

# **Transfer learning and Domain adaptation**

## **Vision & Language**

# Transfer from ImageNet (source)

## Transfer as generic features

Brut Deep features (learned from ImageNet)

(== a learned embedding from Image to vector representation)

Retrieval



## Transfer learning (from source to target)

Frozen features + SVM => solution to small datasets

Frozen features + Deep

Fine tuning not easy in that case (small datasets)

# Transfer from source(=ImageNet task) to target task

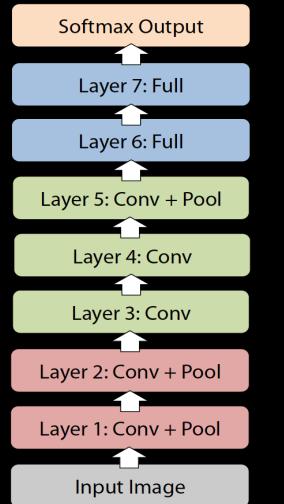
**Source:** ImageNet (dataset + 100 classes) => AlexNet trained

**Target:** new dataset Cal-101 and new classification task with 101 classes => Chopped

AlexNet (layer i) + SVM trained on

## Architecture of Krizhevsky et al.

- 8 layers total
- Trained on Imagenet dataset [Deng et al. CVPR'09]
- 18.2% top-5 error
- Our reimplementation:  
18.1% top-5 error



## Tapping off Features at each Layer

Plug features from each layer into linear SVM or soft-max

	Cal-101 (30/class)	Cal-256 (60/class)
SVM (1)	$44.8 \pm 0.7$	$24.6 \pm 0.4$
SVM (2)	$66.2 \pm 0.5$	$39.6 \pm 0.3$
SVM (3)	$72.3 \pm 0.4$	$46.0 \pm 0.3$
SVM (4)	$76.6 \pm 0.4$	$51.3 \pm 0.1$
SVM (5)	<b><math>86.2 \pm 0.8</math></b>	$65.6 \pm 0.3$
SVM (7)	<b><math>85.5 \pm 0.4</math></b>	<b><math>71.7 \pm 0.2</math></b>
Softmax (5)	$82.9 \pm 0.4$	$65.7 \pm 0.5$
Softmax (7)	<b><math>85.4 \pm 0.4</math></b>	<b><math>72.6 \pm 0.1</math></b>

=> Results better than SoA CV methods on Cal-101!

