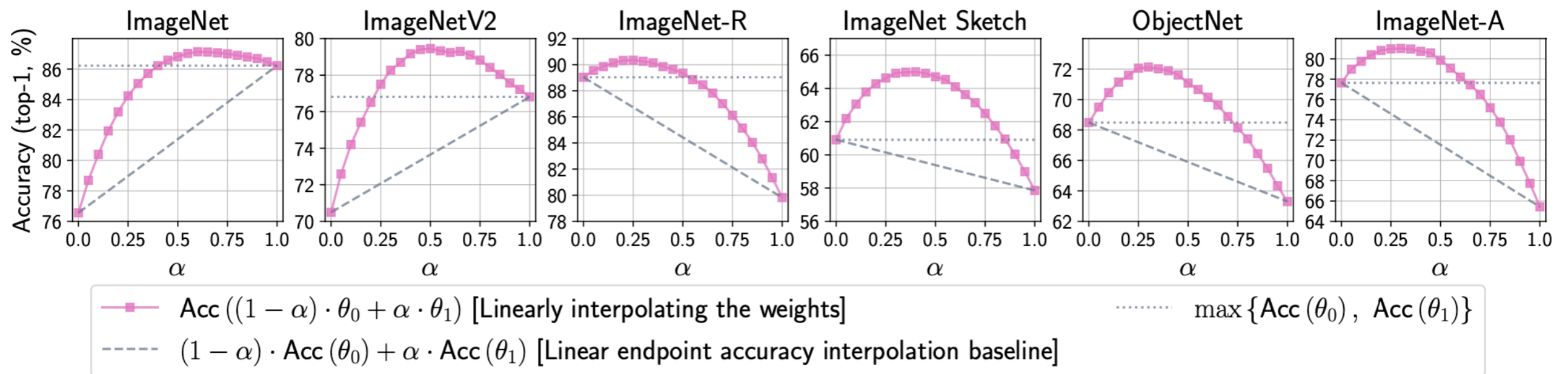
New Paradigms in Transfer Learning

Weight interpolation of model parameters



New Paradigms in Transfer Learning

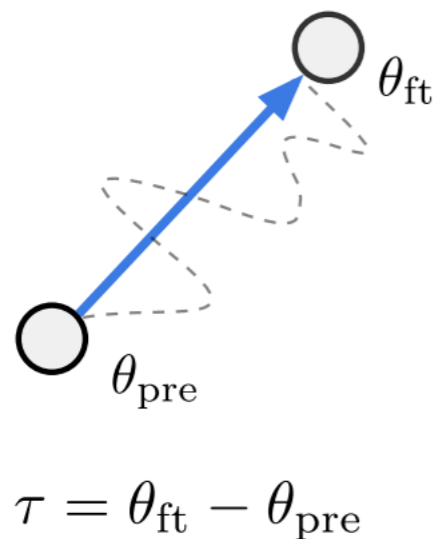
Weight interpolation of model parameters



