

# Domain Adaptation

CS 330

# Plan for Today

## **Domain Adaptation**

- Problem statements
- Algorithms
  - Data reweighting
  - Feature alignment
  - Domain translation

**Goal for by the end of lecture:** Understand different domain adaptation methods and when to use one vs. another

# Problem Settings Recap

## Multi-Task Learning

Solve multiple tasks  $\mathcal{T}_1, \dots, \mathcal{T}_T$  at once.

$$\min_{\theta} \sum_{i=1}^T \mathcal{L}_i(\theta, \mathcal{D}_i)$$

## Transfer Learning

Solve target task  $\mathcal{T}_b$  after solving source task(s)  $\mathcal{T}_a$   
by *transferring* knowledge learned from  $\mathcal{T}_a$

## Meta-Learning Problem

Transfer Learning with Many Source Tasks

Given data from  $\mathcal{T}_1, \dots, \mathcal{T}_n$ , solve new task  $\mathcal{T}_{\text{test}}$  more quickly / proficiently / stably

# What is domain adaptation?

Perform well on target domain  $p_T(x, y)$ ,  
using training data from source domain(s)  $p_S(x, y)$

A form of **transfer learning**, with access to target domain data during training  
("transductive" learning)

Unsupervised domain adaptation: access to unlabeled target domain data

Semi-supervised domain adaptation: access to unlabeled and labeled target domain data

Supervised domain adaptation: access to labeled target domain data.

We will focus on *unsupervised domain adaptation*.

# What is domain adaptation?

Perform well on target domain  $p_T(x, y)$ ,  
using training data from source domain(s)  $p_S(x, y)$

A form of **transfer learning**, with access to target domain data during training  
("transductive" learning)

Unsupervised domain adaptation: access to unlabeled target domain data

Common assumptions:

- Source and target domain only differ in domain of the function, i.e.  $p_S(y | x) = p_T(y | x)$
- There exists a single hypothesis with low error.

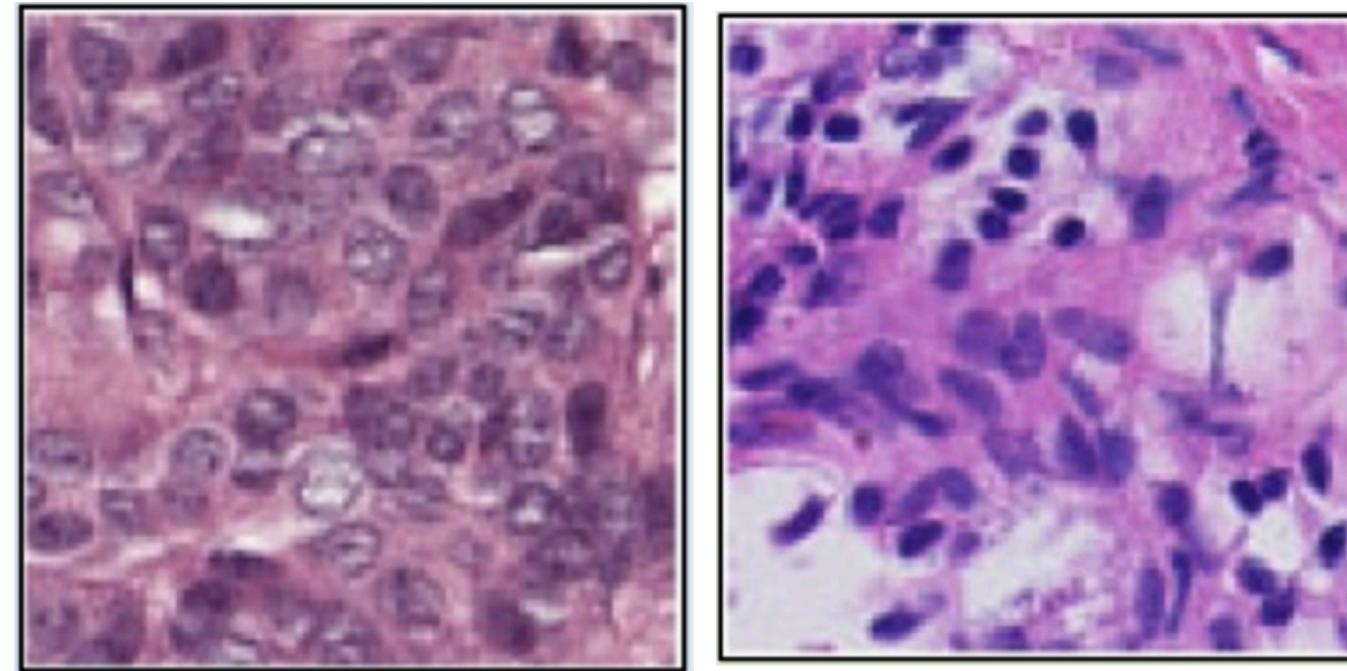
A "domain" is a special case of a "task"

A task:  $\mathcal{T}_i \triangleq \{p_i(\mathbf{x}), p_i(\mathbf{y} | \mathbf{x}), \mathcal{L}_i\}$     A domain:  $d_i \triangleq \{p_i(\mathbf{x}), p(\mathbf{y} | \mathbf{x}), \mathcal{L}\}$

# Example domain adaptation problems

## Tumor detection & classification

Source hospital Target hospital



varying imaging techniques,  
different demographics

Domains can also be:

- people/users
- points in time
- institutions  
(schools, companies, universities)

## Land use classification

Source region Target region



appearance of buildings, plants;  
weather conditions, pollution

## Text classification, generation

Source corpus Target corpus



*Simple English*  
**WIKIPEDIA**

differing sentence structure,  
vocabulary, word use



## Revisiting assumptions:

- Access to target domain data during training.
- There exists a single hypothesis  $f(y|x)$  with low error.

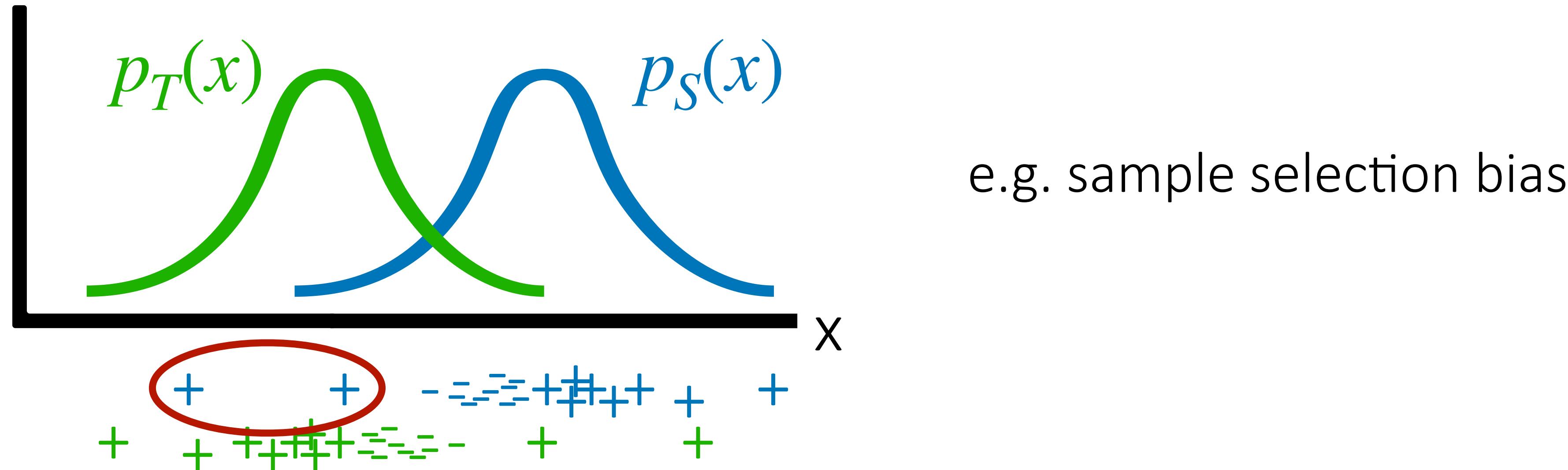
# Plan for Today

## Domain Adaptation

- Problem statements
- Algorithms
  - **Data reweighting**
  - Feature alignment
  - Domain translation

**Goal for by the end of lecture:** Understand different domain adaptation methods and when to use one vs. another

# Toy domain adaptation problem

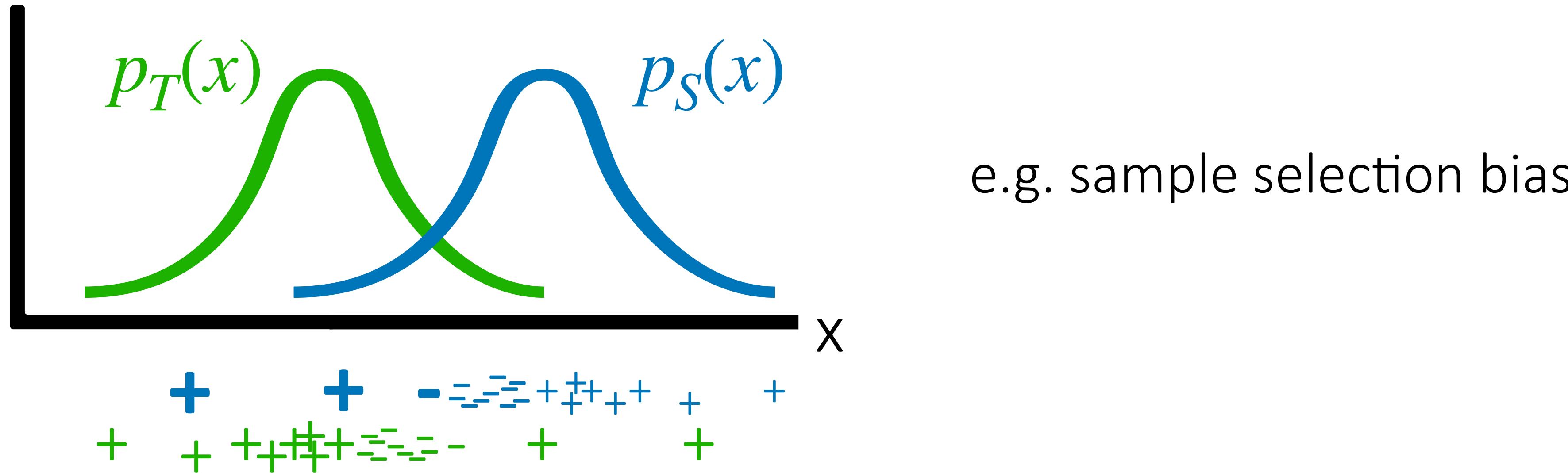


**Problem:** Classifier trained on  $p_S(x)$  pays little attention to examples with high probability under  $p_T(s)$

How can we learn a classifier that does well on  $p_T(x)$ ?

(using labeled data from  $p_S(x)$  & unlabeled data from  $p_T(s)$ )

# Toy domain adaptation problem



**Problem:** Classifier trained on  $p_S(x)$  pays little attention to examples with high probability under  $p_T(s)$

**Solution:** Upweight examples with high  $p_T(x)$  but low  $p_S(x)$

Why does this make sense mathematically?

# Domain adaptation via importance sampling

Empirical risk minimization on **source data**:  $\min_{\theta} \mathbb{E}_{p_S(x,y)}[L(f_\theta(x), y)]$

**Goal:** ERM on **target distribution**:  $\min_{\theta} \mathbb{E}_{p_T(x,y)}[L(f_\theta(x), y)]$

$$\begin{aligned}\mathbb{E}_{p_T(x,y)}[L(f_\theta(x), y)] &= \int p_T(x, y) L(f_\theta(x), y) dx dy \\ &= \int p_T(x, y) \frac{p_S(x, y)}{p_S(x, y)} L(f_\theta(x), y) dx dy \\ &= \mathbb{E}_{p_S(x,y)} \left[ \frac{p_T(x, y)}{p_S(x, y)} L(f_\theta(x), y) \right]\end{aligned}$$

Note:  $p(y|x)$  cancels out if it is the same for source & target

**Solution:** Upweight examples with high  $p_T(x)$  but low  $p_S(x)$

# Domain adaptation via importance sampling

$$\min_{\theta} \mathbb{E}_{p_S(x,y)} \left[ \frac{p_T(x)}{p_S(x)} L(f_\theta(x), y) \right]$$

How to estimate the importance weights  $\frac{p_T(x)}{p_S(x)}$ ?

Option 1: Estimate likelihoods  $p_T(x)$  and  $p_S(x)$ , then divide. But, difficult to estimate accurately.

Can we estimate the ratio *without* training a generative model?

Bayes rule:

$$p(x \mid \text{target}) = \frac{p(\text{target} \mid x)p(x)}{p(\text{target})}$$

$$p(x \mid \text{source}) = \frac{p(\text{source} \mid x)p(x)}{p(\text{source})}$$

$$\frac{p_T(x)}{p_S(x)} = \frac{p(x \mid \text{target})}{p(x \mid \text{source})} = \frac{p(\text{target} \mid x)p(\text{source})}{p(\text{source} \mid x)p(\text{target})}$$

↑                      ↑  
can estimate with    a constant  
binary classifier!

# Domain adaptation via importance sampling

$$\min_{\theta} \mathbb{E}_{p_S(x,y)} \left[ \frac{p_T(x)}{p_S(x)} L(f_\theta(x), y) \right]$$

$$\frac{p_T(x)}{p_S(x)} = \frac{p(x \mid \text{target})}{p(x \mid \text{source})} = \frac{p(\text{target} \mid x)p(\text{source})}{p(\text{source} \mid x)p(\text{target})}$$

↑  
↑  
a constant

can estimate with  
binary classifier!

## Full algorithm:

1. Train binary classifier  $c(\text{source} \mid x)$  to discriminate between source and target data.
2. Reweight or resample data  $\mathcal{D}_S$  according to  $\frac{1 - c(\text{source} \mid x)}{c(\text{source} \mid x)}$ .
3. Optimize loss  $L(f_\theta(x), y)$  on reweighted or resampled data.

# What assumption does this make?

$$\min_{\theta} \mathbb{E}_{p_S(x,y)} \left[ \frac{p_T(x)}{p_S(x)} L(f_\theta(x), y) \right]$$

Source  $p_S(x)$  needs to cover the target  $p_T(x)$ .

Formally: if  $p_T(x) \neq 0$ , then  $p_S(x) \neq 0$ .

Text classification, generation

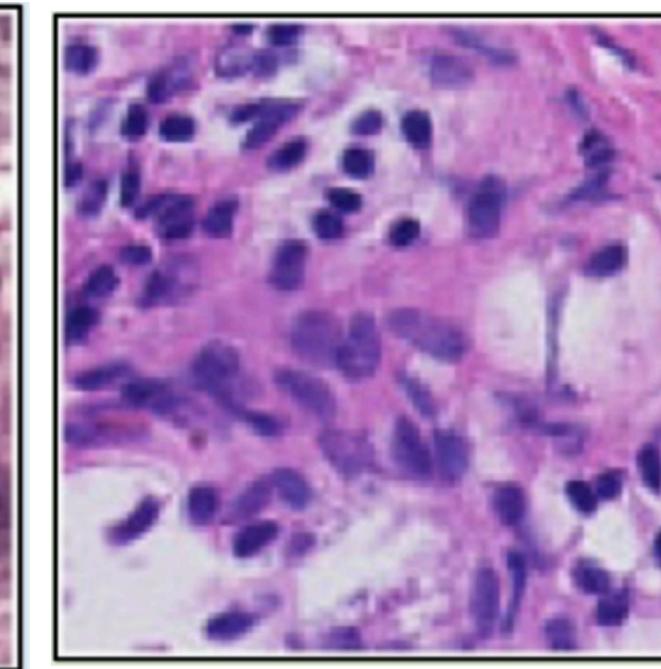
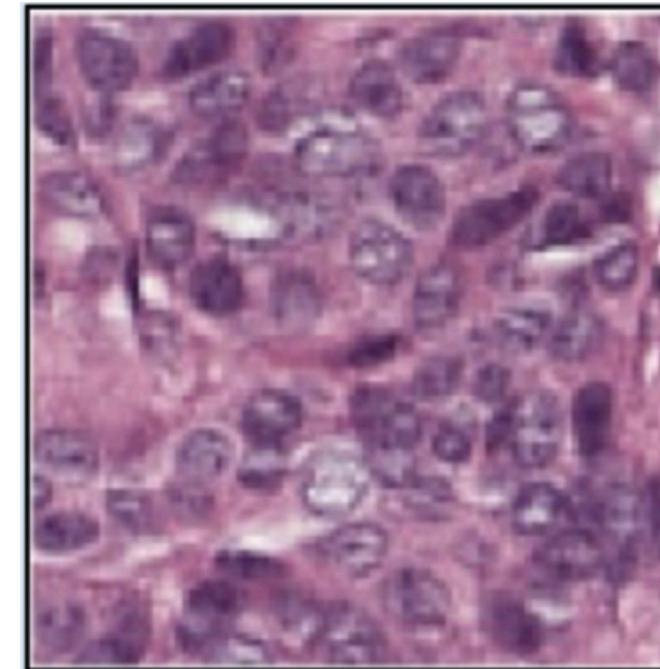
Source corpus    Target corpus



→ May have enough coverage of distr.

Tumor detection & classification

Source hospital    Target hospital



→ Source probably won't cover target distr!

# Plan for Today

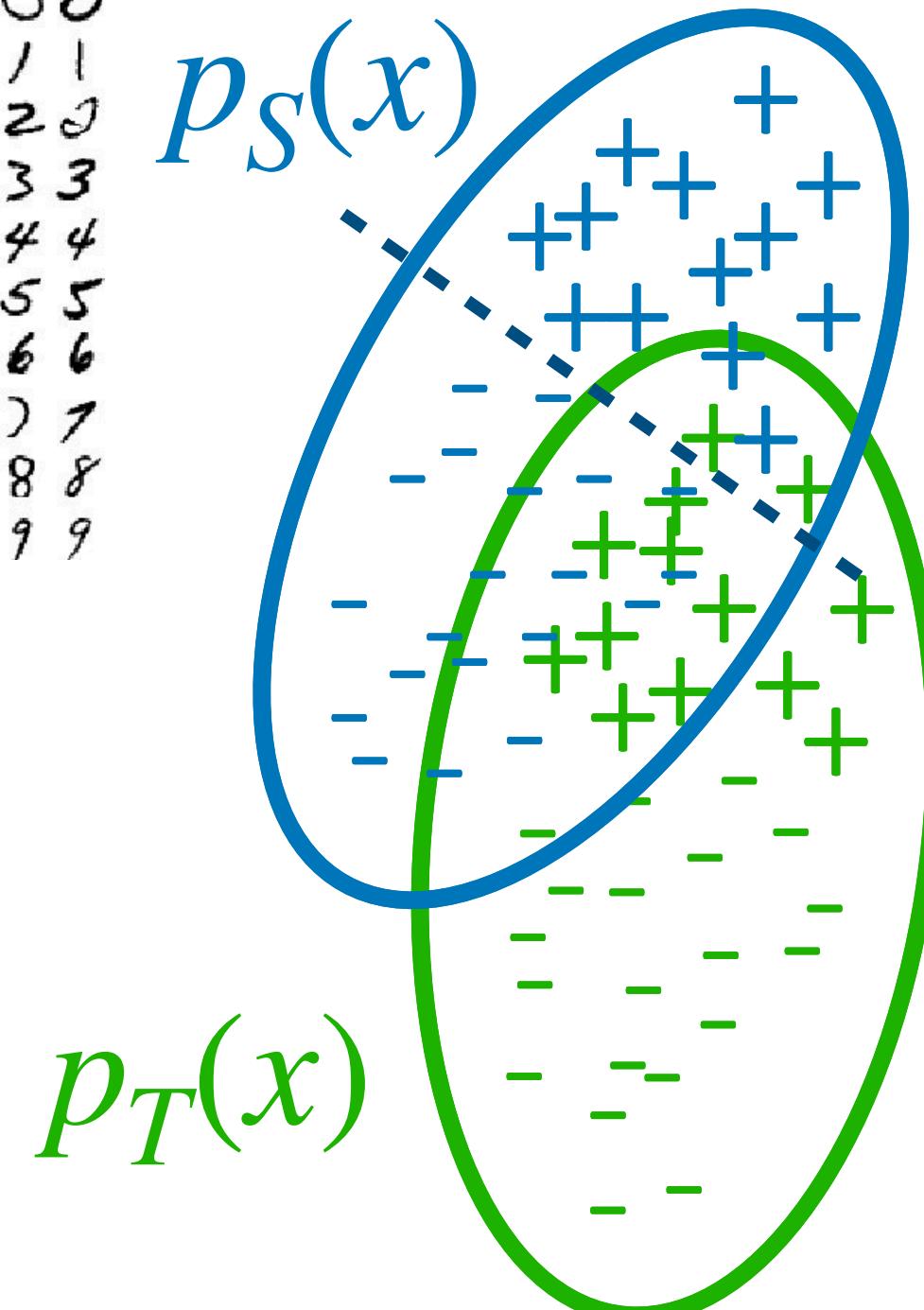
## **Domain Adaptation**

- Problem statements
- Algorithms
  - Data reweighting
  - **Feature alignment**
  - Domain translation

**Goal for by the end of lecture:** Understand different domain adaptation methods and when to use one vs. another

# Domain adaptation if support is not shared?

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



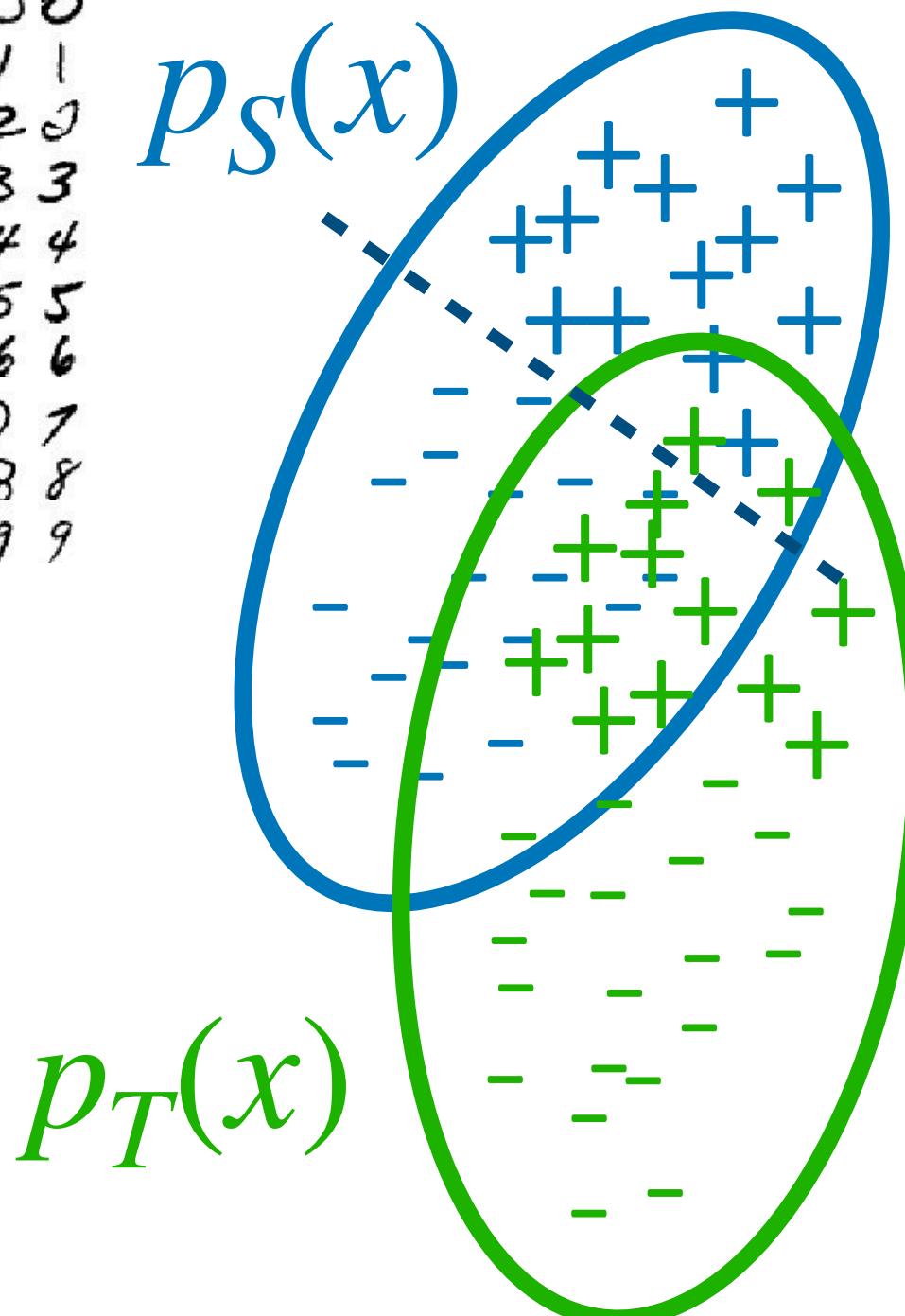
Can we align the features?