# Transfer: fine-tuning of a deep model on target task

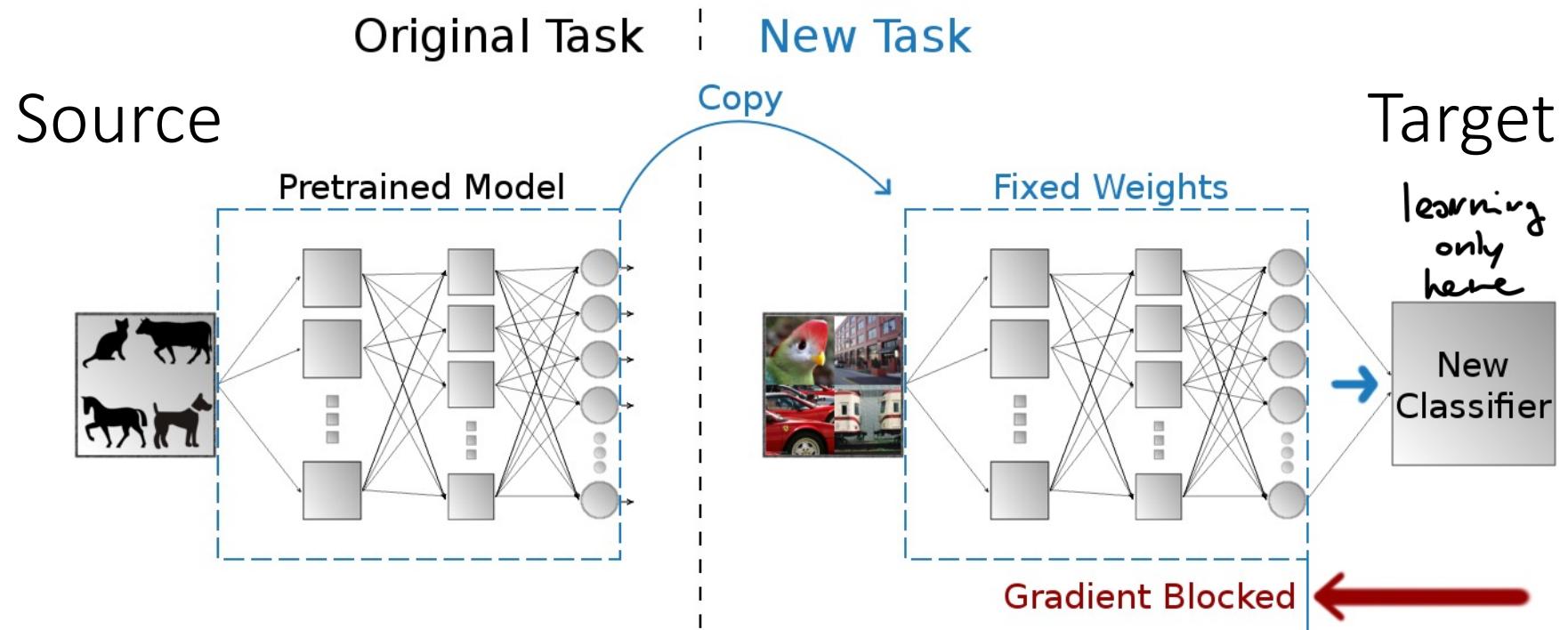
Train a deep (AlexNet) on source (ImageNet)

Keep the deep params. for target and complete with a small deep on top (fully trained on target task)

Fine-tune the whole model on target data

Challenge: only limited target data, careful about overfitting

Solution: Freeze the gradient's update for AlexNet part



# Transfer: fine-tuning of a deep model on target task

Train a deep (AlexNet) on source (ImageNet)

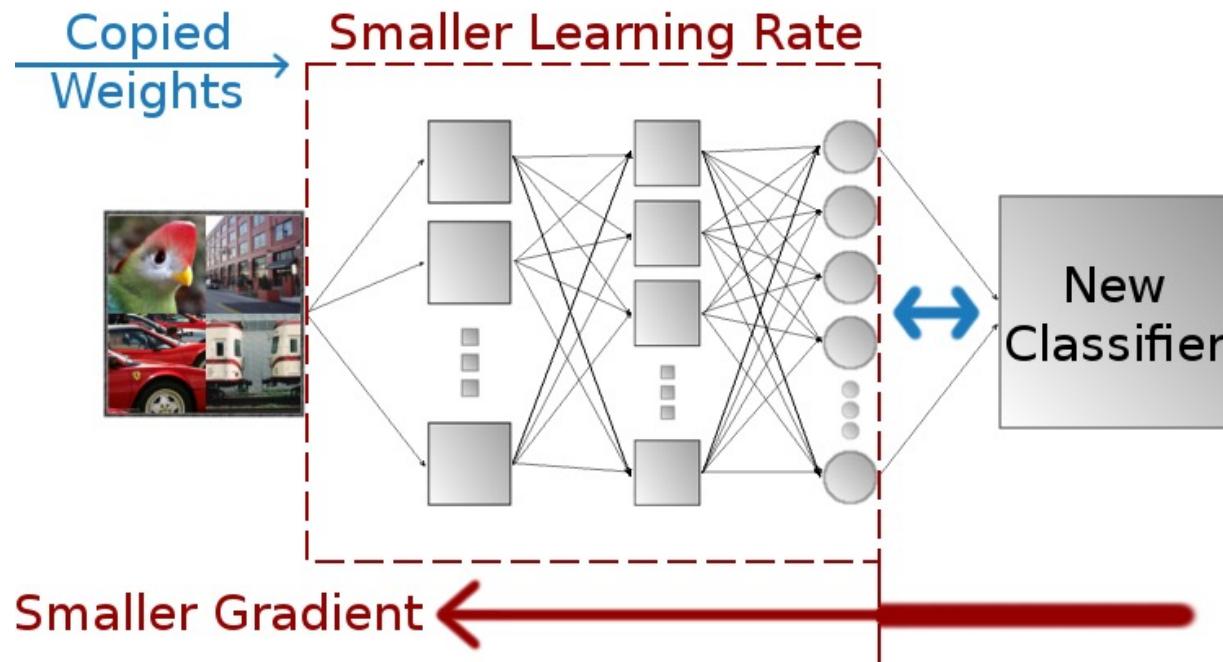
Keep the deep params. for target and complete with a small deep on top (fully trained on target task)