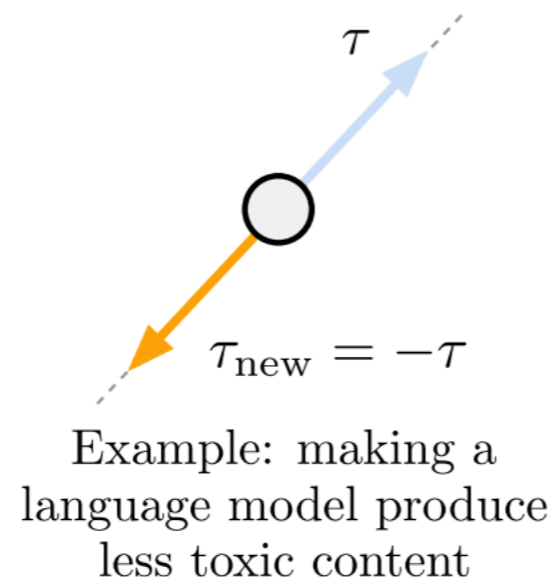
New Paradigms in Transfer Learning

Interpret each task as the difference of model parameters

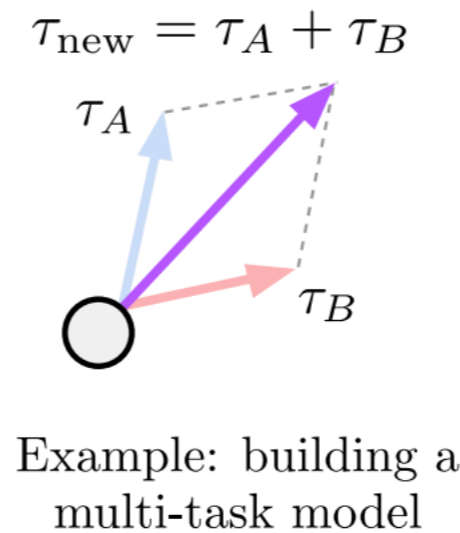
a) Task vectors



b) Forgetting via negation



c) Learning via addition



d) Task analogies

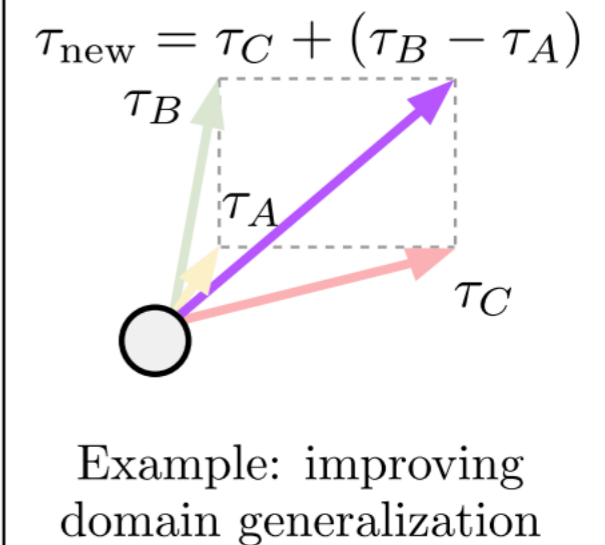


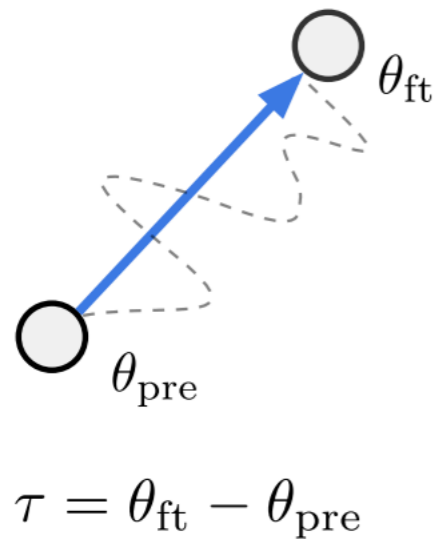
Table 3: Improving performance on target tasks with external task vectors. For four text classification tasks from the GLUE benchmark, adding task vectors downloaded from the Hugging Face Hub can improve accuracy of fine-tuned T5 models. Appendix D.6 shows additional details.

Method	MRPC	RTE	CoLA	SST-2	Average
Zero-shot	74.8	52.7	8.29	92.7	57.1
Fine-tuned	88.5	77.3	52.3	94.5	78.1
Fine-tuned + task vectors	89.3 (+0.8)	77.5 (+0.2)	53.0 (+0.7)	94.7 (+0.2)	78.6 (+0.5)