Source classifier in *aligned feature space*  
is more accurate in target domain.

How to align the features?

# Domain adaptation if support is not shared?

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



How to align the features?

Source encoder  $f_{\theta_S}$       Target encoder  $f_{\theta_T}$

Need to match features at *population-level*.

i.e. make encoded samples  $f_{\theta_S}(x), x \sim p_S(\cdot)$

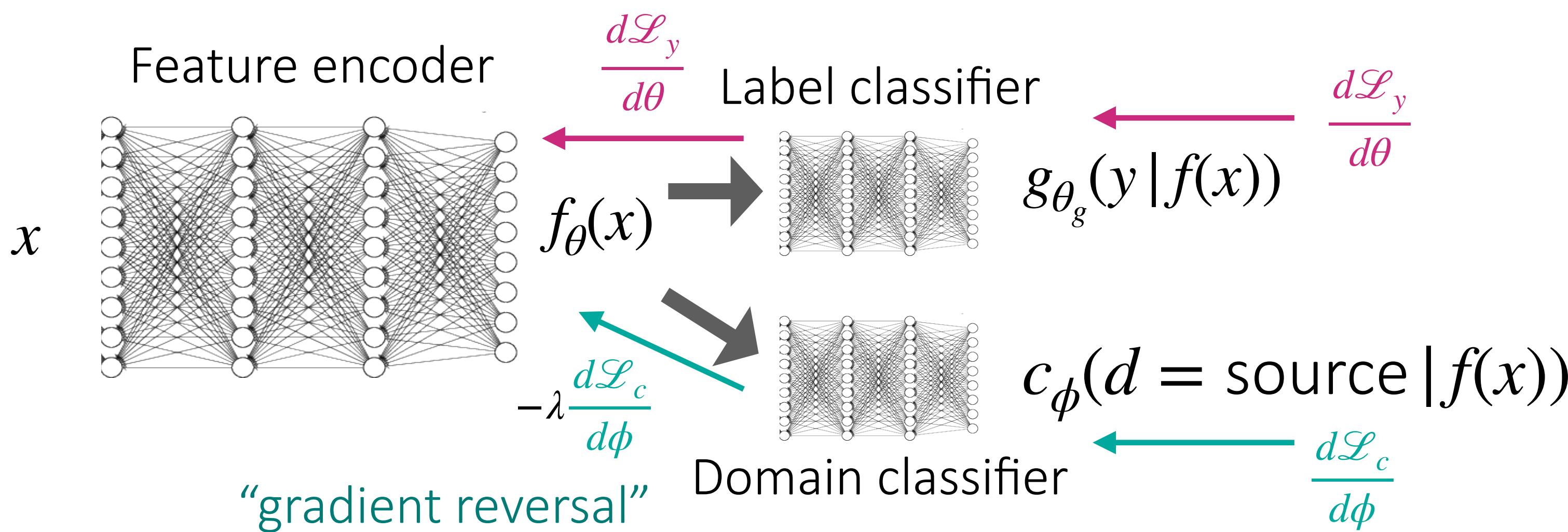
indistinguishable from  $f_{\theta_T}(x), x \sim p_T(\cdot)$

Key idea: Try to fool a domain classifier  $c(d = \text{source} | f(x))$ .

If samples are indistinguishable to discriminator, then distributions are the same.

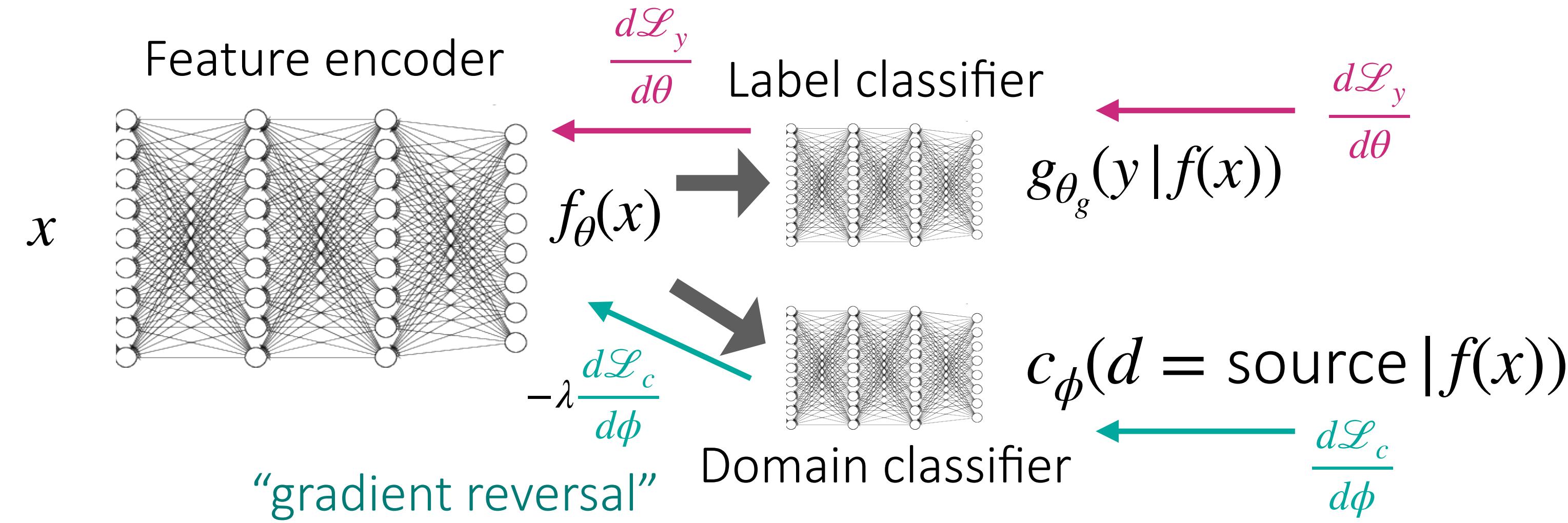
# Domain adaptation via feature alignment

Key idea: Try to fool a domain classifier  $c(d = \text{source} | f(x))$ .



Minimize label prediction error & maximize “domain confusion”

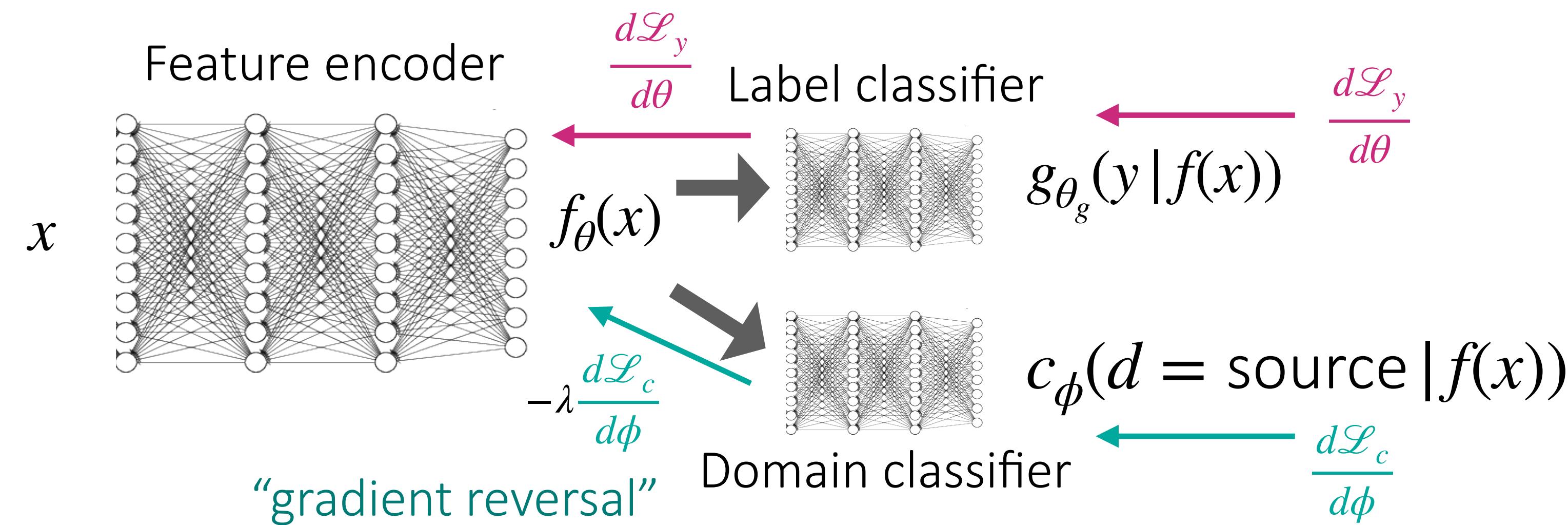
# Domain adaptation via feature alignment



Full algorithm:

1. Randomly initialize encoder(s)  $f_{\theta}$ , label classifier  $g_{\theta_g}$ , domain classifier  $c_{\phi}$
2. Update domain classifier:  $\min_{\phi} \mathcal{L}_c = -\mathbb{E}_{x \sim D_S}[\log c_{\phi}(f(x))] - \mathbb{E}_{x \sim D_T}[1 - \log c_{\phi}(f(x))]$ .
3. Update label classifier & encoder:  $\min_{\theta, \theta_g} \mathbb{E}_{(x,y) \sim D_S}[L(g_{\theta_g}(f_{\theta}(x)), y)] - \lambda \mathcal{L}_c$
4. Repeat steps 2 & 3.

# Domain adaptation via feature alignment



Can learn separate source and target encoder

Source encoder  $f_{\theta_S}$

Target encoder  $f_{\theta_T}$

Make encoded samples  $f_{\theta_S}(x), x \sim p_S(\cdot)$

indistinguishable from  $f_{\theta_T}(x), x \sim p_T(\cdot)$

—> can give model more flexibility

Different forms of domain adversarial training.

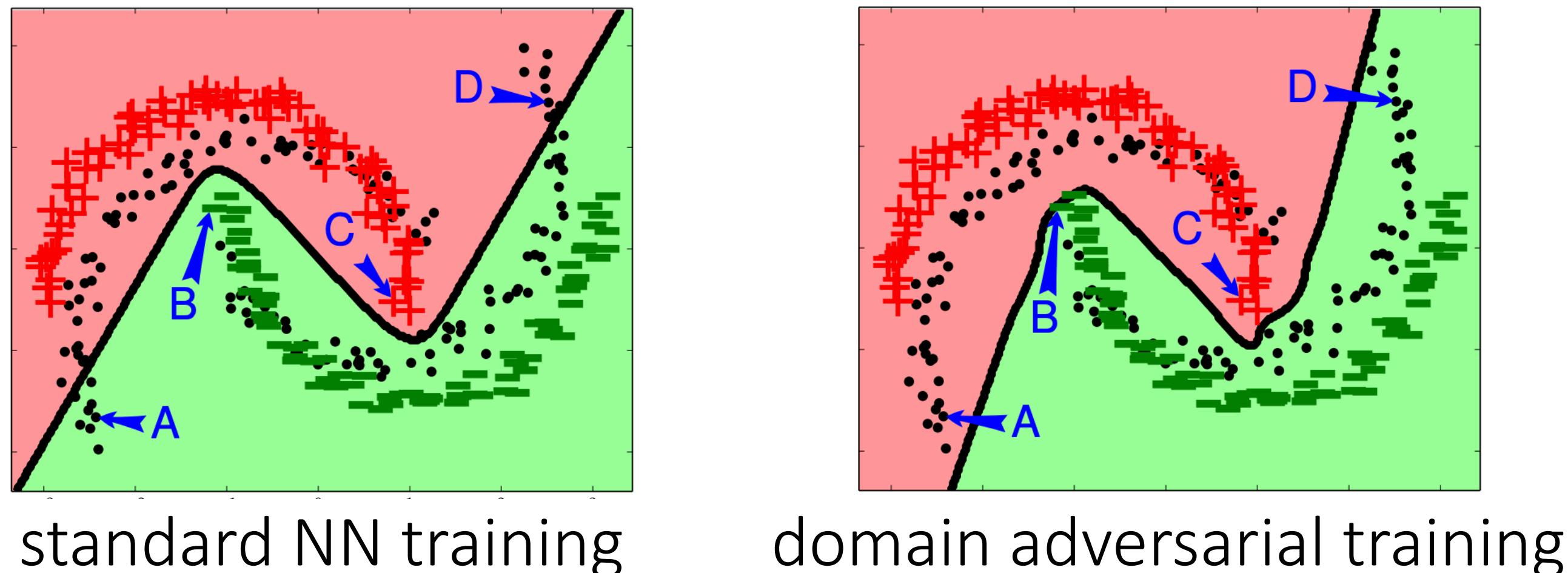
Option 1: Maximize domain classifier loss  
(gradient reversal, same as GANs)

Option 2: Optimize for 50/50 guessing

# Domain adaptation via feature alignment

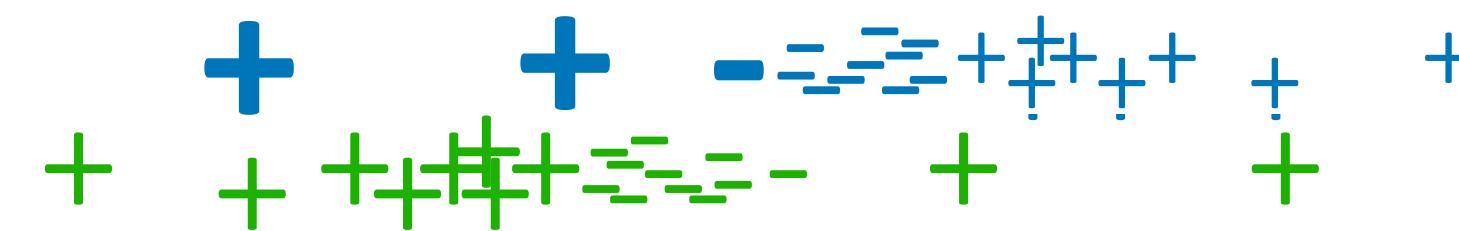
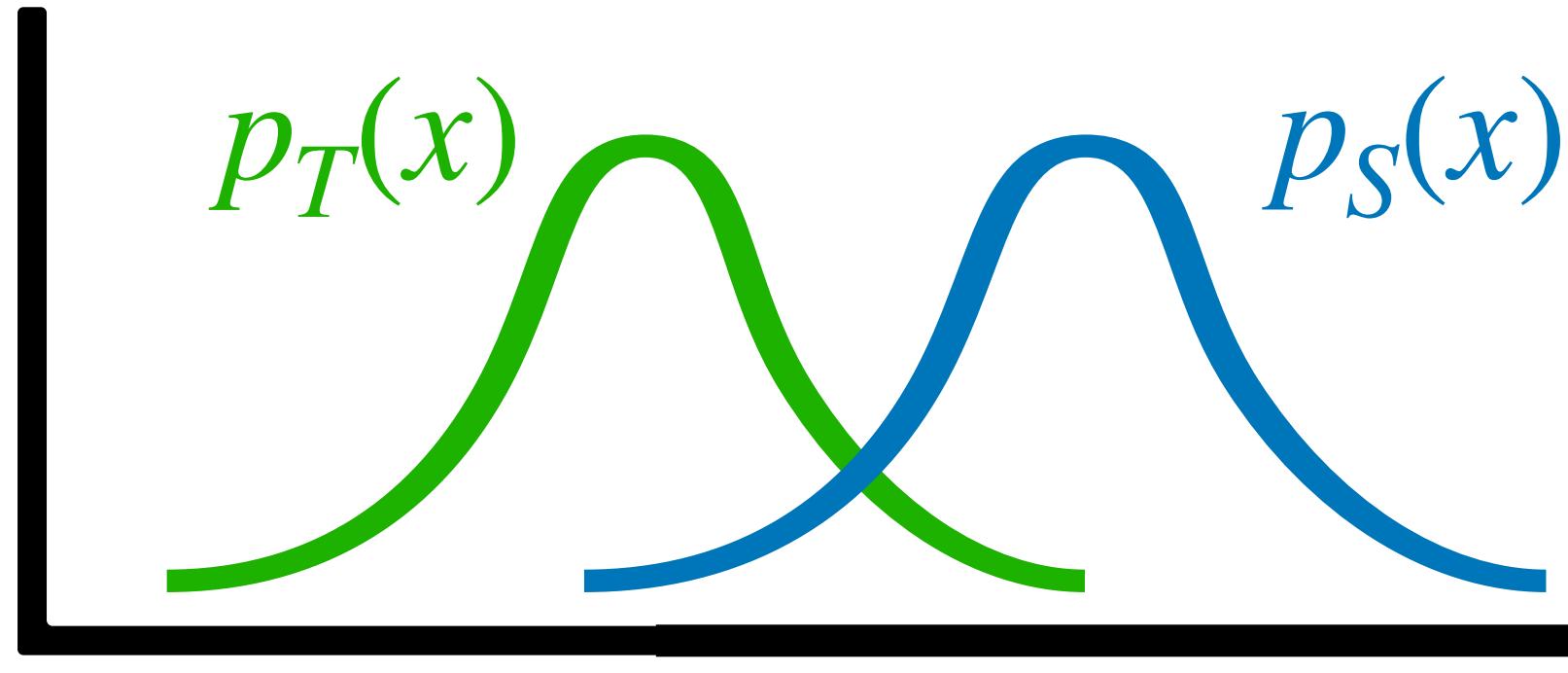
Toy example

source domain: +, —  
target domain data: •



SOURCE		MNIST	SYN NUMBERS	SVHN	SYN SIGNS
METHOD	SOURCE TARGET	MNIST-M	SVHN	MNIST	GTSRB
SOURCE ONLY		.5225	.8674	.5490	.7900
DANN		<b>.7666 (52.9%)</b>	<b>.9109 (79.7%)</b>	<b>.7385 (42.6%)</b>	<b>.8865 (46.4%)</b>
TRAIN ON TARGET		.9596	.9220	.9942	.9980

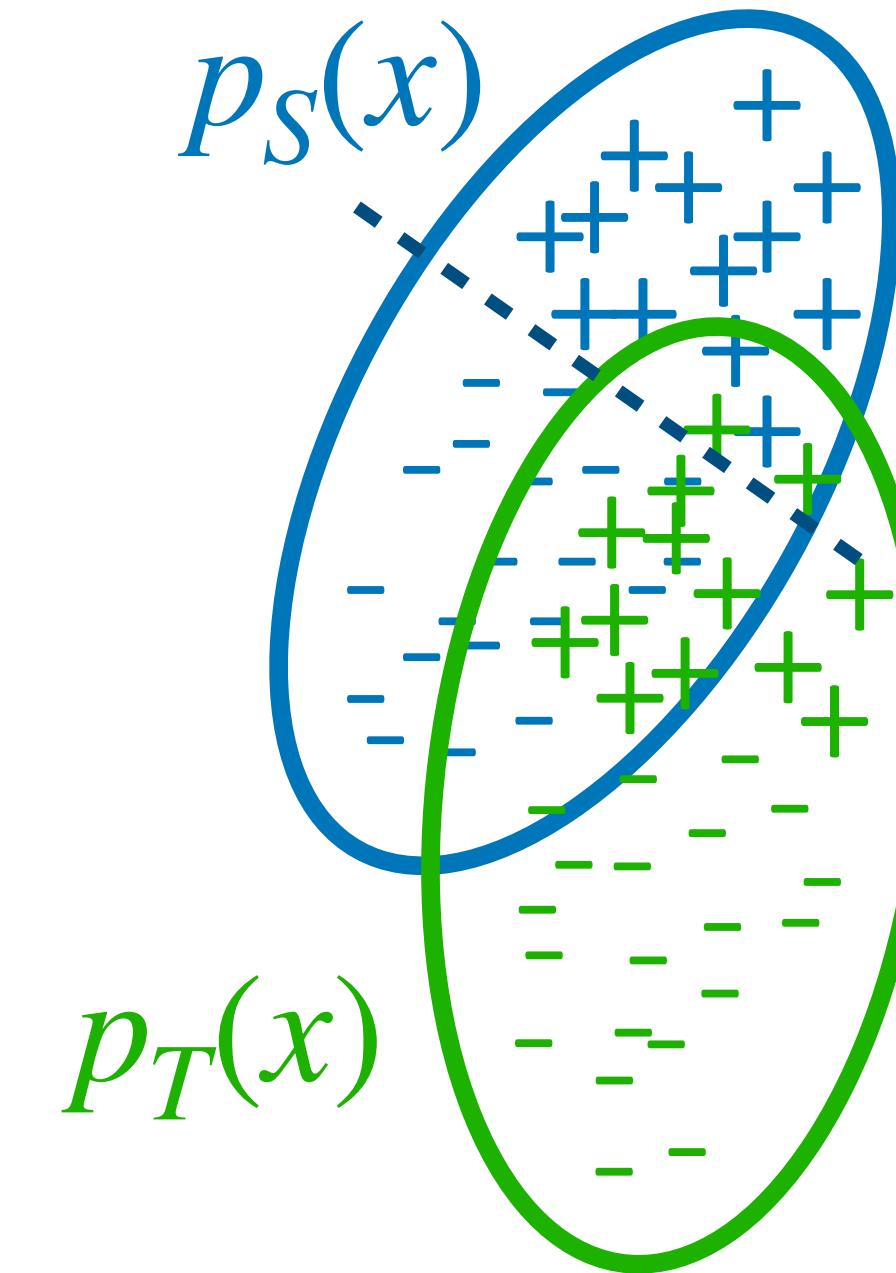
## Importance weighting



$$\min_{\theta} \mathbb{E}_{p_S(x,y)} \left[ \frac{p_T(x)}{p_S(x)} L(f_\theta(x), y) \right]$$

- + simple, can work well
- requires source distr. to cover target

## Feature alignment



- + fairly simple to implement, can work quite well
- + doesn't require source data coverage
- involves adversarial optimization
- requires clear alignment in data

# Plan for Today

## **Domain Adaptation**

- Problem statements
- Algorithms
  - Data reweighting
  - Feature alignment
  - **Domain translation**

**Goal for by the end of lecture:** Understand different domain adaptation methods and when to use one vs. another

# What if it is hard to align features?

Idea: translate between domains

i.e.  $F : X_S \rightarrow X_T$  or  $G : X_T \rightarrow X_S$

If you could translate source examples to target examples:

1. Translate labeled **source** dataset to **target** domain with  $F$ .
2. Train predictor on translated dataset.
3. Deploy predictor.

Alternatively, if you could translate from target to source:

1. Train predictor on **source** dataset.
2. Translate **target** example to **source** domain with  $G$ .
3. Evaluate predictor on translated example.

**Key question:** How to translate between domains?

# Domain Translation with CycleGAN

Idea: translate between domains

$$\text{i.e. } F : X_S \rightarrow X_T \quad \text{or } G : X_T \rightarrow X_S$$

Key question: How to translate between domains?

**Step 1:** Train  $F$  to generate images from  $p_T(x)$   
and  $G$  to generate images from  $p_S(x)$

Using GAN objective:  $\mathcal{L}_{\text{GAN}} = \mathbb{E}_{x \sim p_T(\cdot)}[\log D_T(x)] + \mathbb{E}_{x \sim p_S(\cdot)}[1 - \log D_T(F(x))]$

**Challenge:** The mapping is underconstrained, can be arbitrary.

Can we encourage models to learn a consistent, bijective mapping?

**Step 2:** Train  $F$  and  $G$  to be cyclically consistent.

$$F(G(x)) \approx x \text{ and } G(F(x)) \approx x$$

# Domain Translation with CycleGAN

Idea: translate between domains

$$\text{i.e. } F : X_S \rightarrow X_T \quad \text{or } G : X_T \rightarrow X_S$$

**Step 1:** Train  $F$  to generate images from  $p_T(x)$   
and  $G$  to generate images from  $p_S(x)$

Using GAN objective:  $\mathcal{L}_{\text{GAN}} = \mathbb{E}_{x \sim p_T(\cdot)}[\log D_T(x)] + \mathbb{E}_{x \sim p_S(\cdot)}[1 - \log D_T(F(x))]$

**Step 2:** Train  $F$  and  $G$  to be cyclically consistent.

$$F(G(x)) \approx x \text{ and } G(F(x)) \approx x$$

$$\text{i.e. } \mathbb{E}_{x \sim p_S(\cdot)} \|G(F(x)) - x\|_1 + \mathbb{E}_{x \sim p_T(\cdot)} \|F(G(x)) - x\|_1$$

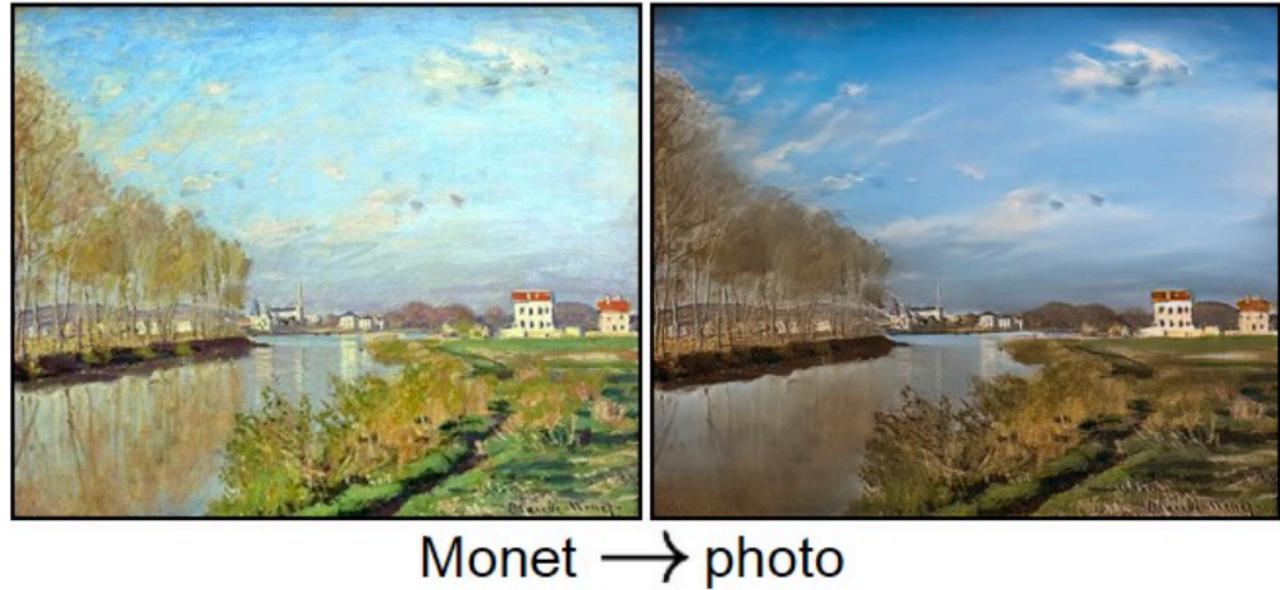
**Full objective:**  $\mathcal{L}_{\text{GAN}}(F, D_T) + \mathcal{L}_{\text{GAN}}(G, D_S) + \lambda \mathcal{L}_{\text{cyc}}(F, G)$

# Domain Translation with CycleGAN

Idea: translate between domains

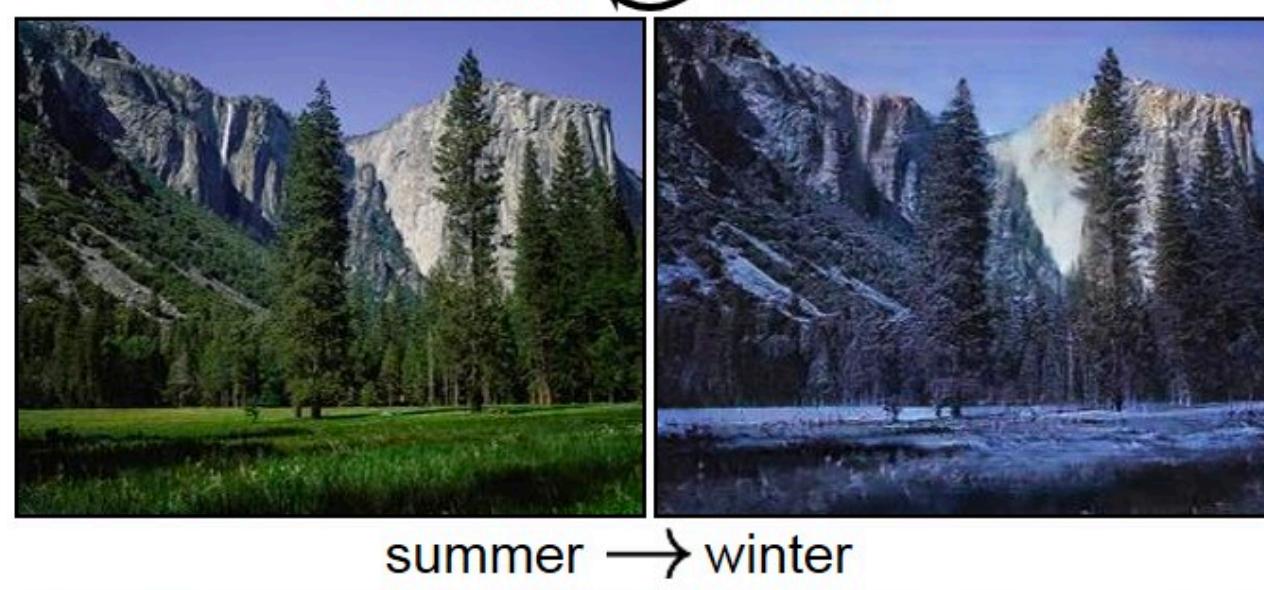
$$\text{i.e. } F : X_S \rightarrow X_T \quad \text{or } G : X_T \rightarrow X_S$$

Monet  $\curvearrowright$  Photos



Monet  $\rightarrow$  photo

Summer  $\curvearrowright$  Winter



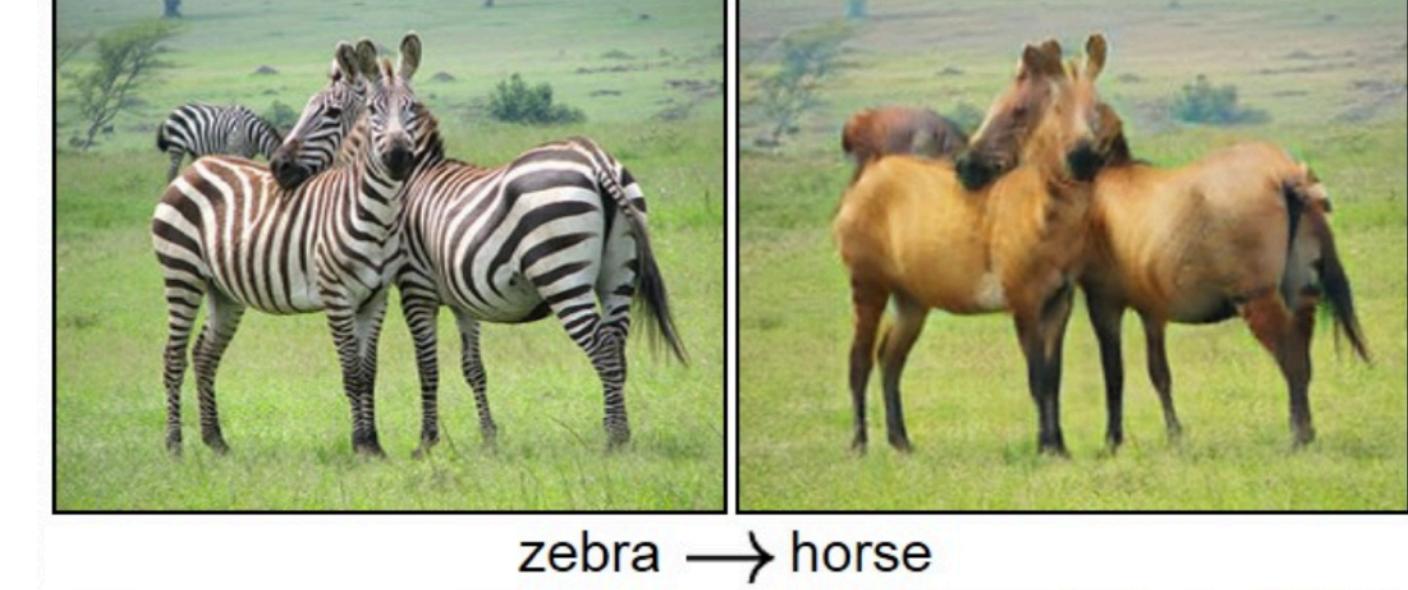
summer  $\rightarrow$  winter

edges  $\rightarrow$  shoes



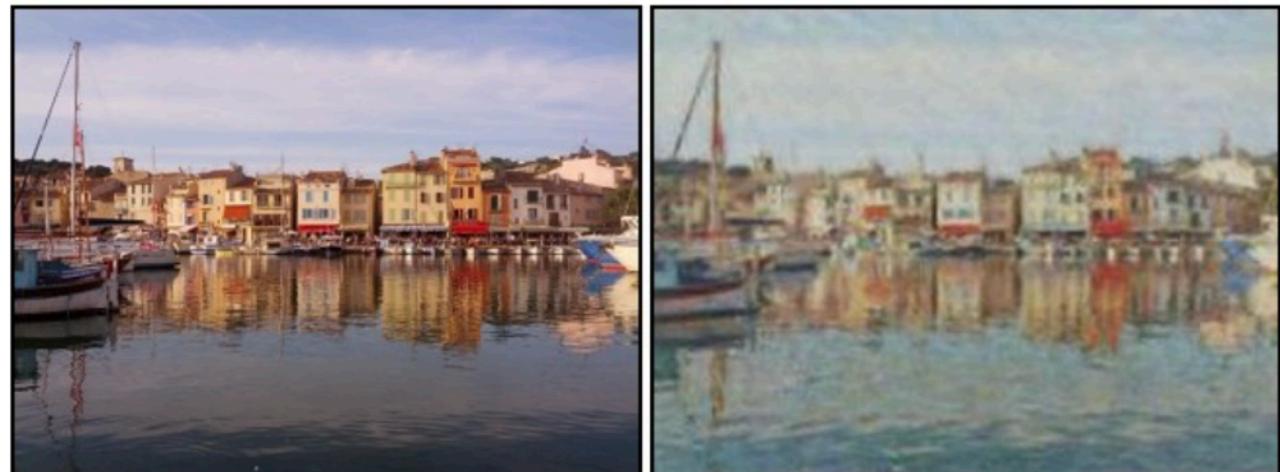
edges  $\rightarrow$  shoes

Zebras  $\curvearrowright$  Horses



zebra  $\rightarrow$  horse

photo  $\rightarrow$  Monet



winter  $\rightarrow$  summer

shoes  $\rightarrow$  edges



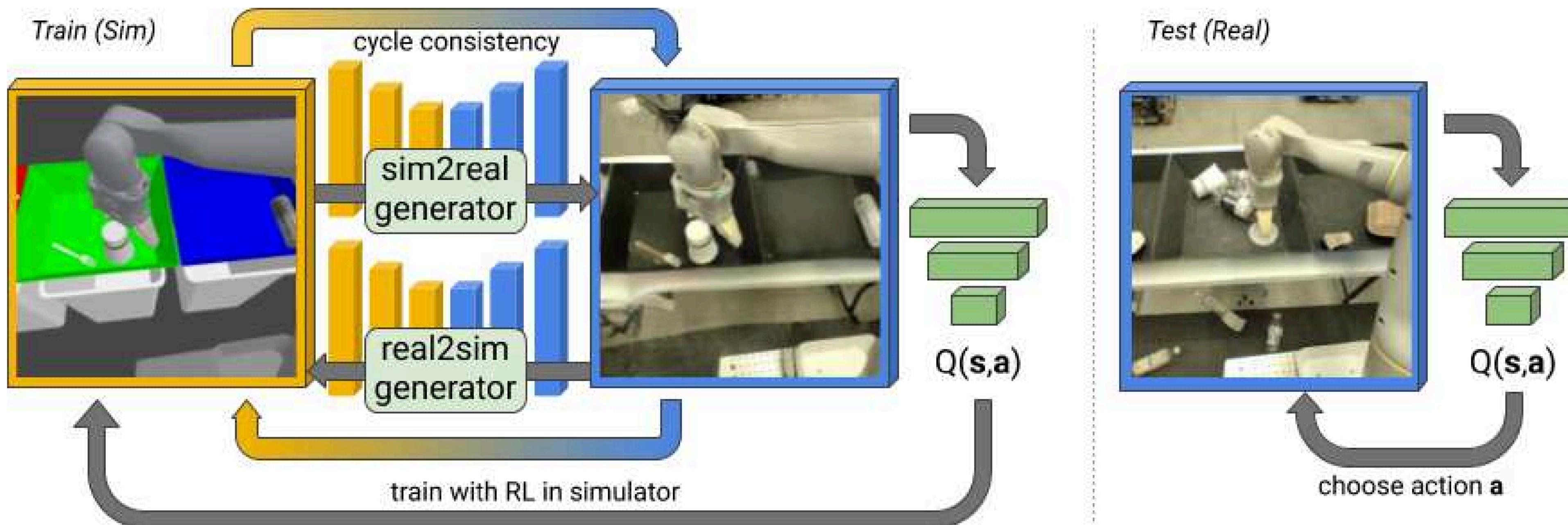
shoes  $\rightarrow$  edges

horse  $\rightarrow$  zebra



# CycleGAN for Domain Adaptation

Robotics sim2real policy adaptation

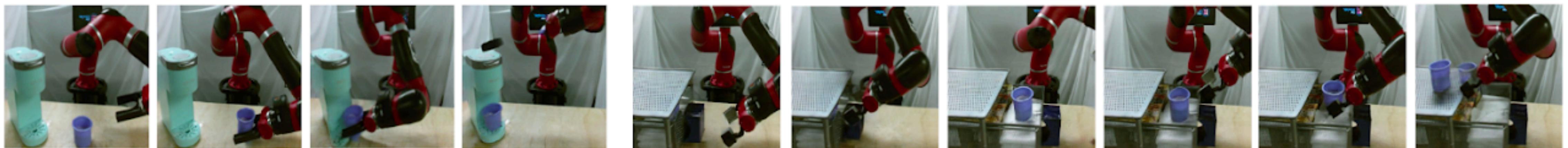


Simulation-to-Real Model	Robot 1 Grasp Success
Sim-Only [19]	21%
Randomized Sim [19]	37%
GAN	29%
CycleGAN	61%
GraspGAN	63%
<b>RL-CycleGAN</b>	<b>70%</b>