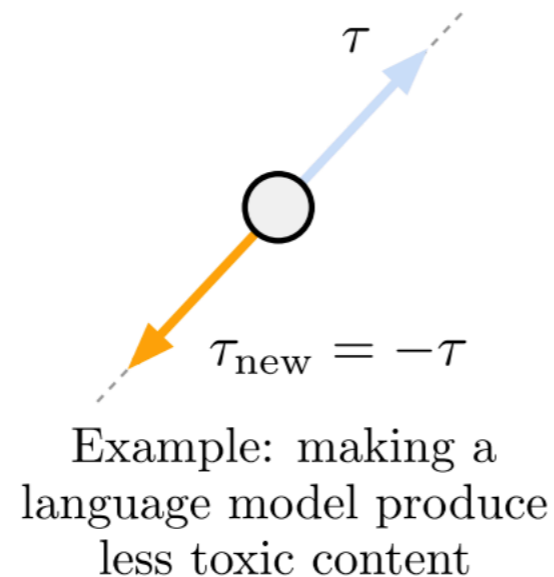
New Paradigms in Transfer Learning

Interpret each task as the difference of model parameters

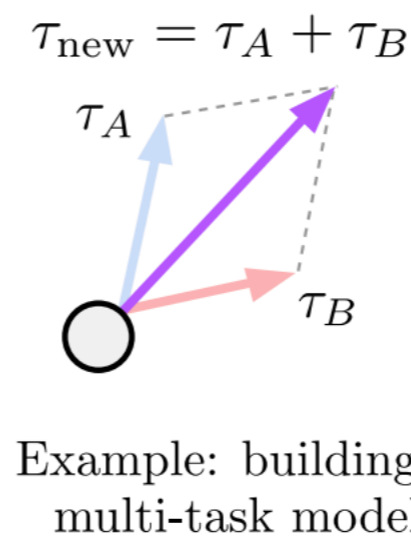
a) Task vectors



b) Forgetting via negation



c) Learning via addition



d) Task analogies

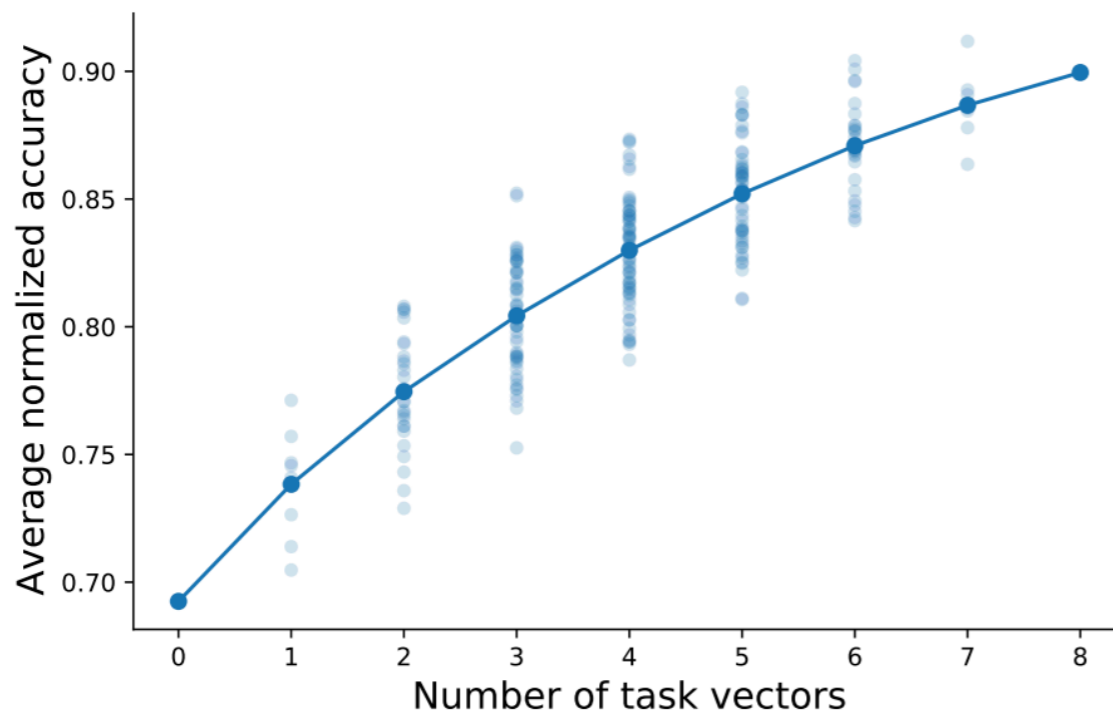
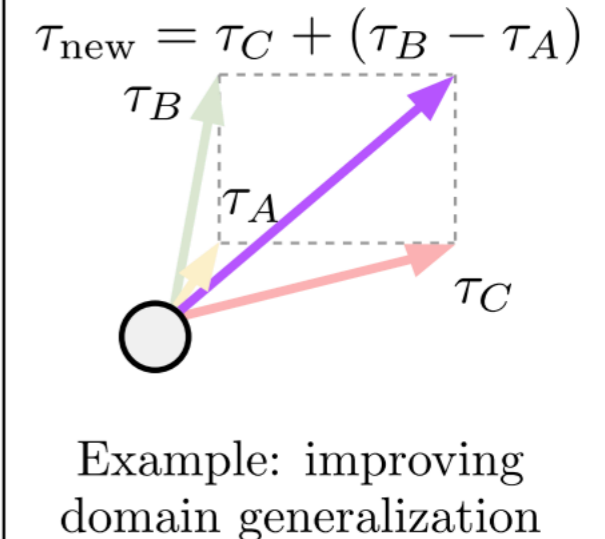


Figure 3: **Adding task vectors builds multi-task models** for image classification tasks. Accuracy is averaged over all downstream tasks. When more task vectors are available, better multi-task vectors can be built. Each point represents an experiment with a subset of the eight tasks we study, and the solid line connects the average performance for each subset size. Recall that the average normalized accuracy of using multiple fine-tuned models is always one. Additional details and experiments are in [Appendix D](#).

Open Questions

Why and when task arithmetic work?

- Does this imply a linear relationship between task and model parameter?
- How to enable more fine-grained task arithmetic?
- Is it specific to “foundation models”?
- Connection to flat-minima? Sharpness-aware minimization?
- ...