# CycleGAN for Domain Adaptation

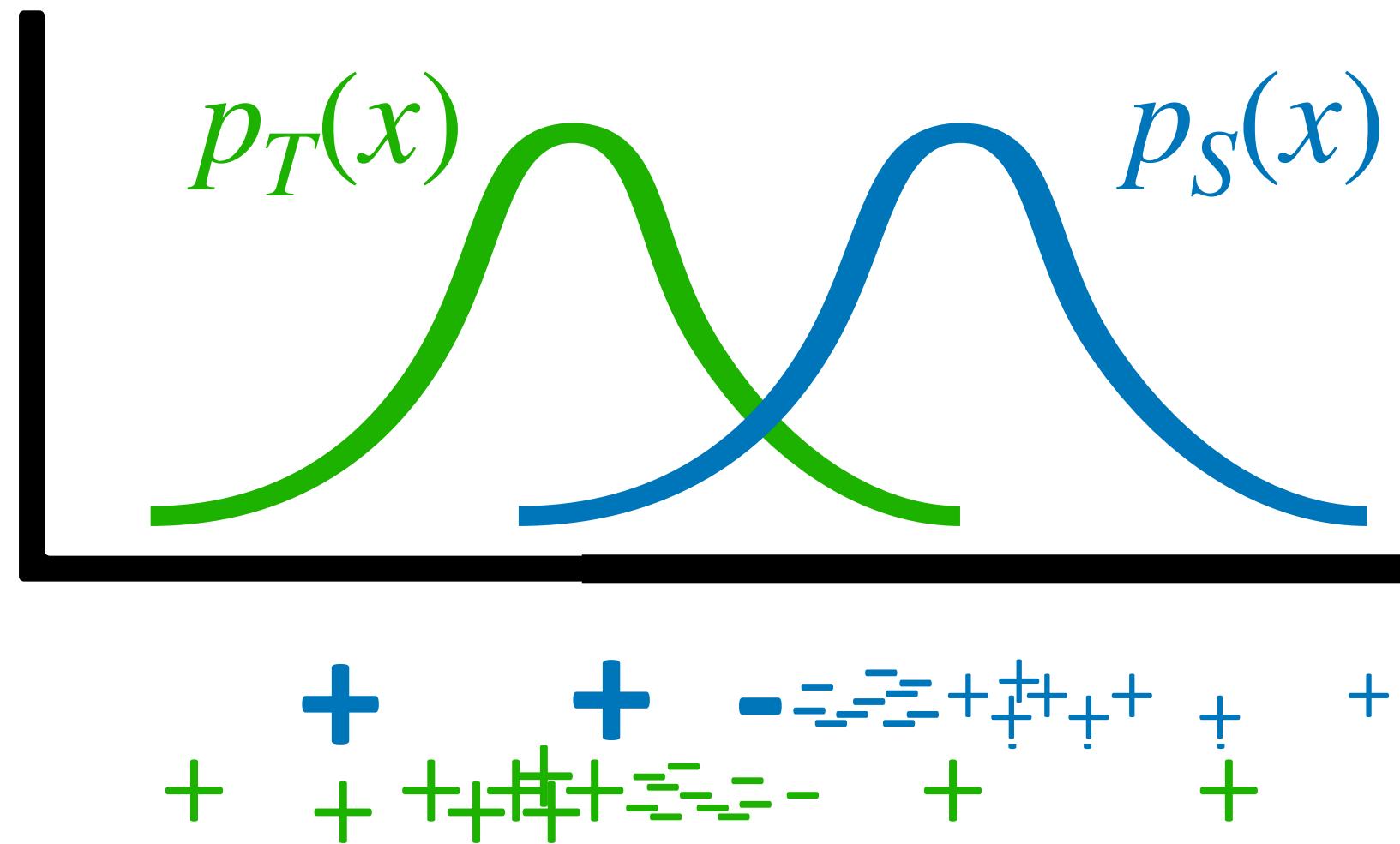
Human-robot domain adaptation

Input human images



Generated images in robot domain

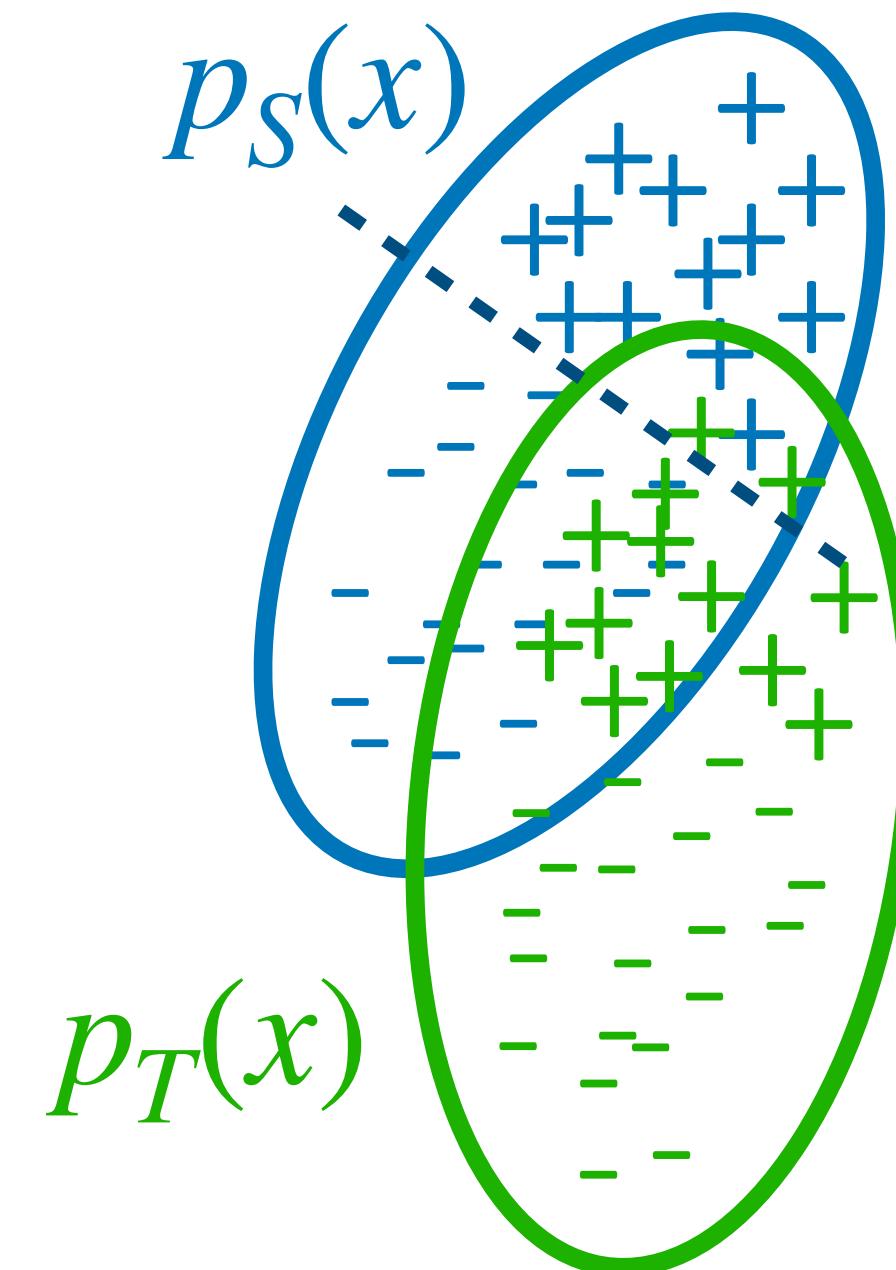
## Importance weighting



$$\min_{\theta} \mathbb{E}_{p_{S(x,y)}} \left[ \frac{p_T(x)}{p_S(x)} L(f_{\theta}(x), y) \right]$$

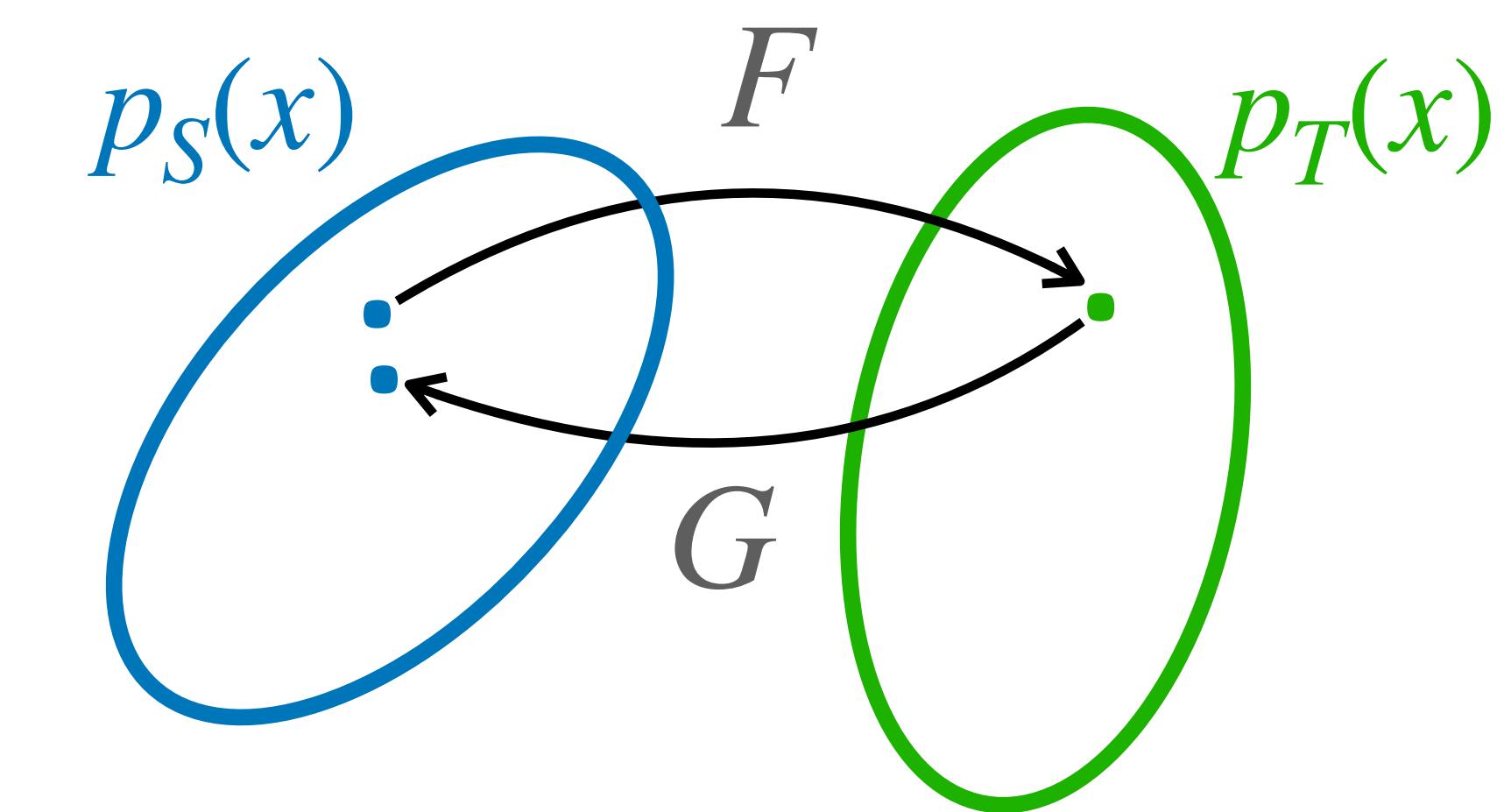
- + simple, can work well
- requires source distr. to cover target

## Feature alignment



- + fairly simple to implement, can work quite well
- + doesn't require source coverage
- involves adversarial optimization
- requires clear alignment in data

## Domain translation



- + conceptually neat, can work quite well
- + interpretable (easier to debug, cool pictures)
- involves generative modeling & adversarial optimization
- requires clear alignment in data

# CycleGAN & DANN for Domain Adaptation

CyCADA: incorporates both cycle consistency & domain adversarial training

## Character recognition

Model	USPS → MNIST	SVHN → MNIST
Source only	$69.6 \pm 3.8$	$67.1 \pm 0.6$
DANN (Ganin et al., 2016)	-	73.6
DTN (Taigman et al., 2017a)	-	84.4
CoGAN (Liu & Tuzel, 2016b)	89.1	-
ADDA (Tzeng et al., 2017)	$90.1 \pm 0.8$	$76.0 \pm 1.8$
PixelDA (Bousmalis et al., 2017b)	-	-
UNIT (Liu et al., 2017)	93.6	<b>90.5*</b>
<b>CyCADA (Ours)</b>	<b>96.5 ± 0.1</b>	<b>90.4 ± 0.4</b>
Target Fully Supervised	$99.2 \pm 0.1$	$99.2 \pm 0.1$

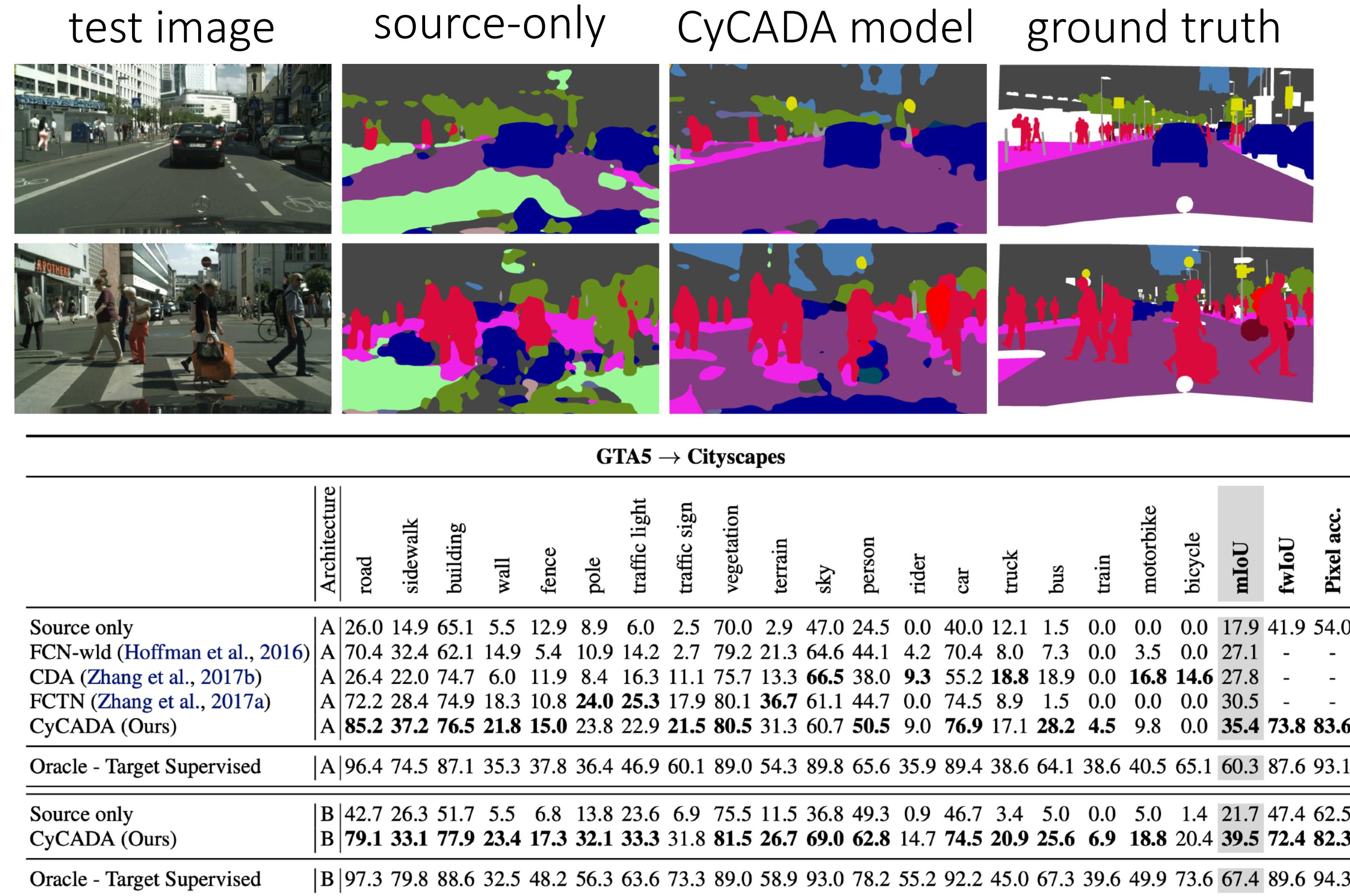


Table 4: Adaptation between GTA5 and Cityscapes, showing IoU for each class and mean IoU, freq-weighted IoU and pixel accuracy. CyCADA significantly outperforms baselines, nearly closing the gap to the target-trained oracle on pixel accuracy.

# Plan for Today

## **Domain Adaptation**

- Problem statements
- Algorithms
  - Data reweighting
  - Feature alignment
  - Domain translation

**Goal for by the end of lecture:** Understand different domain adaptation methods and when to use one vs. another

# Plan for Today

## Domain Generalization

- Problem formulation
- Algorithms
  - Adding explicit regularizers
  - Data augmentation

## Goals for this lecture:

- Understand [domain generalization](#): intuition, problem formulation
- Familiarize mainstream DG approaches: [regularization-based](#), [augmentation-based](#)

# Recap: Domain Adaptation

Perform well on target domain  $p_T(x, y)$ ,  
using training data from source domain(s)  $p_S(x, y)$

A form of **transfer learning**, with access to target domain data during training  
("transductive" learning)

Unsupervised domain adaptation: access to unlabeled target domain data

Common assumptions:

- Source and target domain only differ in domain of the function, i.e.  $p_S(y | x) = p_T(y | x)$
- There exists a single hypothesis with low error on both source and target domains.

Revisiting: A "domain" is a special case of a "task"

A task:  $\mathcal{T}_i \triangleq \{p_i(\mathbf{x}), p_i(\mathbf{y} | \mathbf{x}), \mathcal{L}_i\}$       A domain:  $d_i \triangleq \{p_i(\mathbf{x}), p(\mathbf{y} | \mathbf{x}), \mathcal{L}\}$

# Recap: Domain Adaptation

Perform well on target domain  $p_T(x, y)$ ,  
using training data from source domain(s)  $p_S(x, y)$

A form of **transfer learning**

Unsupervised domain adaptation

Can we always access unlabeled data  
from the target domain?

Common assumptions:

- Source and target domain only differ in domain of the function, i.e.  $p_S(y | x) = p_T(y | x)$
- There exists a single hypothesis with low error on both source and target domains.

Revisiting: A “domain” is a special case of a “task”

A task:  $\mathcal{T}_i \triangleq \{p_i(\mathbf{x}), p_i(\mathbf{y} | \mathbf{x}), \mathcal{L}_i\}$       A domain:  $d_i \triangleq \{p_i(\mathbf{x}), p(\mathbf{y} | \mathbf{x}), \mathcal{L}\}$

# Recap: Domain Adaptation

Perform well on target domain  $p_T(x, y)$ ,  
using training data from source domain(s)  $p_S(x, y)$

A form of ~~transfer learning~~ with access to target domain data during training

- Real-time deployment and don't have time to collect target domain data
- Obtaining target data may be restricted by privacy policy

Common

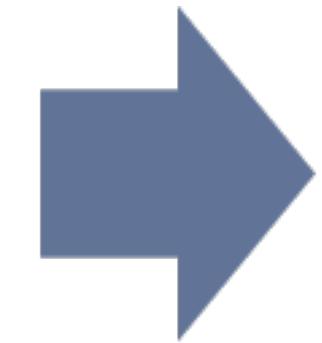
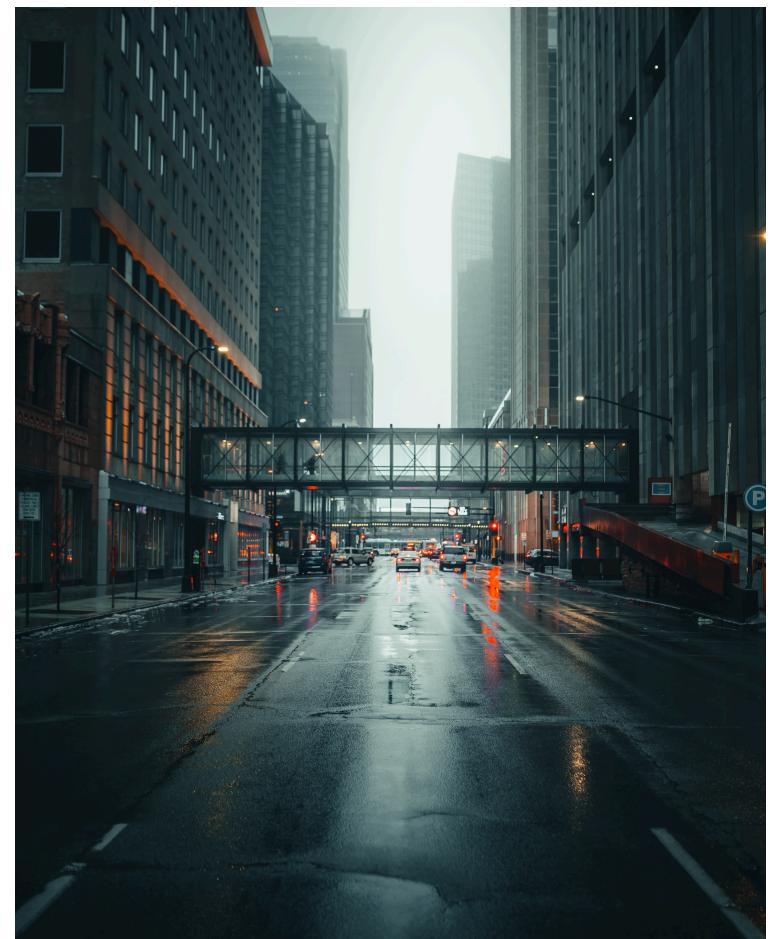
- Source and target domain only differ in domain of the function, i.e.  $p_S(y | x) = p_T(y | x)$
- There exists a single hypothesis with low error on both source and target domains.

Revisiting: A "domain" is a special case of a "task"

$$\text{A task: } \mathcal{T}_i \triangleq \{p_i(\mathbf{x}), p_i(\mathbf{y} | \mathbf{x}), \mathcal{L}_i\} \quad \text{A domain: } d_i \triangleq \{p_i(\mathbf{x}), p(\mathbf{y} | \mathbf{x}), \mathcal{L}\}$$

# Real-Time Deployment

Real-time deployment and don't have time to collect data



Deploy to a new road

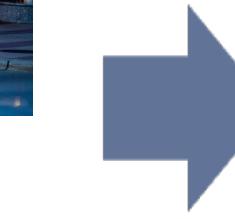
Trained on three types of roads

# Privacy Concerns

## GDPR



Trained on 3 hospitals



Deploy to a new hospital



Can't access training data

# Plan for Today

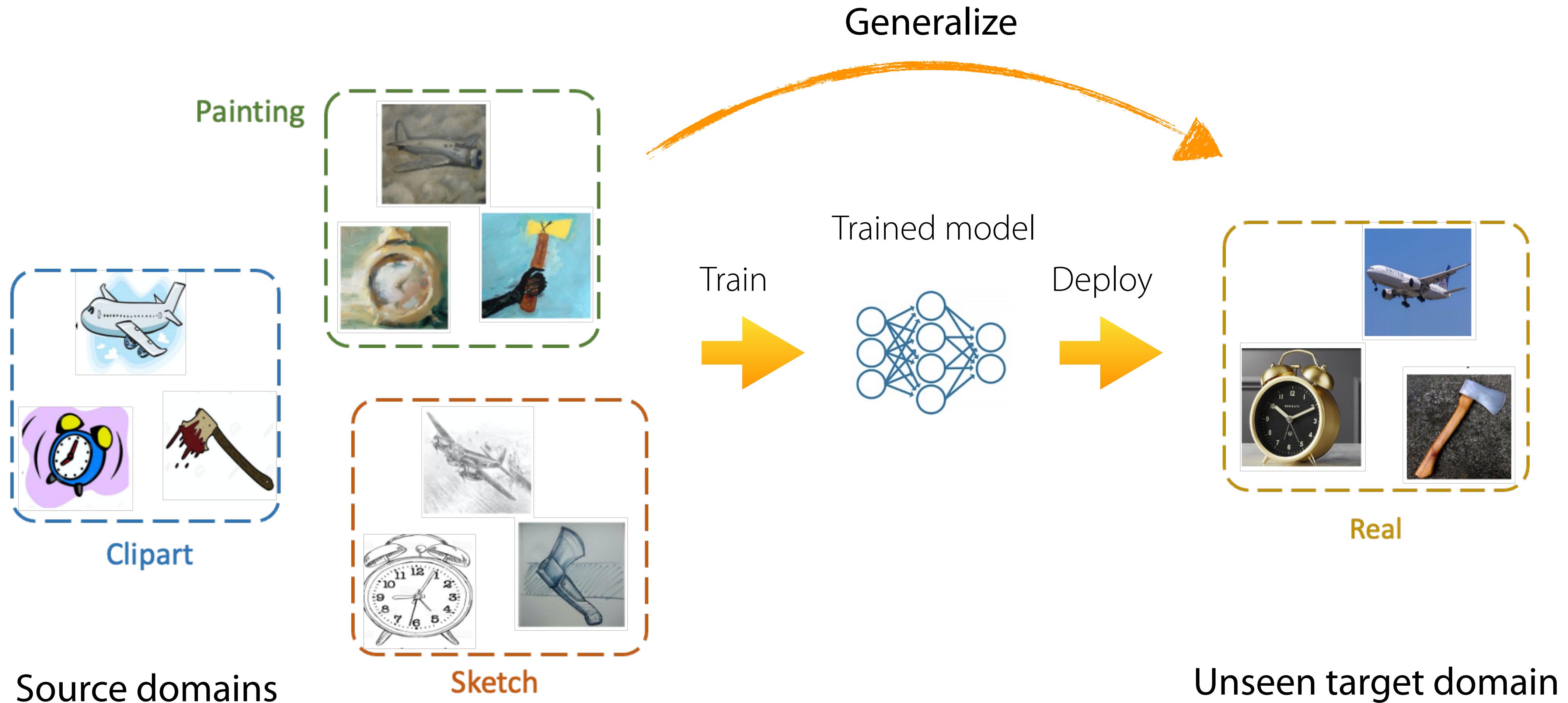
## Domain Generalization

- **Problem formulation**
- Algorithms
  - Adding explicit regularizers
  - Data augmentation

## Goals for this lecture:

- Understand [domain generalization](#): intuition, problem formulation
- Familiarize mainstream DG approaches: [regularization-based](#), [augmentation-based](#)

# Domain Generalization



Source: DomainNet dataset

# Domain Generalization Problem

Given source domains  $p_1(x, y), \dots, p_n(x, y)$ , solve unseen target domain  $p_T(x, y)$  without accessing the data from it.

Common assumptions

- All domains only differ in domain of the function, i.e.,  $p_1(y|x) = \dots = p_n(y|x) = p_T(y|x)$ .
- Only  $p(x)$  can change
- There exists a single hypothesis with low error in all domains.

Revisiting: A “domain” is a special case of a “task”

$$\text{A task: } \mathcal{T}_i \triangleq \{p_i(\mathbf{x}), p_i(\mathbf{y}|\mathbf{x}), \mathcal{L}_i\} \quad \text{A domain: } d_i \triangleq \{p_i(\mathbf{x}), p(\mathbf{y}|\mathbf{x}), \mathcal{L}\}$$

# Meta-Learning v.s. Domain Generalization

Revisiting: A “domain” is a special case of a “task”

$$\text{A task: } \mathcal{T}_i \triangleq \{p_i(\mathbf{x}), p_i(\mathbf{y} \mid \mathbf{x}), \mathcal{L}_i\} \quad \text{A domain: } d_i \triangleq \{p_i(\mathbf{x}), p(\mathbf{y} \mid \mathbf{x}), \mathcal{L}\}$$

## Meta-Learning Problem

Transfer learning with many source tasks

Given data from  $\mathcal{T}_1, \dots, \mathcal{T}_n$ , solve new task  $\mathcal{T}_t$  more quickly / proficiently / stably

## Domain Generalization

A special case of meta-learning

Given data from domains  $d_1, \dots, d_n$ , perform well on new domain  $d_t$

- Only  $p_i(x)$  changes across tasks
- direct generalization/no adaptation

# Domain Adaptation v.s. Domain Generalization

## Domain Adaptation

“transductive” setting

Given labeled data from source domain  $p_S(x, y)$  and unlabeled data from target domain  $p_T(x, y)$ , perform well on this target domain

Target data access during training  (unlabeled data)

Only one source domain

The model is specialized for the target domain

## Domain Generalization

“inductive” setting

Given labeled data from a set of source domains  $p_1(x, y), \dots, p_n(x, y)$ , perform well on target domain  $p_T(x, y)$

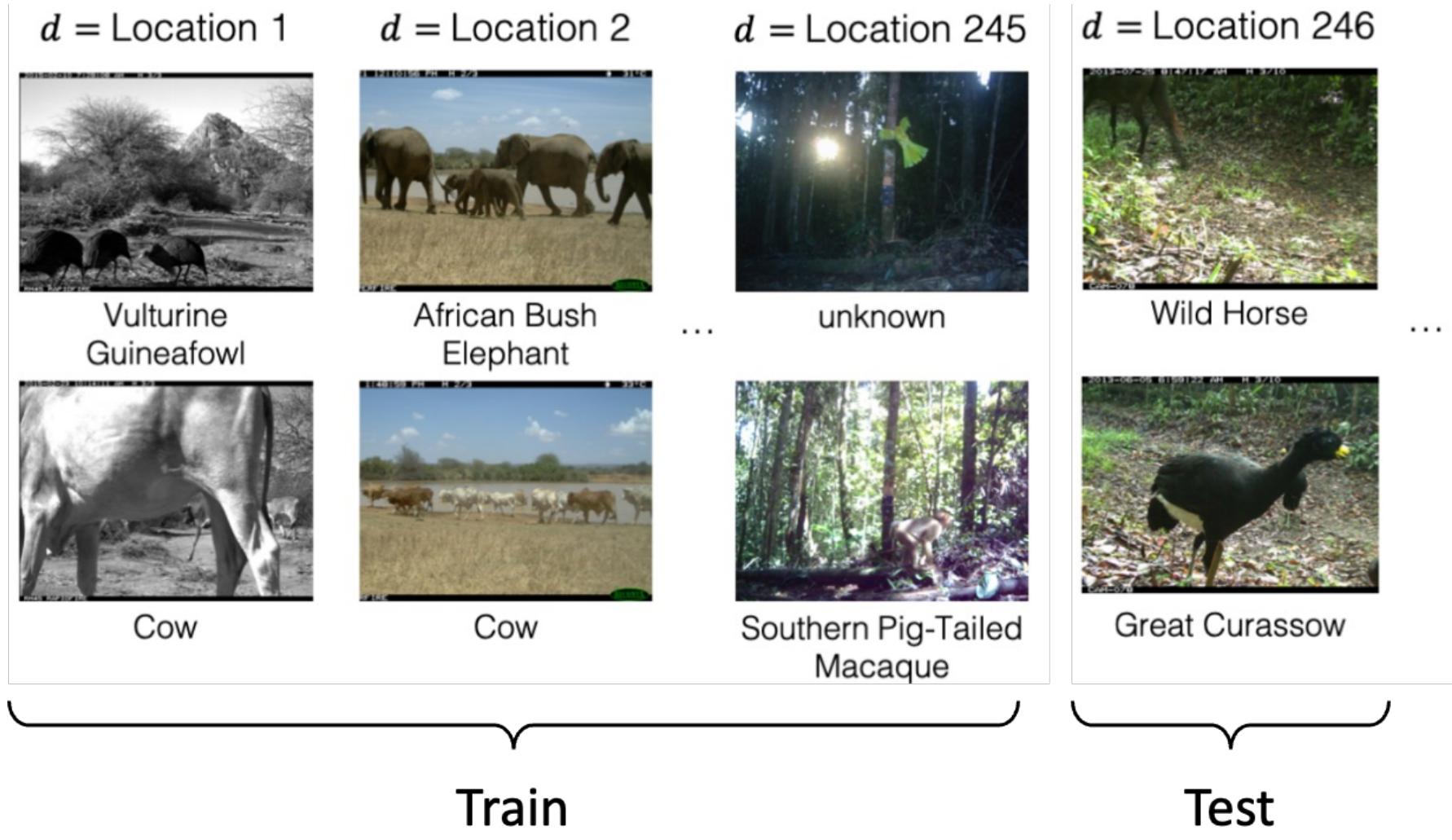
Test data access during training 

Need more than one source domain

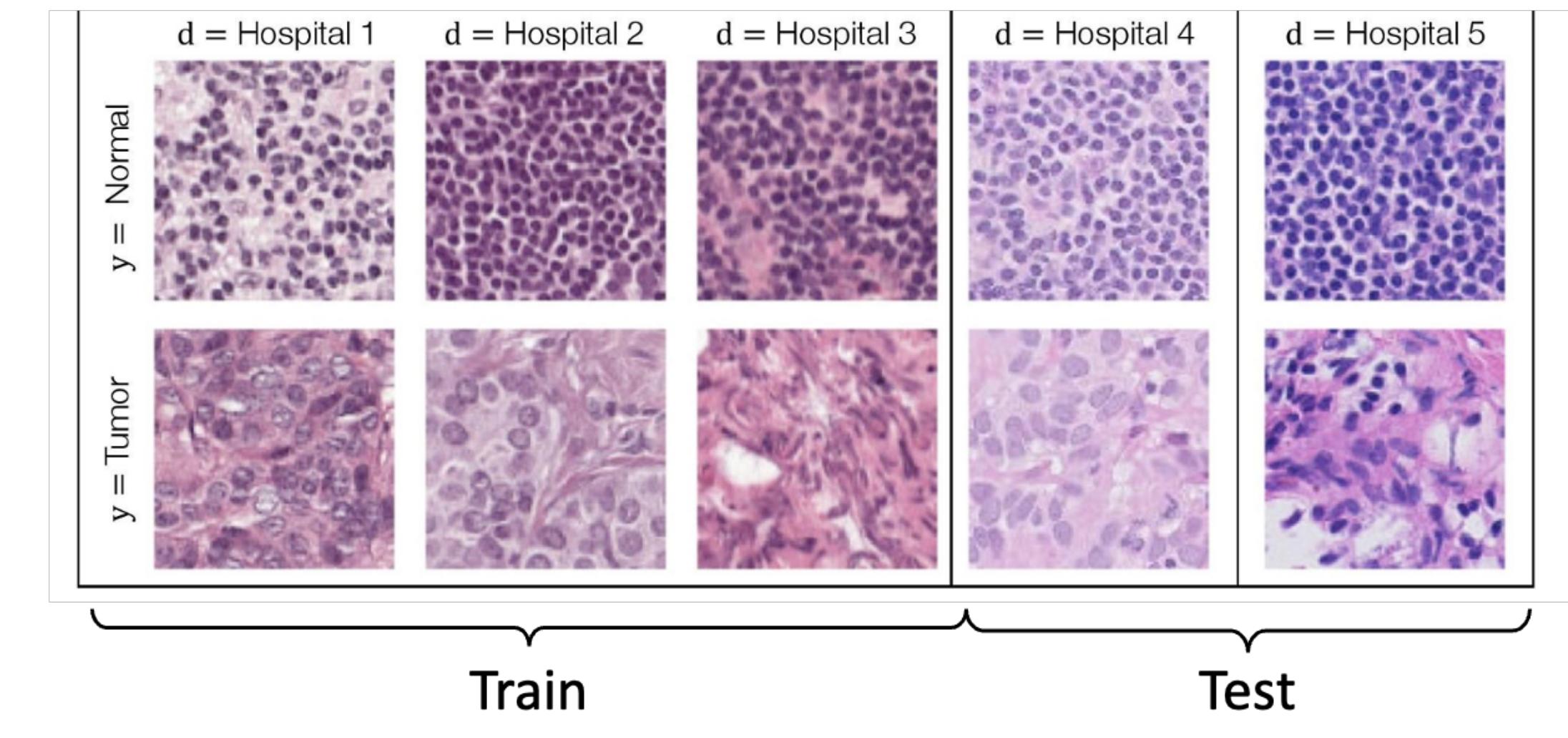
The model can be applied to all domains

# Domain Generalization: Applications

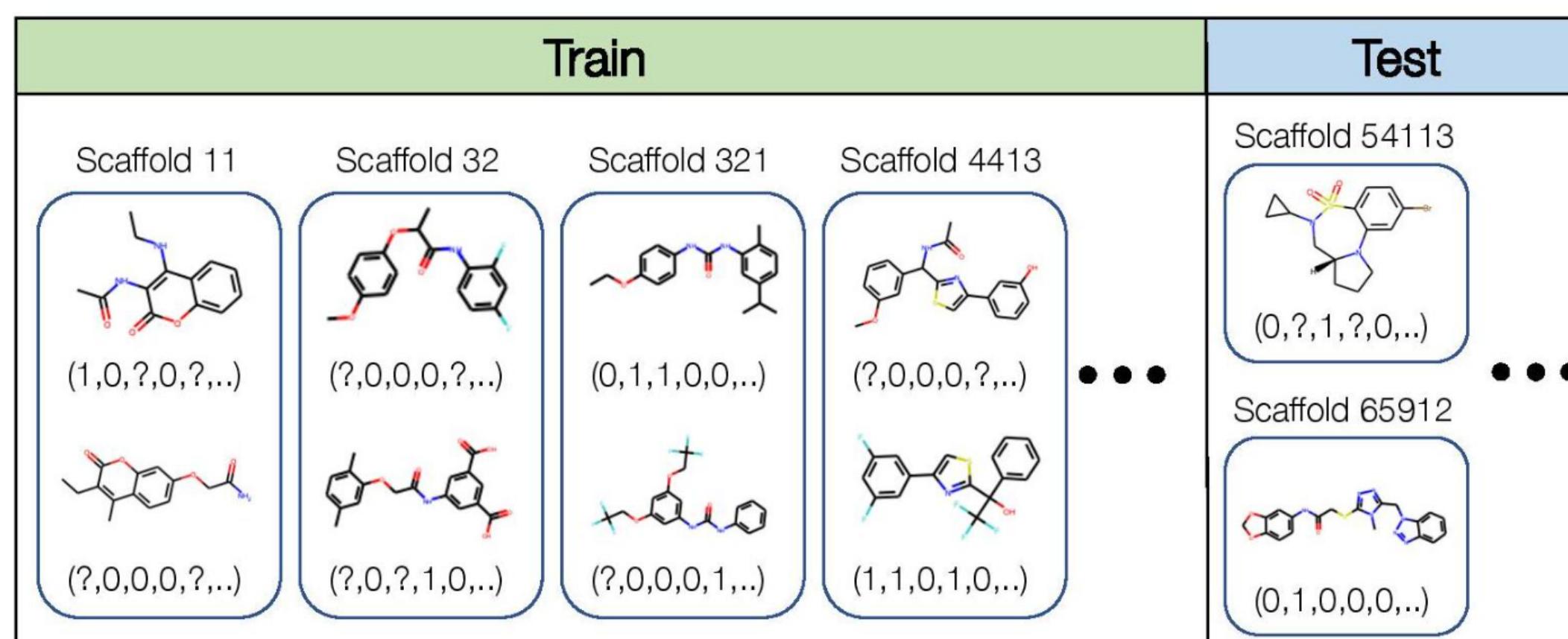
## Wildlife recognition



## Tissue classification



## Molecule property prediction



## Code completion

	Repository ID ( $d$ )	Source code context ( $x$ )	Next tokens ( $y$ )
Train	Repository 1	... from easyrec.gateway import EasyRec <EOL> gateway = EasyRec('tenant','key') <EOL> item_type = gateway. <span style="background-color: pink;">_____</span>	get_item_type
	Repository 2	... response = gateway.get_other_users() <EOL> get_params = HTTPPretty. <span style="background-color: pink;">_____</span>	last_request
Test	Repository 6,001	import numpy as np ... <EOL> if np.linalg.norm(target - prev_target) > far_threshold: <EOL> norm = np. <span style="background-color: pink;">_____</span>	linalg
		... new_trans = np.zeros((n_beats + max_beats, n_beats) <EOL> new_trans[:n_beats,:n_beats] = np. <span style="background-color: pink;">_____</span>	max
⋮			

# Plan for Today

## Domain Generalization

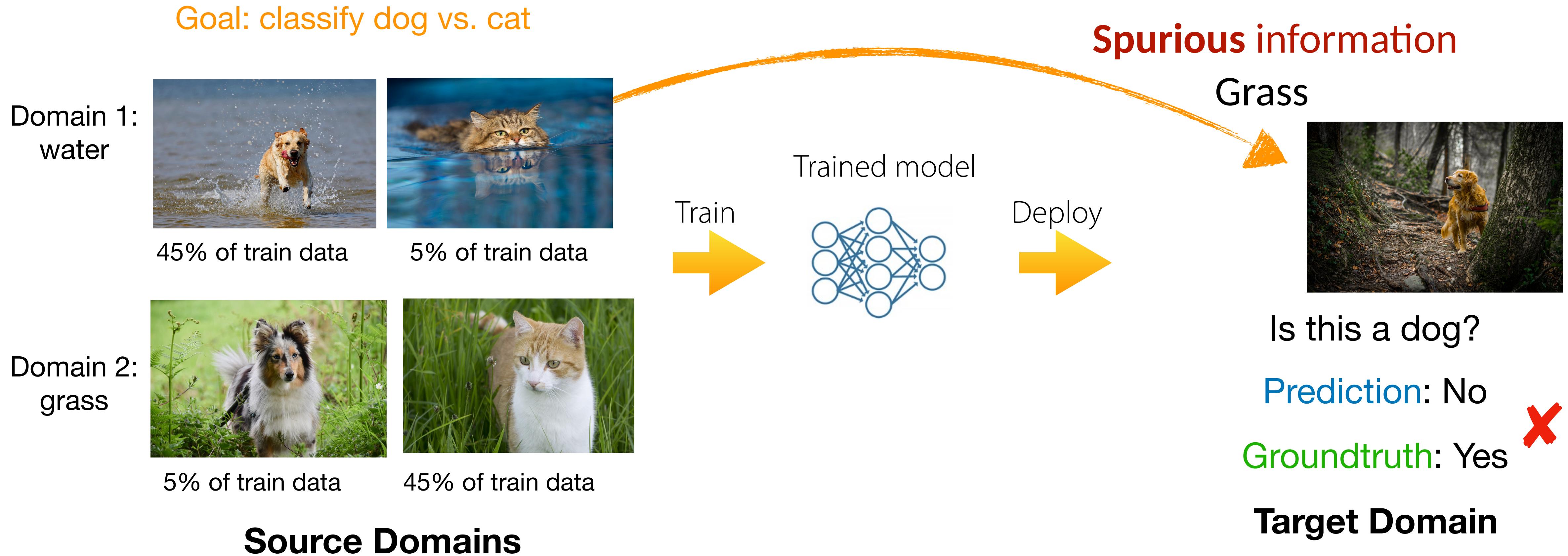
- Problem formulation
- Algorithms
  - **Adding explicit regularizers**
  - Data augmentation

## Goals for this lecture:

- Understand [domain generalization](#): intuition, problem formulation
- Familiarize mainstream DG approaches: [regularization-based](#), [augmentation-based](#)

# How to Learn Generalizable Representations?

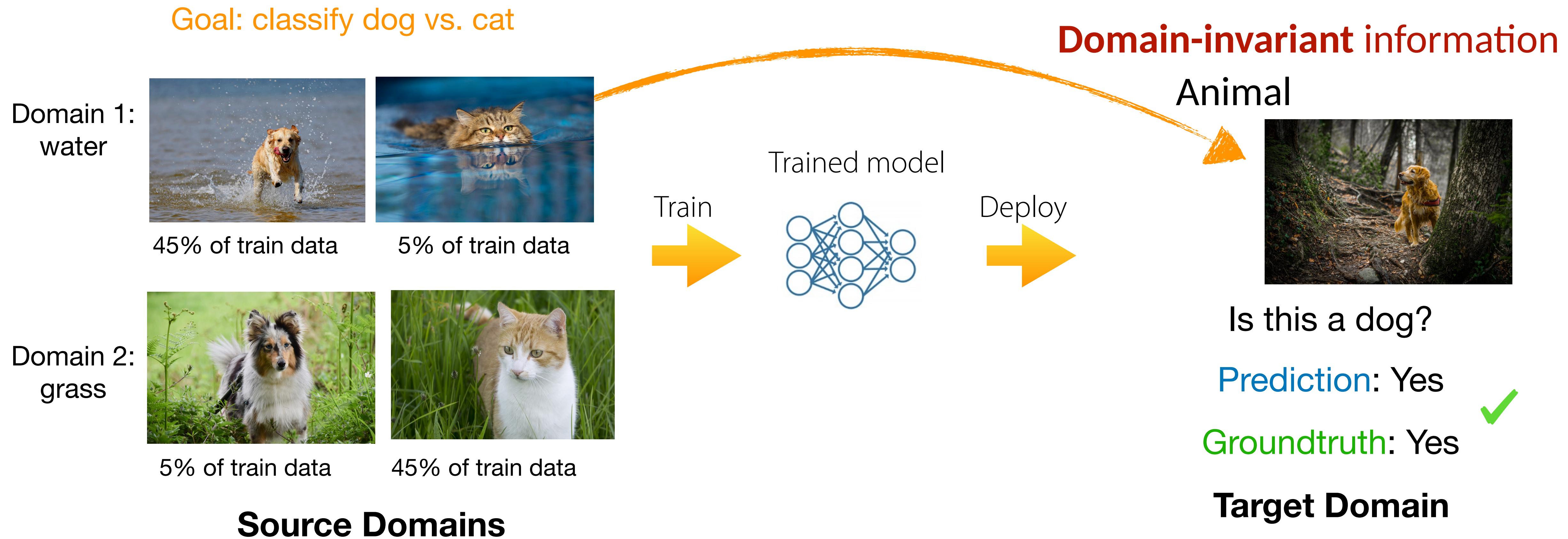
Why do machine learning models fail to generalize?



# How to Learn Generalizable Representations?

To overcome spurious correlation —> train a neural network to learn **domain invariance**

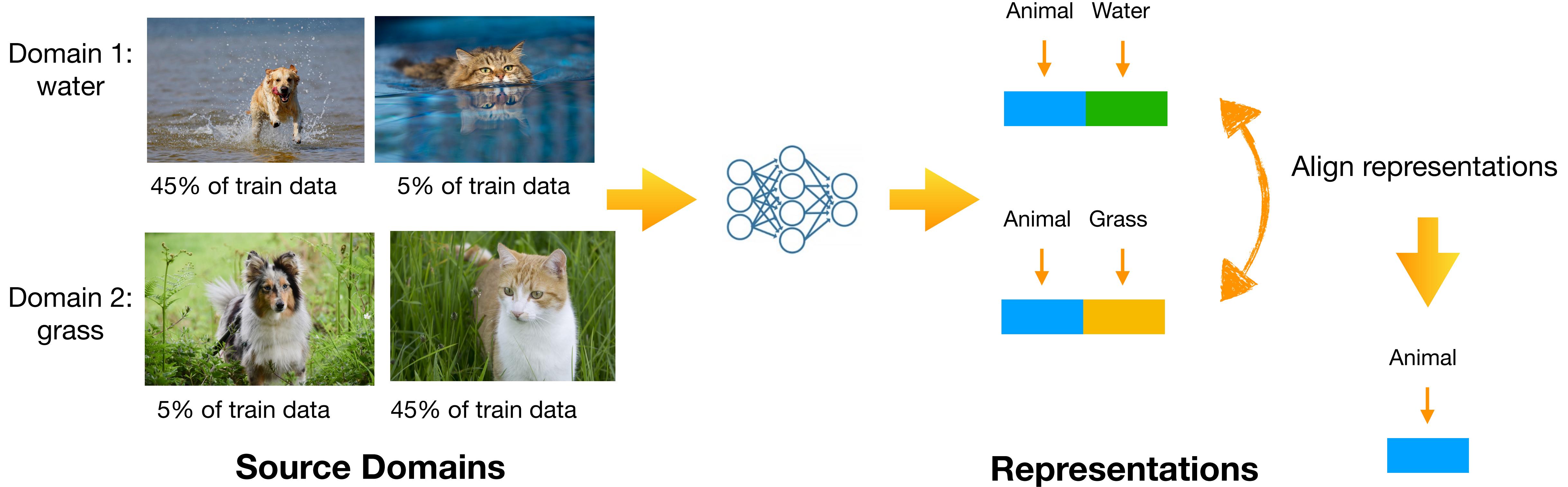
Domain invariance: we want to learn **features that don't change across domains**



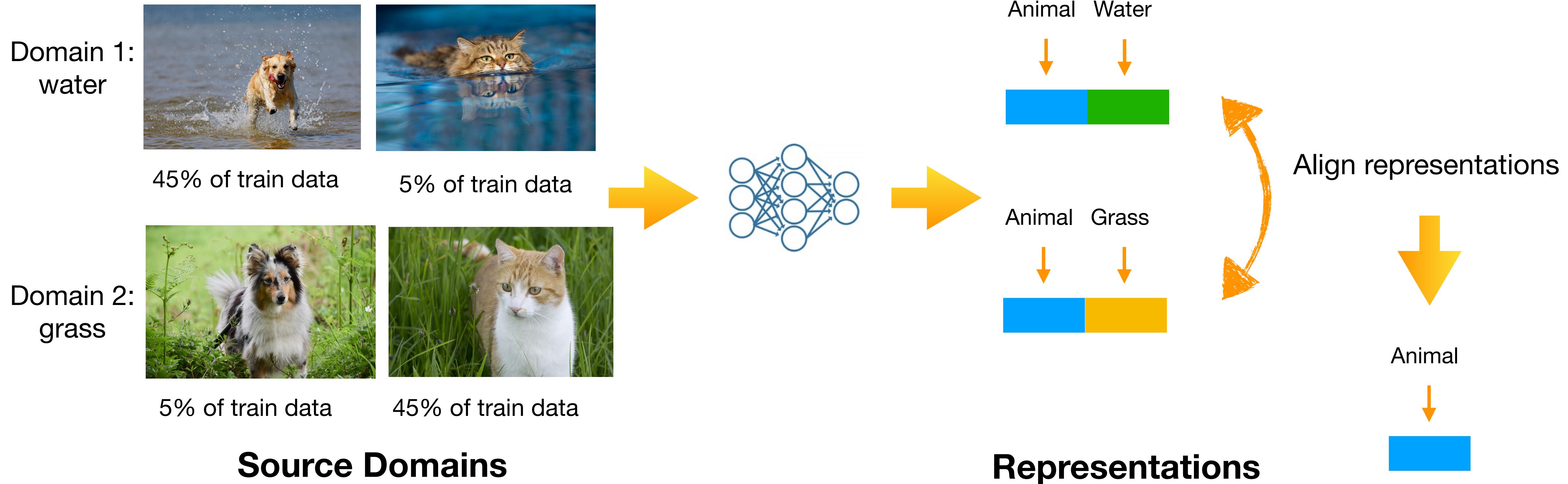
# Regularization-based Method

**Key idea:** Use a regularizer to align representations across domains

—> get domain-invariant representation



# Regularization-based Method



## Source Domains

## Representations

Label classification loss

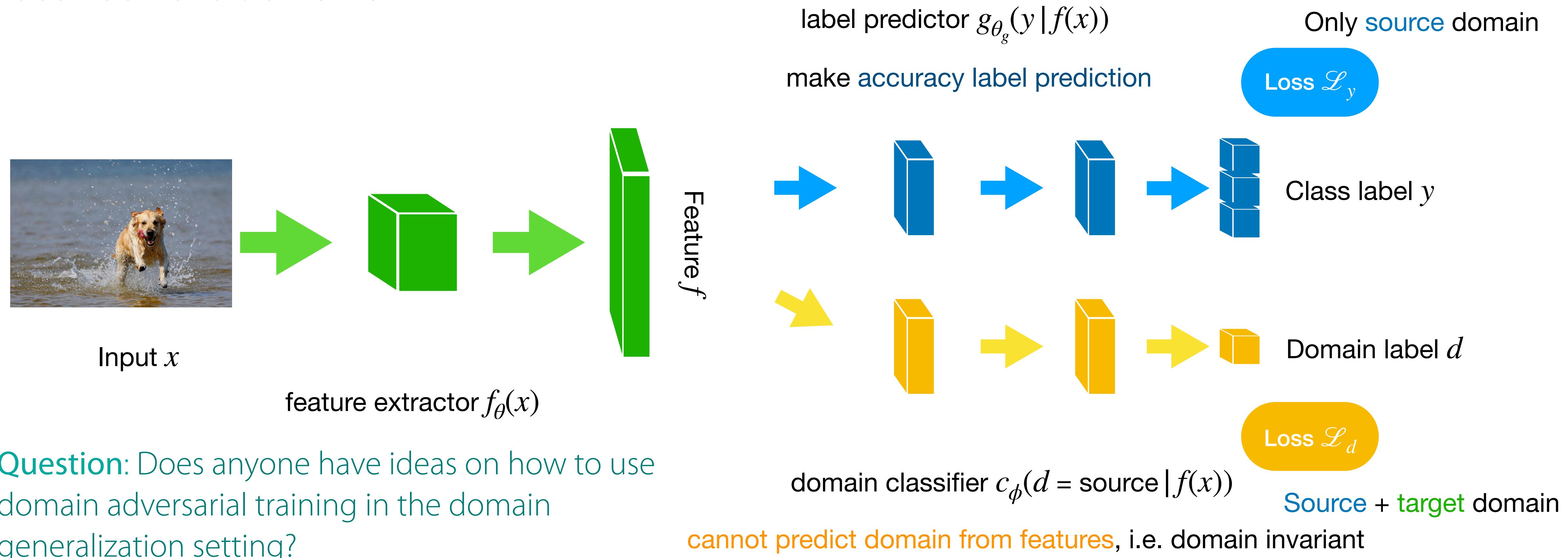
$$\min_{\theta} \mathbb{E}_{(x,y)}[\ell(f_{\theta}(x), y)] + \lambda \mathcal{L}_{reg}$$

Average over training examples

Explicit regularizer to learn  
domain-invariant representation

# Recap: Domain Adversarial Training in DA

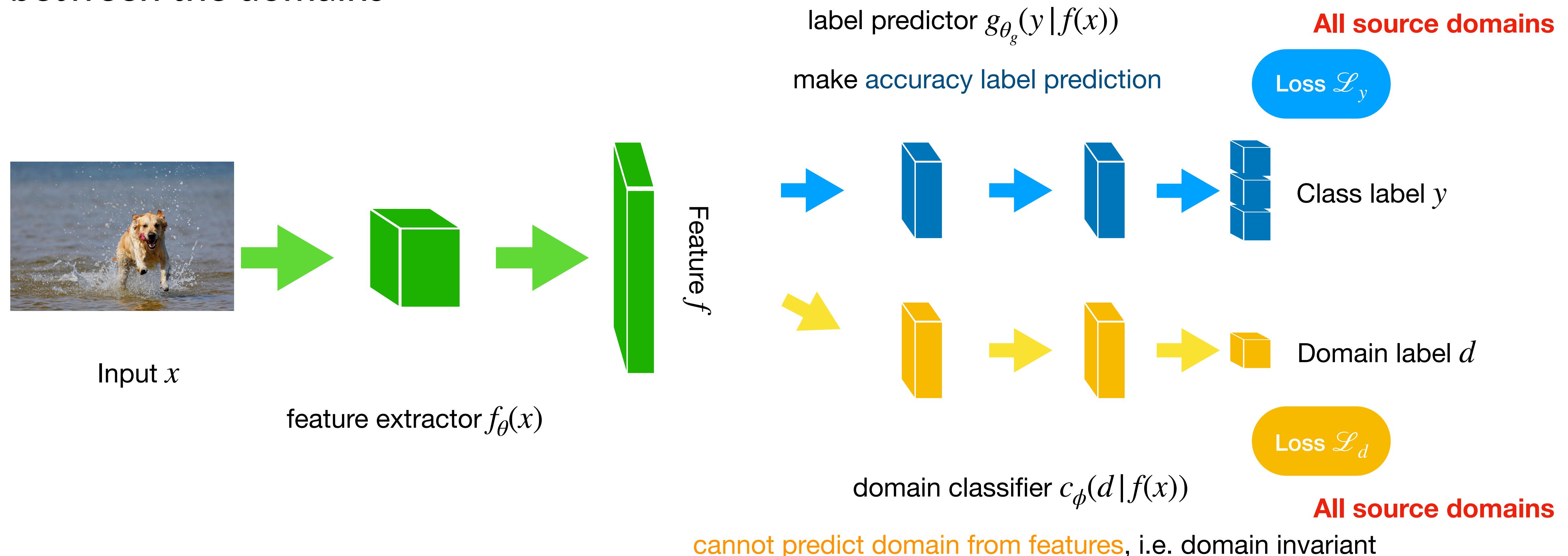
**Key idea:** predictions must be made based on features that cannot be discriminated between the domains



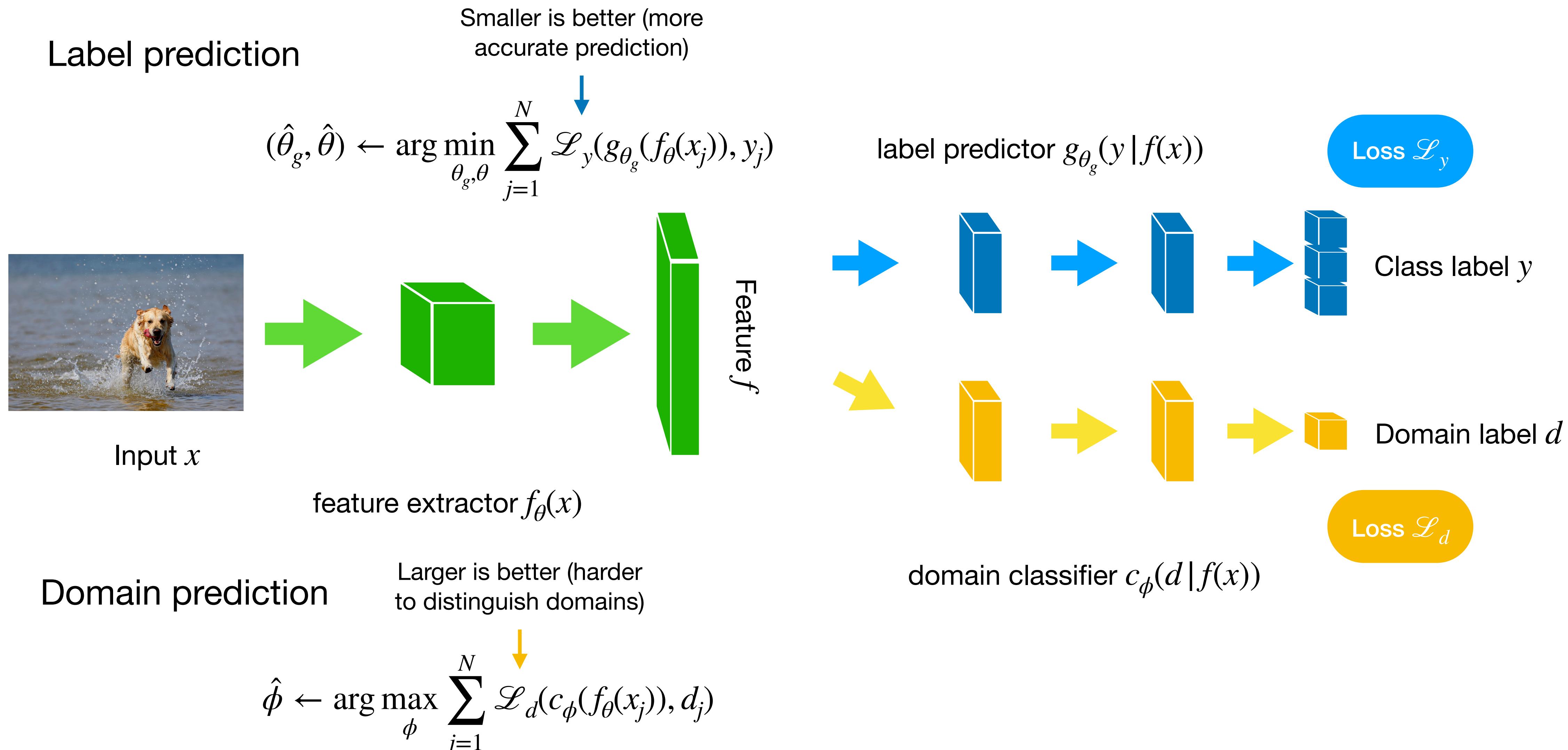
**Question:** Does anyone have ideas on how to use domain adversarial training in the domain generalization setting?

# Domain Adversarial Training in DG

**Key idea:** predictions must be made based on features that cannot be discriminated between the domains



# Domain Adversarial Training in DG



# Domain Adversarial Training in DG

DANN loss in DG

Full algorithm

Label classification loss

$$\min_{\theta} \mathbb{E}_{(x,y)}[\ell(f_{\theta}(x), y)] + \lambda \mathcal{L}_{reg}$$

Explicit regularizer to learn domain-invariant representation

$$\mathcal{L} = \sum_{j=1}^N \mathcal{L}_y(g_{\theta_g}(f_{\theta}(x_j)), y_j) - \lambda \mathcal{L}_d(c_{\phi}(f_{\theta}(x_j)), d_j)$$

1. Randomly initialize encoder  $f_{\theta}$ , label classifier  $g_{\theta_g}$ , domain classifier  $c_{\phi}$

2. Update domain classifier:  $\min_{\phi} \mathcal{L} = \sum_{i=1}^n \mathcal{L}_d(c_{\phi}(f_{\theta}(x_i)), d_i)$

3. Update label classifier & encoder:  $\min_{\theta, \theta_g} \mathcal{L} = \sum_{i=1}^n \mathcal{L}_y(g_{\theta_g}(f_{\theta}(x_i)), y_i) - \lambda \mathcal{L}_d(c_{\phi}(f_{\theta}(x_i)), d_i)$

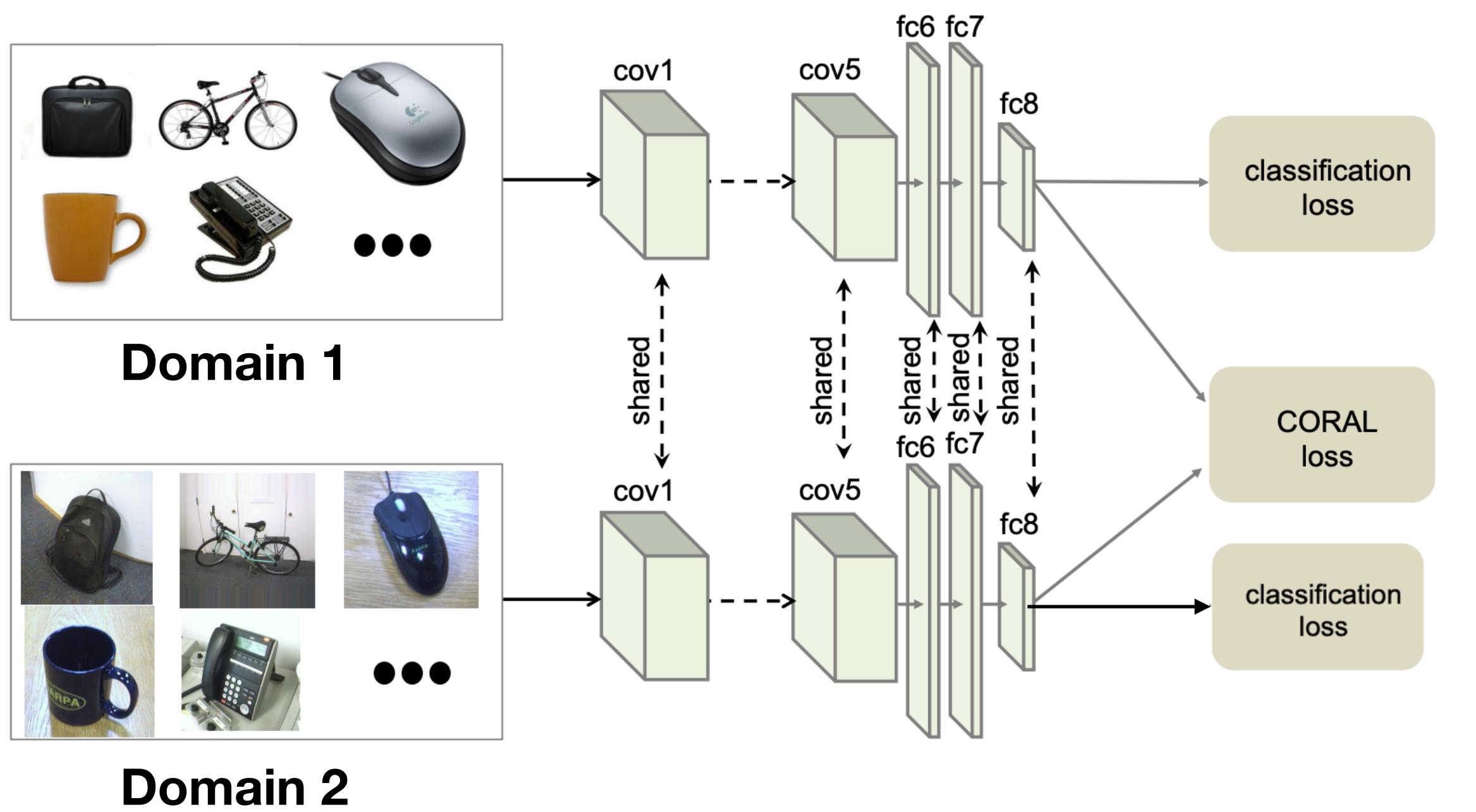
4. Repeat steps 2 & 3

Are there any other ways to [learn domain-invariant features without adversarial optim](#)?

# Alternative Approach — CORAL

**Key idea:** directly aligning representations between different domains with some similarity metrics

CORAL: Correlation Alignment for Domain Adaptation (usually also used in DG)



Notations

$k$ : num of features

$$\mathbf{X}_1 \in \mathbb{R}^{n_1 \times k}$$

$$\mu_1 = \frac{1}{n_1} \mathbf{1}^T \mathbf{X}_1 \in \mathbb{R}^{1 \times k}$$

$$\mathbf{X}_2 \in \mathbb{R}^{n_2 \times k}$$

$$\mu_2 = \frac{1}{n_2} \mathbf{1}^T \mathbf{X}_2 \in \mathbb{R}^{1 \times k}$$

$$C_1 = \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (\mathbf{X}_1 - \mu_1)^T (\mathbf{X}_1 - \mu_1)$$

$$C_2 = \frac{1}{n_2 - 1} \sum_{i=1}^{n_2} (\mathbf{X}_2 - \mu_2)^T (\mathbf{X}_2 - \mu_2)$$

$$\mathcal{L}_{coral} = \frac{1}{4k^2} \|C_1 - C_2\|_F^2$$

Calculate covariance matrices

CORAL loss

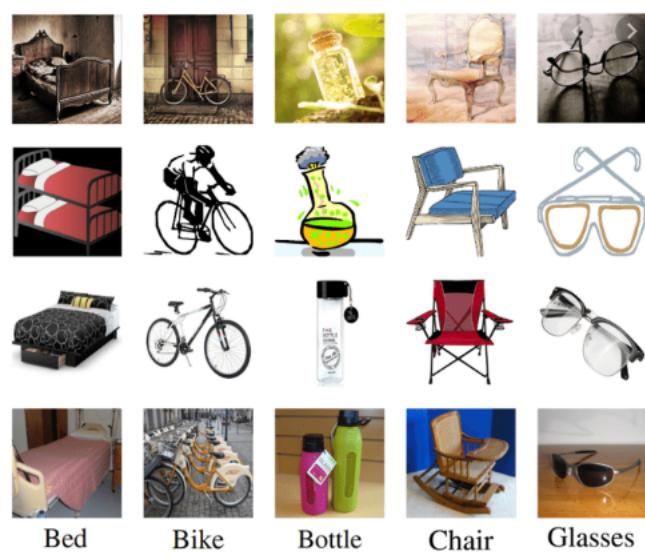
Classification loss

$$\mathcal{L} = \sum_{j=1}^{n_1+n_2} \mathcal{L}_c(f_\theta(x_i), y_i) + \lambda \mathcal{L}_{coral}$$

Explicit regularizer to learn domain-invariant representation

# Results

OfficeHome



ERM

66.5%

CORAL

**68.7%**

DANN

65.9%

DomainNet



40.9%

**41.5%**

38.3%

iWildCam



30.8%

**32.7%**

n/a

# Pros and Cons of Regularization-based Methods

+ General to all kinds of data and networks

+ Some theoretical guarantee

- The regularizer being too harsh / too constraining on the representation

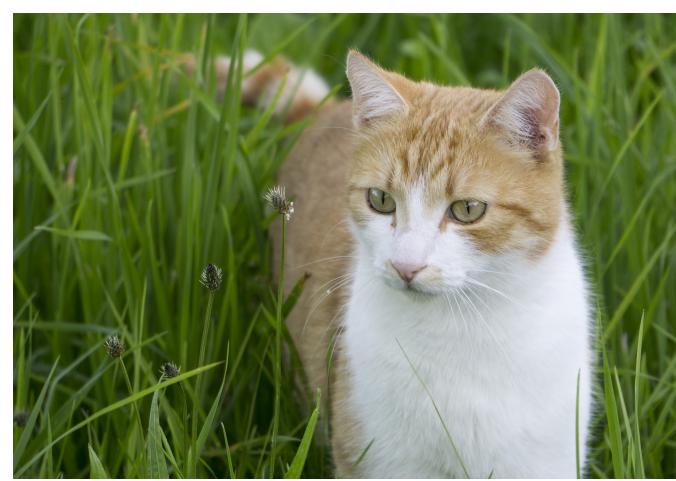
Domain 1:  
water



45% of train data

5% of train data

Domain 2:  
grass



5% of train data

45% of train data

$$\min_{\theta} \mathbb{E}_{(x,y)}[\ell(f_{\theta}(x), y)] + \lambda \mathcal{L}_{reg}$$

←  
Explicit regularizer encourages internal representation to contain **no info about the background**

# Pros and Cons of Regularization-based Methods

+ General to all kinds of data and networks

+ Some theoretical guarantee

- The regularizer being too harsh / too constraining on the representation

These methods can help the performance, but do not always work

Empirical Risk  
Minimization

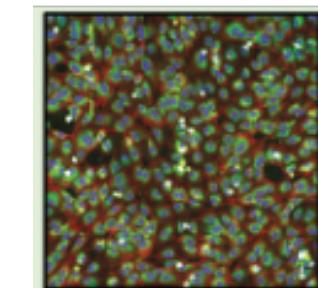


30.8%

iWildCam

Regularization-  
based methods

32.7%  
CORAL



29.9%

RxRx1

28.4%  
CORAL

Are there any other approaches to relax the dependency of the regularizer?

# Plan for Today

## Domain Generalization

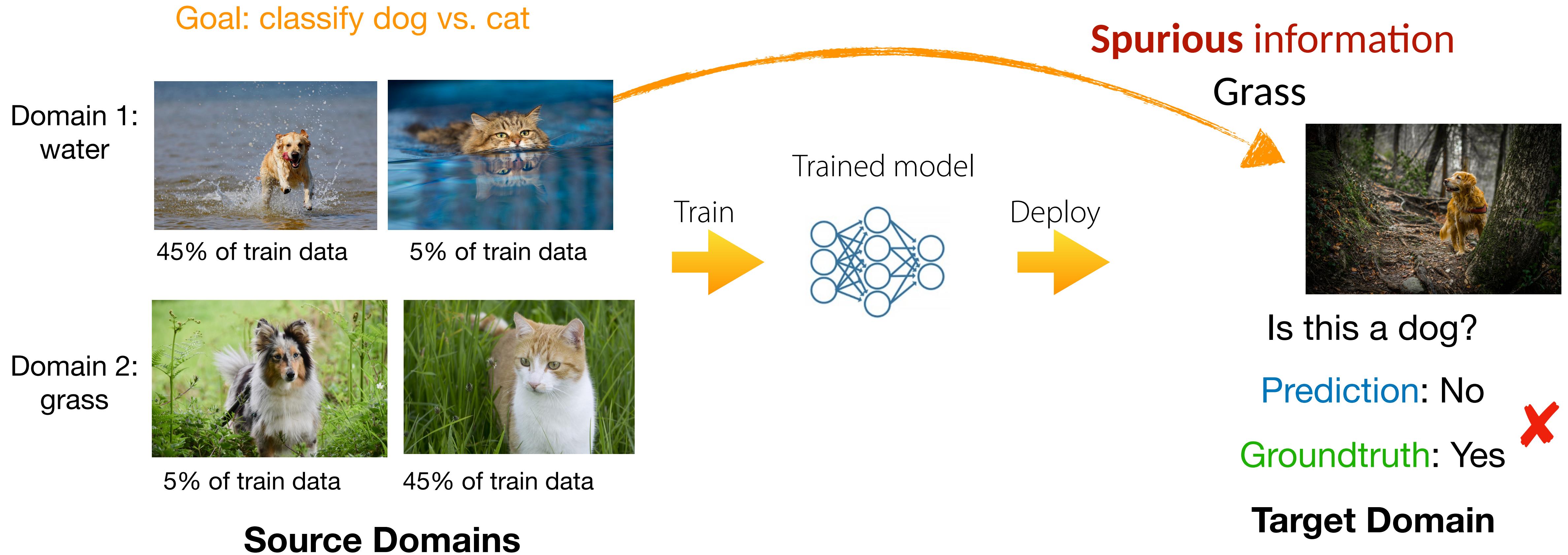
- Problem formulation
- Algorithms
  - Adding explicit regularizers
  - **Data augmentation**

## Goals for this lecture:

- Understand [domain generalization](#): intuition, problem formulation
- Familiarize mainstream DG approaches: [regularization-based](#), [augmentation-based](#)

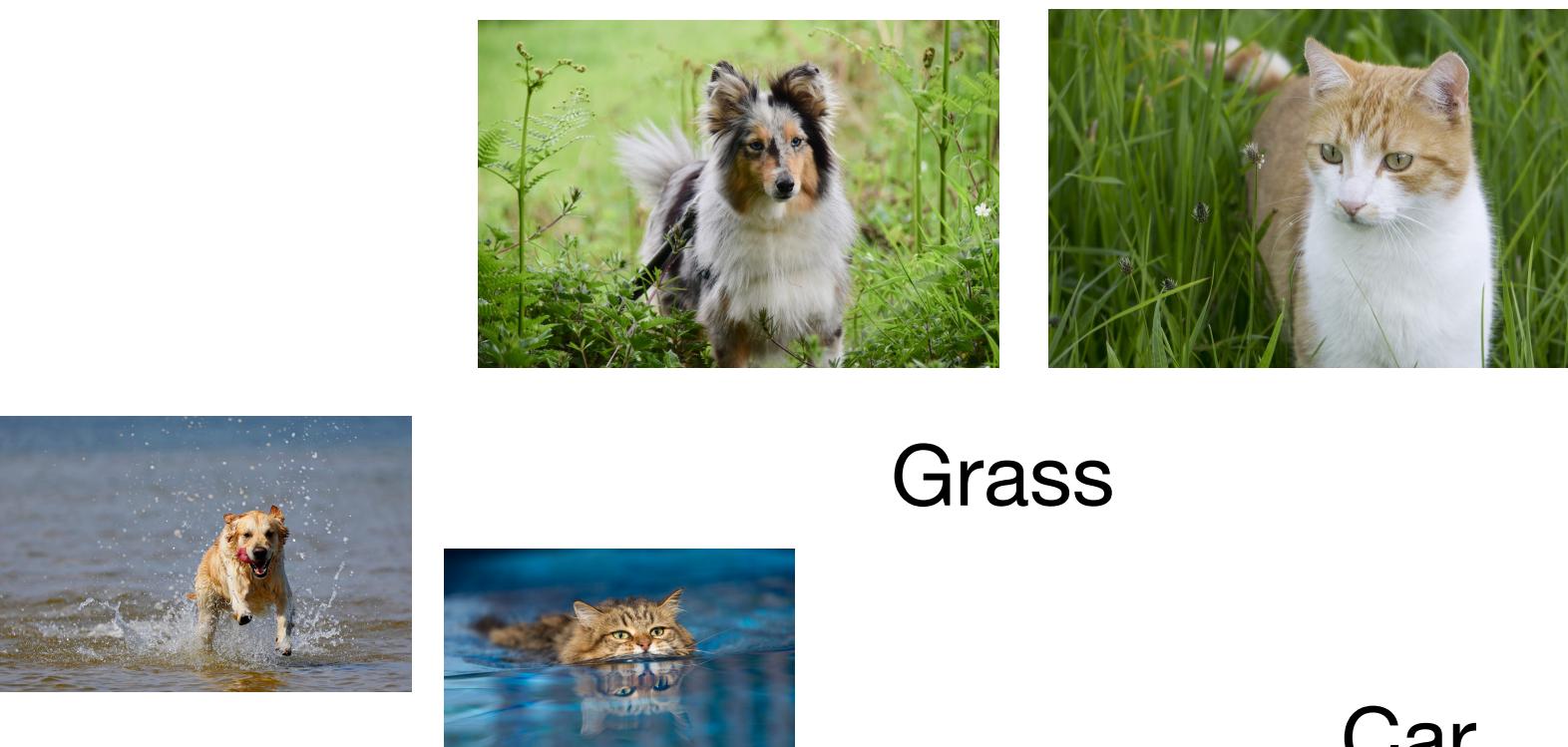
# Recap: Spurious Correlation

**Recap:** spurious correlation between domains and labels



# Data Augmentation

If we can collect more data



Water

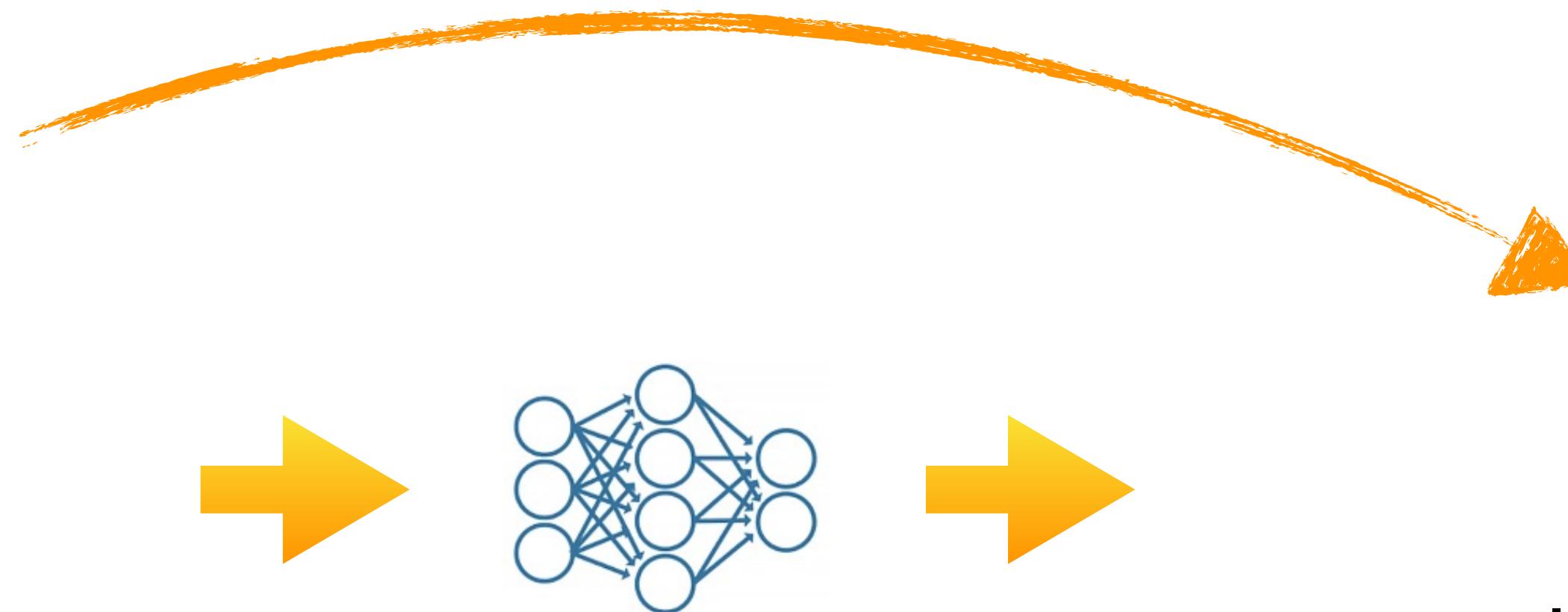
Water

Grass

Car

Keyboard

**Source Domains**

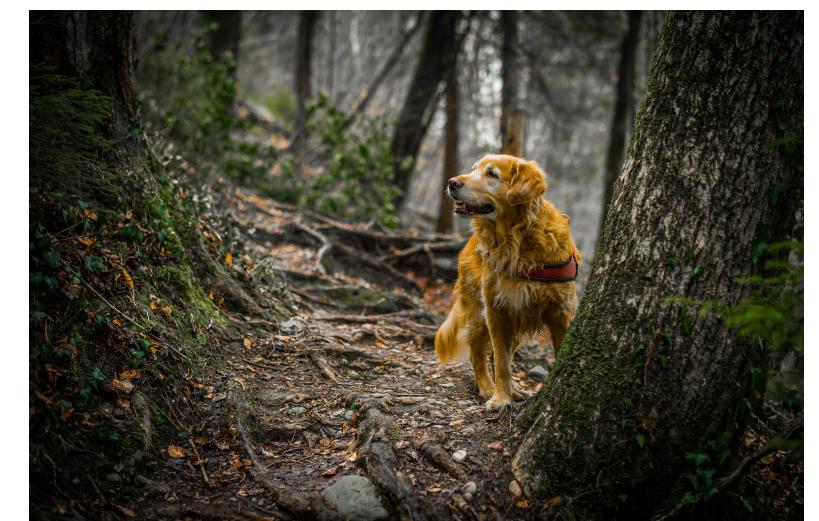


**Challenge**

We can not collect more data —> Let's generate data!

30

30



Is this a dog?

Prediction: Yes

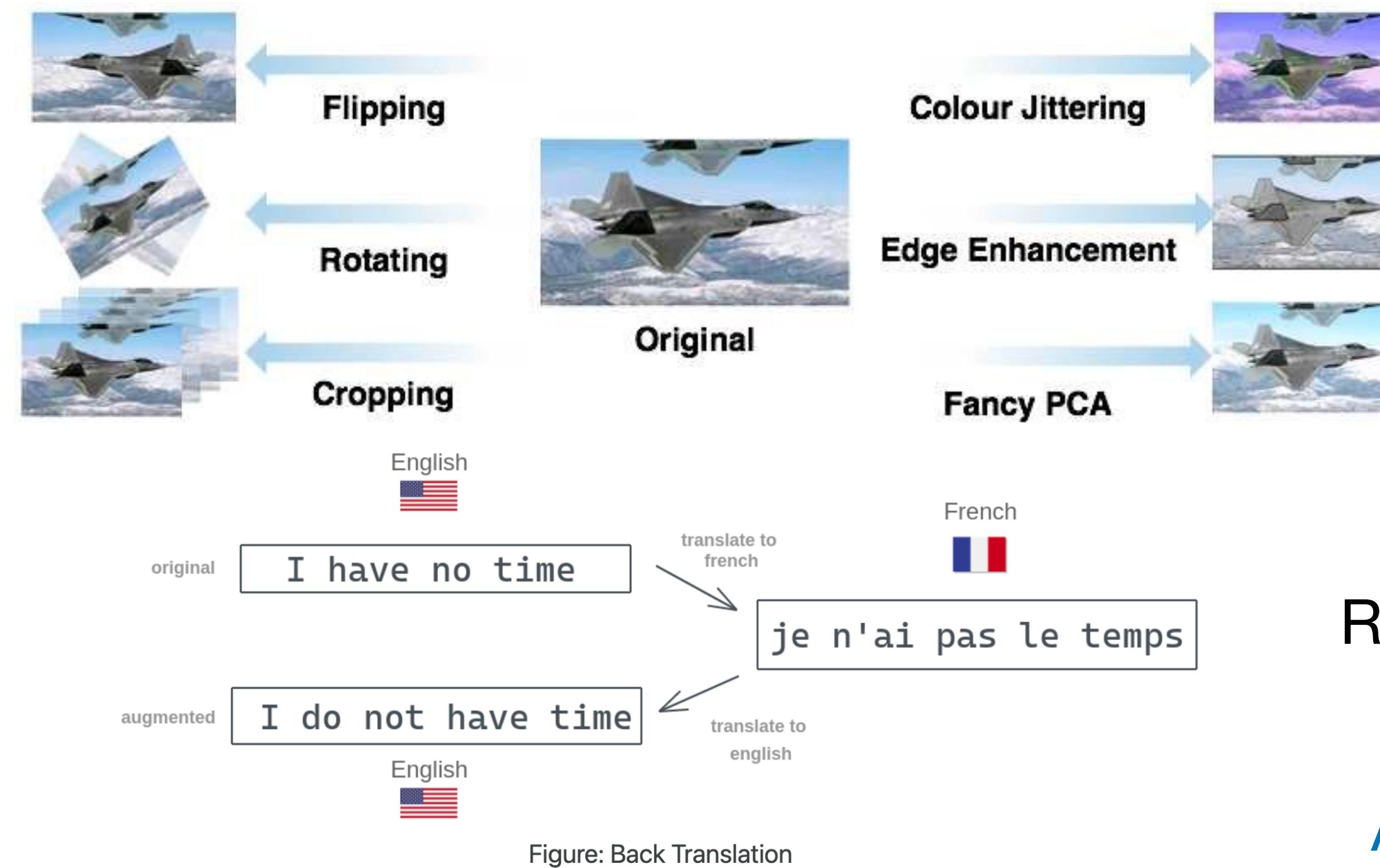


Groundtruth: Yes

**Target Domain**

# Data Augmentation

## Generating data with **simple operators**



Requires knowledge of the problem domain

Any general approaches?

<https://amitness.com/2020/02/back-translation-in-google-sheets/>

# Data Augmentation — Mixup

**Interpolating** training examples

A learning model

$$\mathcal{D}_{tr} = \{x_i, y_i\}_{i=1}^N \rightarrow \text{Classifier},$$

Mixup

$$\tilde{\mathcal{D}}_{tr} = \{\tilde{x}_i, \tilde{y}_i\}_{i=1}^N \rightarrow \text{Classifier},$$

where

$$\tilde{x}_i = \lambda \mathbf{x}_i + (1 - \lambda) \mathbf{x}_j, \tilde{y}_i = \lambda \mathbf{y}_i + (1 - \lambda) \mathbf{y}_j$$

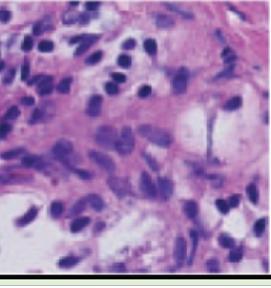
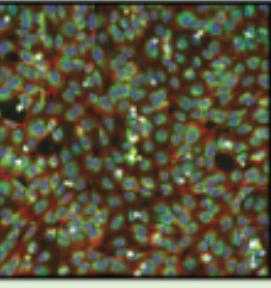
$$\lambda \sim \text{Beta}(\alpha, \beta)$$

Generating some virtual  
examples between two  
classes



# Data Augmentation — Mixup

Mixup can improve the performance on domain generalization

	Empirical Risk Minimization	mixup
	70.3%	71.2%
Camelyon17		
	32.8%	34.2%
FMoW		
	29.9%	26.5%
RxRx1		

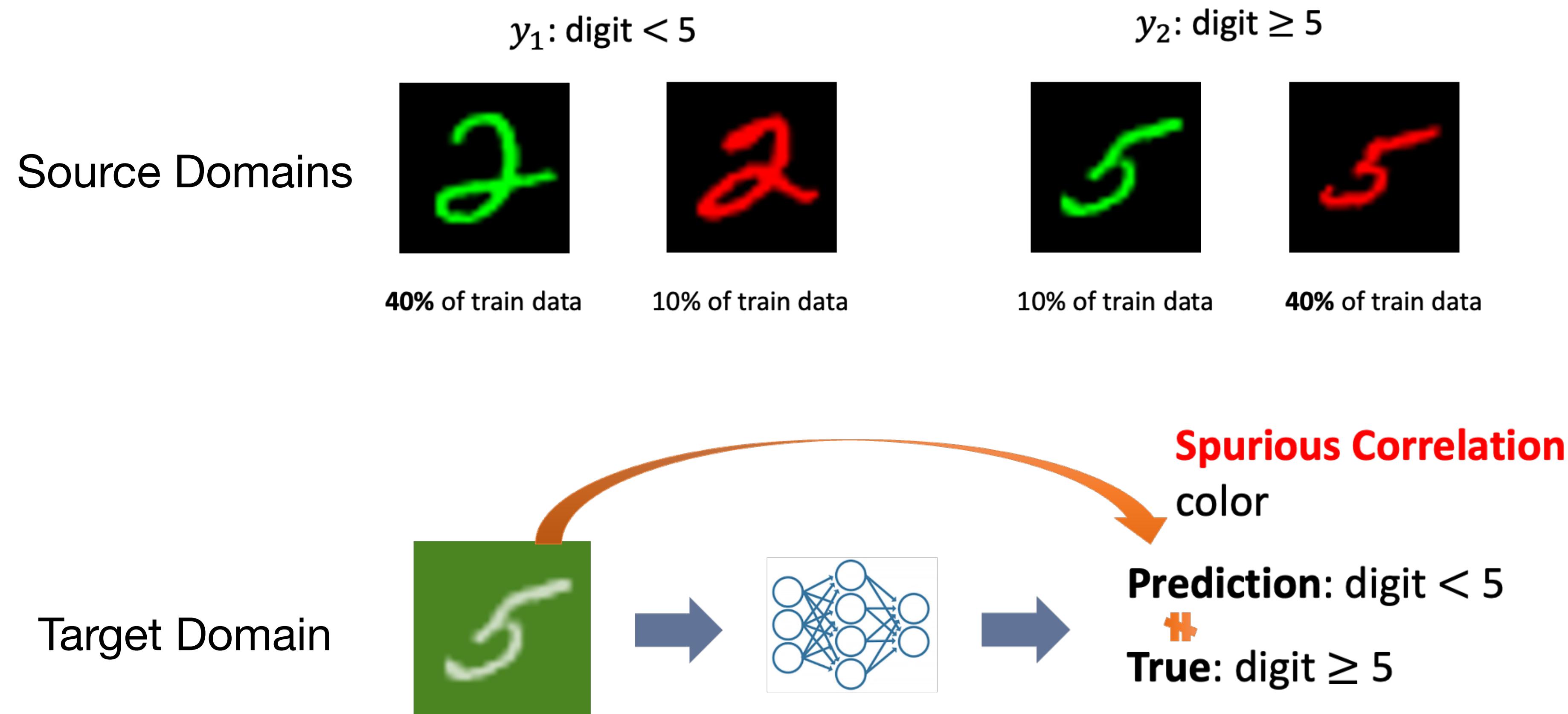
But it is not always good!

Original mixup only focuses  
on data augmentation instead  
of learning domain invariance.

How to Improve it?

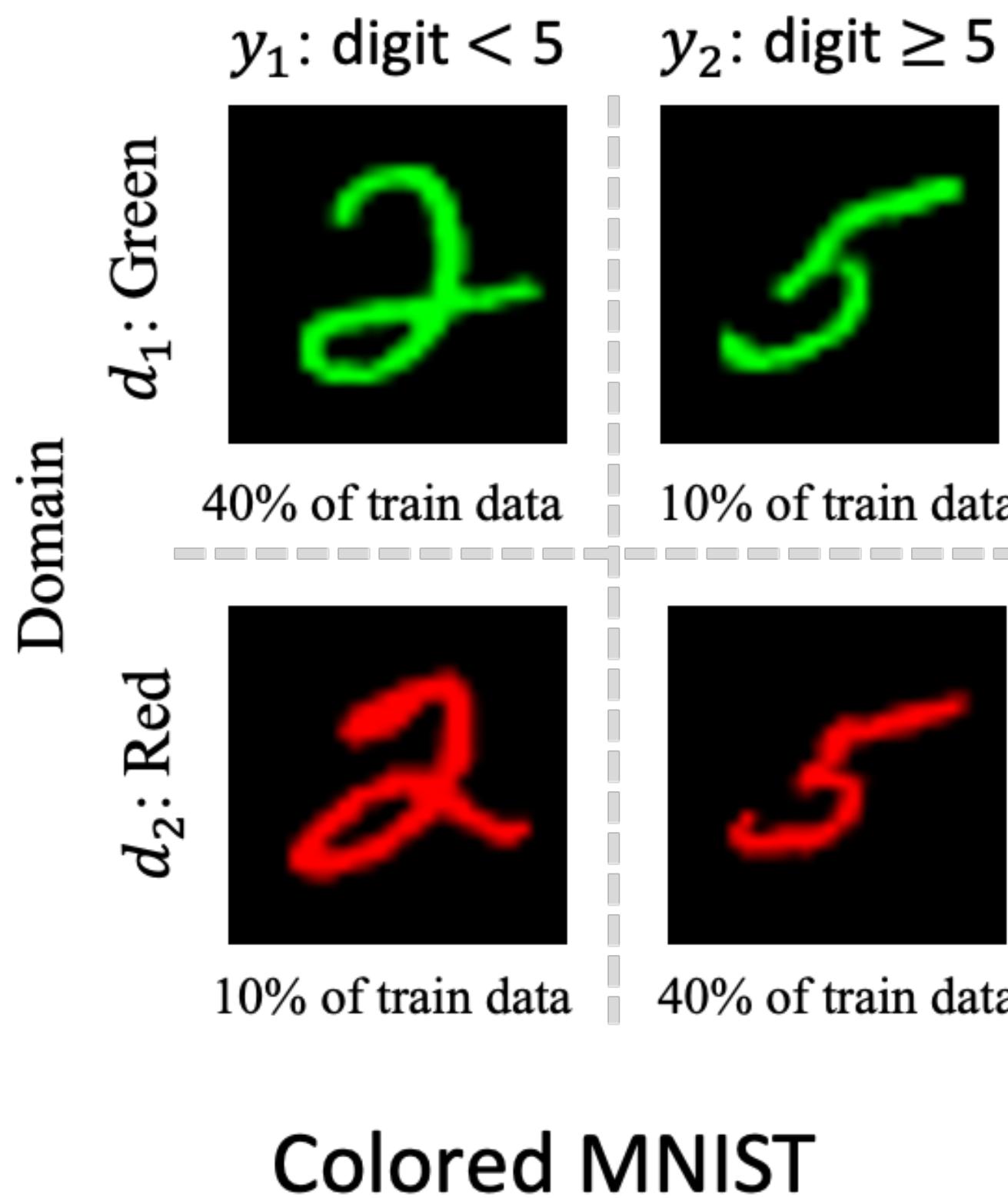
# Data Augmentation — Mixup

A simpler example with spurious correlation



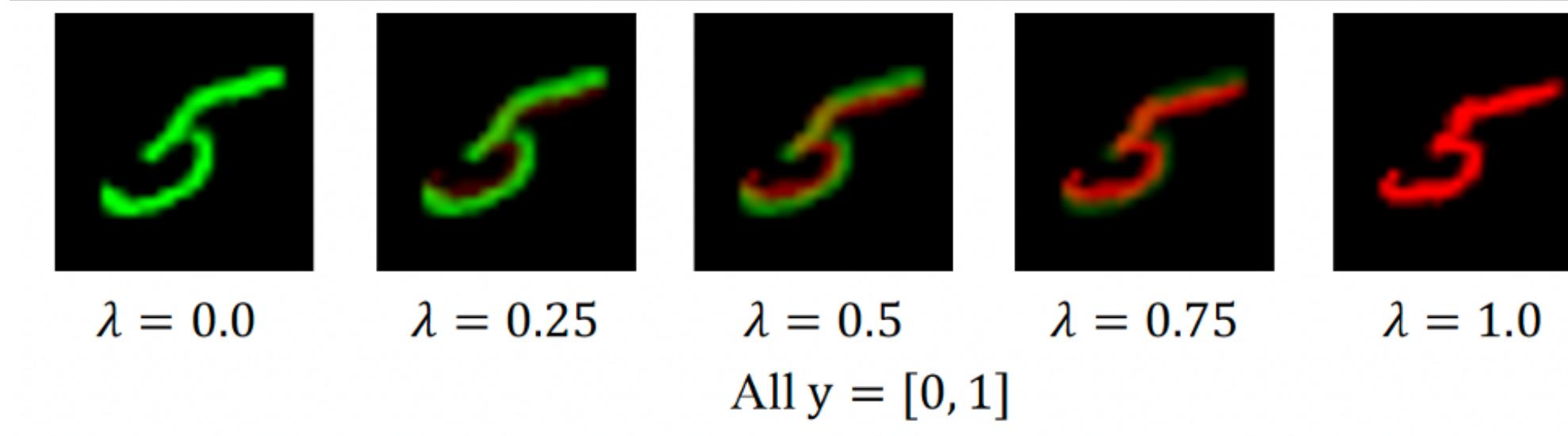
# Can we Improve Mixup? — LISA

**Key idea:** selective interpolate examples to emphasize invariant information



**Mixup:**  $x_{mix} = \lambda \mathbf{x}_i + (1 - \lambda) \mathbf{x}_j, y_{mix} = \lambda \mathbf{y}_i + (1 - \lambda) \mathbf{y}_j$   
 $\lambda \sim \text{Beta}(\alpha, \beta)$

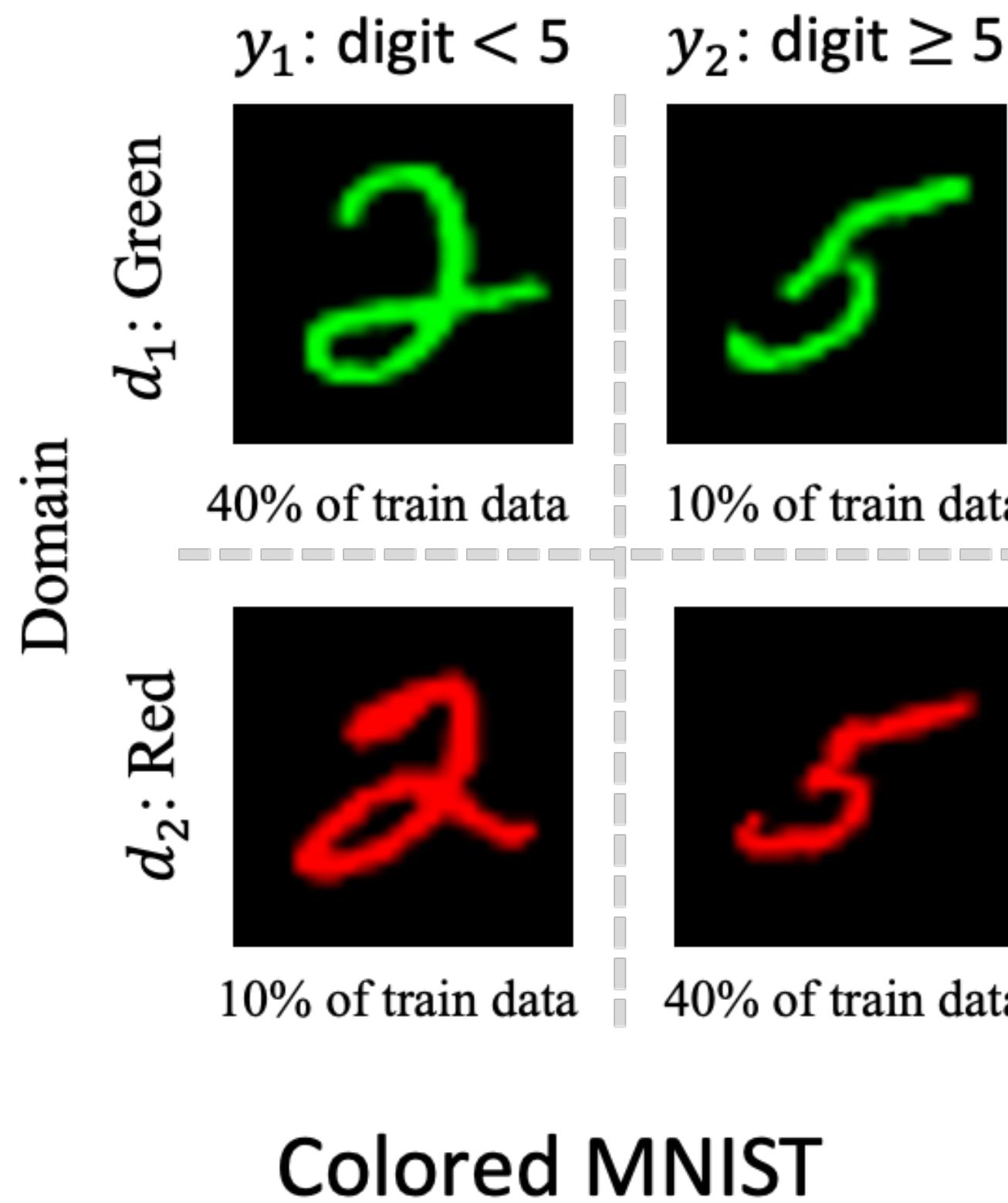
**Intra-label LISA** – Interpolates samples with the **same label** but **different domains** ( $d_i \neq d_j, y_i = y_j$ )



Different background, same label

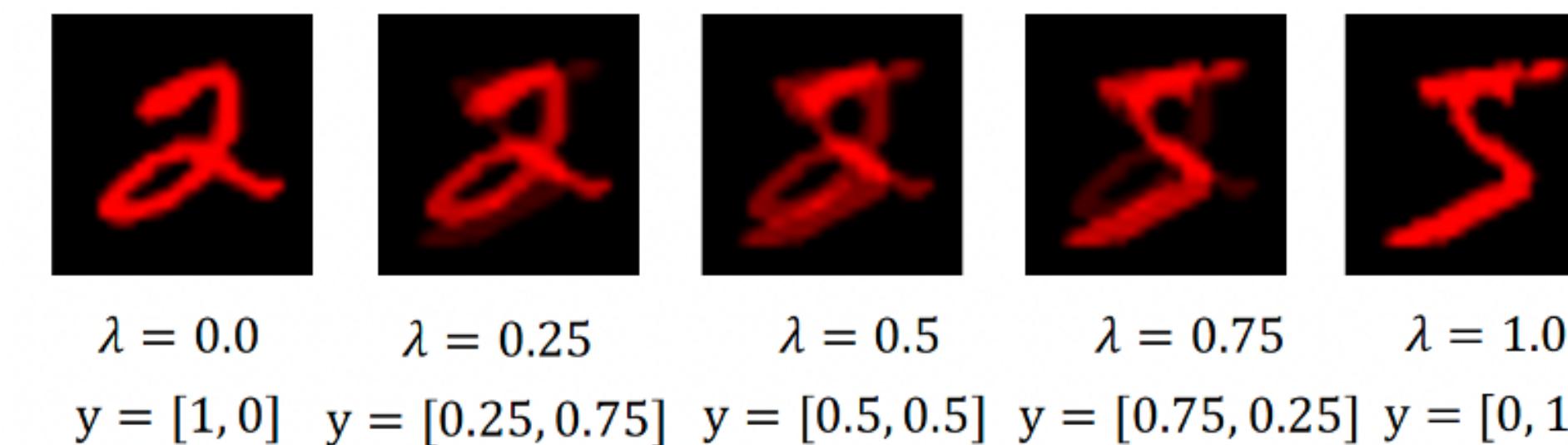
# Can we Improve Mixup? — LISA

**Key idea:** selective interpolate examples to emphasize invariant information



**Mixup:**  $x_{mix} = \lambda \mathbf{x}_i + (1 - \lambda) \mathbf{x}_j, y_{mix} = \lambda \mathbf{y}_i + (1 - \lambda) \mathbf{y}_j$   
 $\lambda \sim \text{Beta}(\alpha, \beta)$

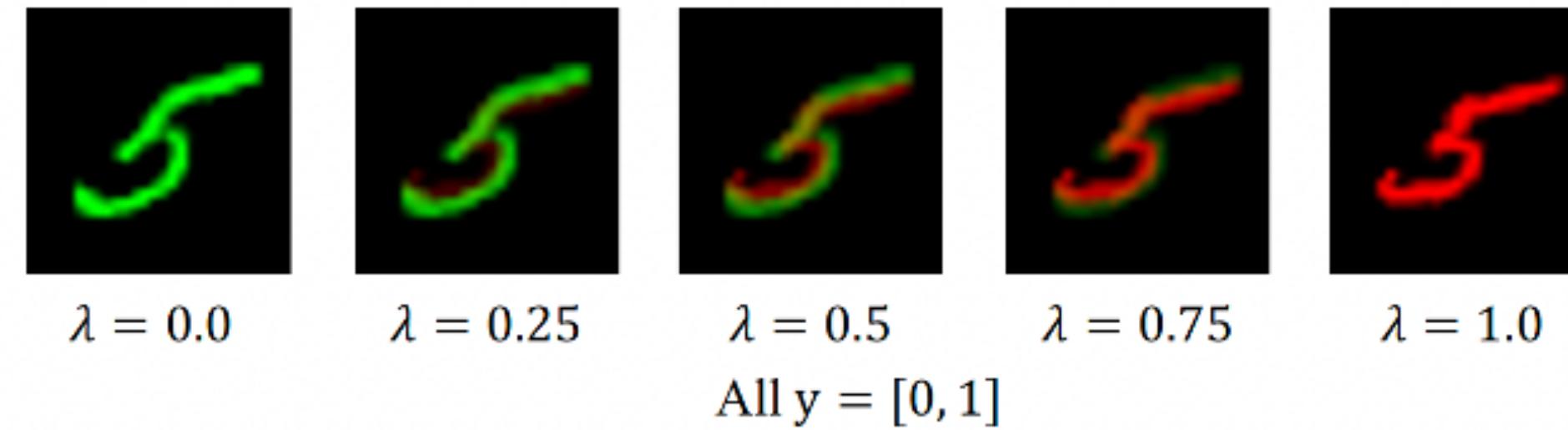
**Intra-domain LISA** – Interpolates samples with the **different label** but **same domains** ( $d_i = d_j, y_i \neq y_j$ )



Domain information is **not** the reason for the label change

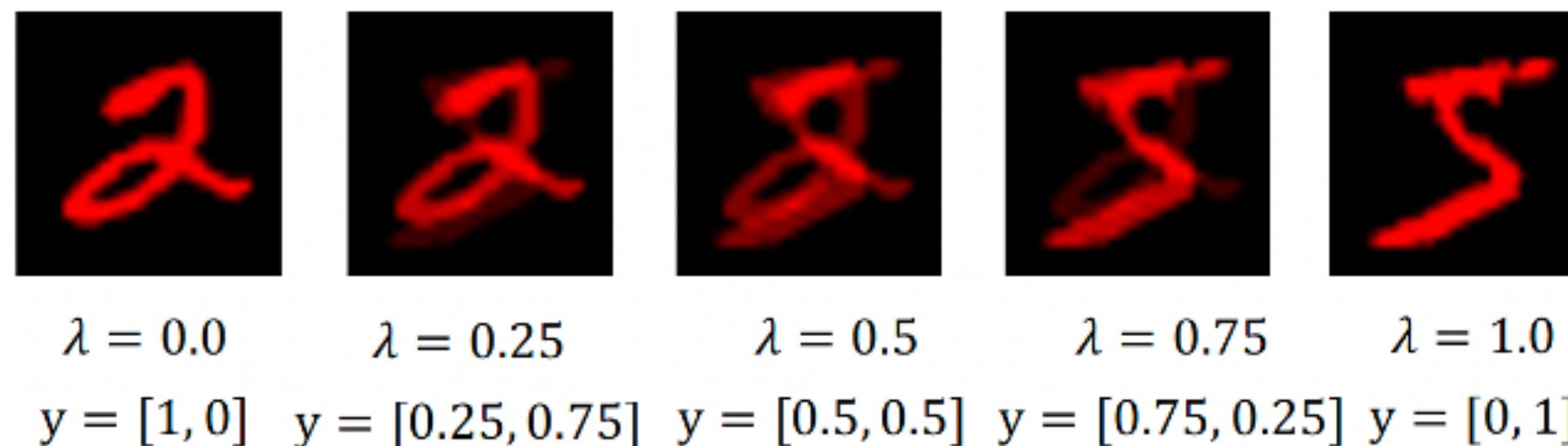
# Can we Improve Mixup? — LISA

Intra-label  
LISA



+ more domains  
+ spurious correlations  
are not very strong

Intra-domain  
LISA



+ domain information is highly  
spuriously correlated with the  
label

$p_{sel}$ : Determine intra-label LISA or intra-domain LISA at each iteration

# Full Algorithm of LISA

1. Randomly initialize the model parameter  $\theta$
2. Sample strategy  $s \sim Bernoulli(p_{sel})$
3. Sample a batch of examples  $\mathcal{B}$ 
  - (i) If  $s=0$ , for each example  $(x_i, y_i)$  in  $\mathcal{B}$ , sample  $(x_j, y_j)$  that satisfies  $(y_i = y_j)$  and  $(d_i \neq d_j)$
  - (ii) If  $s=1$ , for each example  $(x_i, y_i)$  in  $\mathcal{B}$ , sample  $(x_j, y_j)$  that satisfies  $(y_i \neq y_j)$  and  $(d_i = d_j)$
4. Use interpolated examples to update the model
5. Repeat steps 3 & 4



Intra-label LISA



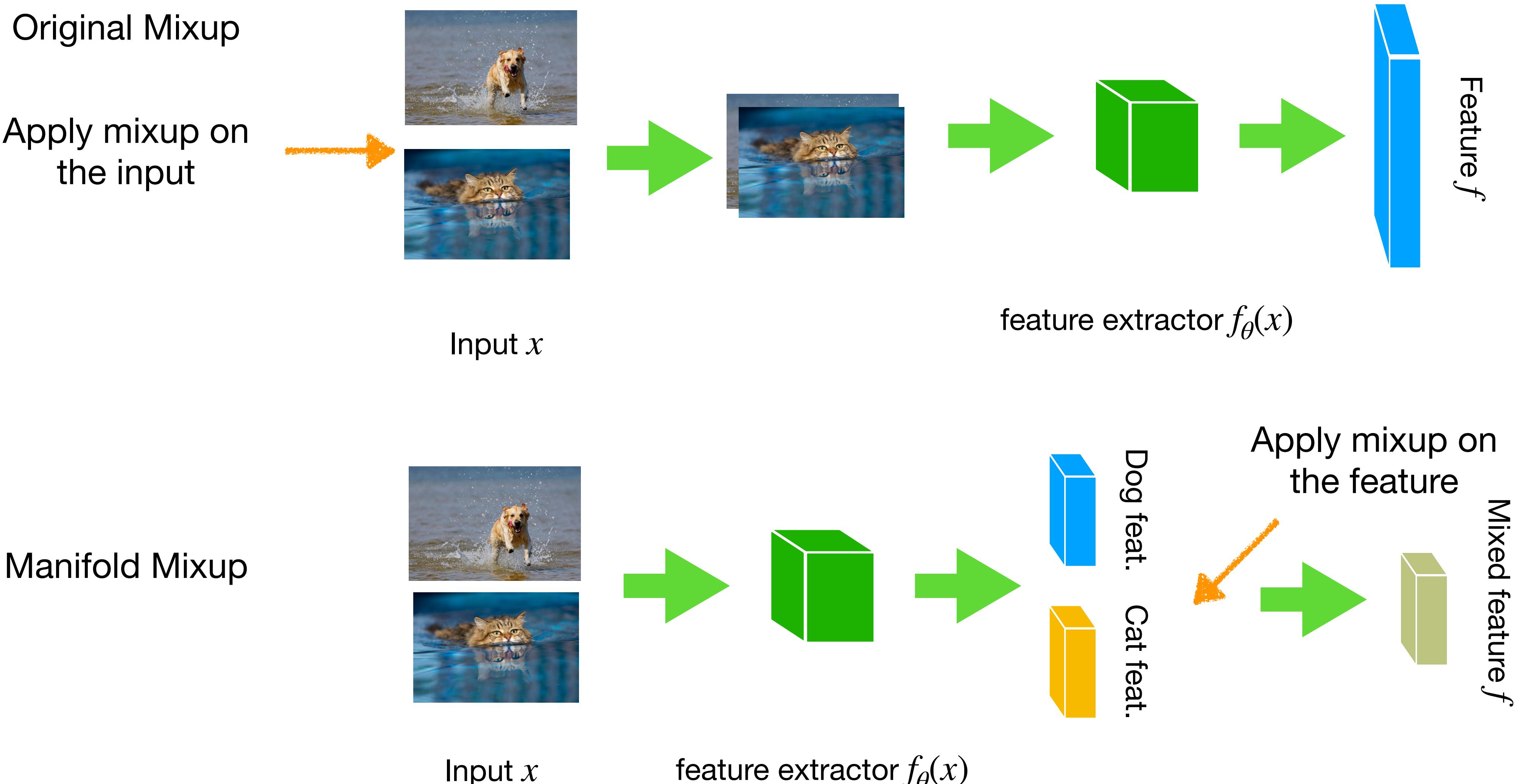
Intra-domain LISA

# Results

	ERM	Regularization-based (CORAL)	Augmentation-based (LISA)
Camelyon17	70.3%	74.7%	<b>77.1%</b>
FMoW	32.3%	34.6%	<b>35.5%</b>
RxRx1	29.9%	28.4%	<b>31.9%</b>
Amazon	53.8%	53.8%	<b>54.7%</b>
iWildCAM	30.8%	<b>32.7%</b>	27.6%
OGB-MolPCBA	<b>28.3%</b>	17.9%	27.5%

LISA can also work on text data, how to apply mixup?

# Manifold Mixup



# Invariance Analysis

Metrics: Accuracy of domain prediction      Divergence of predictions among domains

	IP <sub>adp</sub> ↓				IP <sub>kl</sub> ↓			
	CMNIST	Waterbirds	Camelyon17	MetaShift	CMNIST	Waterbirds	Camelyon17	MetaShift
ERM	82.85%	94.99%	49.43%	67.98%	6.286	1.888	1.536	1.205
Vanilla mixup	92.34%	94.49%	52.79%	69.36%	4.737	2.912	0.790	1.171
IRM	69.42%	95.12%	47.96%	67.59%	7.755	1.122	0.875	1.148
IB-IRM	74.72%	94.78%	48.37%	67.39%	1.004	3.563	0.756	1.115
V-REx	63.58%	93.32%	61.38%	68.38%	3.190	3.791	1.281	1.094
<b>LISA (ours)</b>	<b>58.42%</b>	<b>90.28%</b>	<b>45.15%</b>	<b>66.01%</b>	<b>0.567</b>	<b>0.134</b>	<b>0.723</b>	<b>1.001</b>

LISA leads to **greater domain invariance** than prior methods with explicit regularizers

# Regularization-based v.s. Augmentation-based Methods

## Regularization-based Method

- + General to all kinds of data and networks
- + Some theoretical guarantee
- Rely on the design of regularizers

## Augmentation-based Method

- + Easy to understand and simple to implement
- + No need to worry about how to design regularizers
- Largely limited to classification

# Plan for Today

## Domain Generalization

- Problem formulation
- Algorithms
  - Adding explicit regularizers
  - Data augmentation

## Goals for this lecture:

- Understand [domain generalization](#): intuition, problem formulation
- Familiarize mainstream DG approaches: [regularization-based](#), [augmentation-based](#)