Fine-tune the whole model on target data

Challenge: only limited target data, careful about overfitting

Solution: Freeze the gradient's update for AlexNet part

Other solution: use smaller gradient's update for AlexNet part

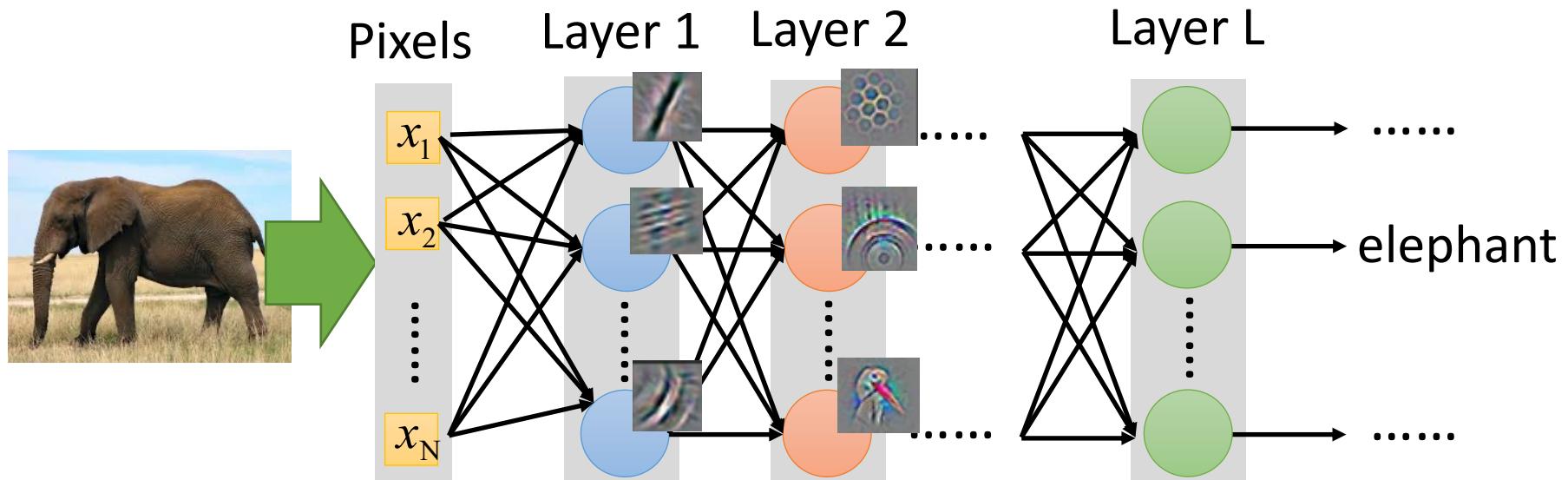


# Transfer: which parts of the deep?

skip

Which layer(s) can be transferred (copied)?

- Speech: usually copy the last few layers
- Image: usually copy the first few layers



# Transfer: which supervision?

- Task description
  - Source data:  $(x^s, y^s)$   A large amount
  - Target data:  $(x^t, y^t)$   (Very) little

Rq: Few/One-shot learning: only a few/one examples in target domain

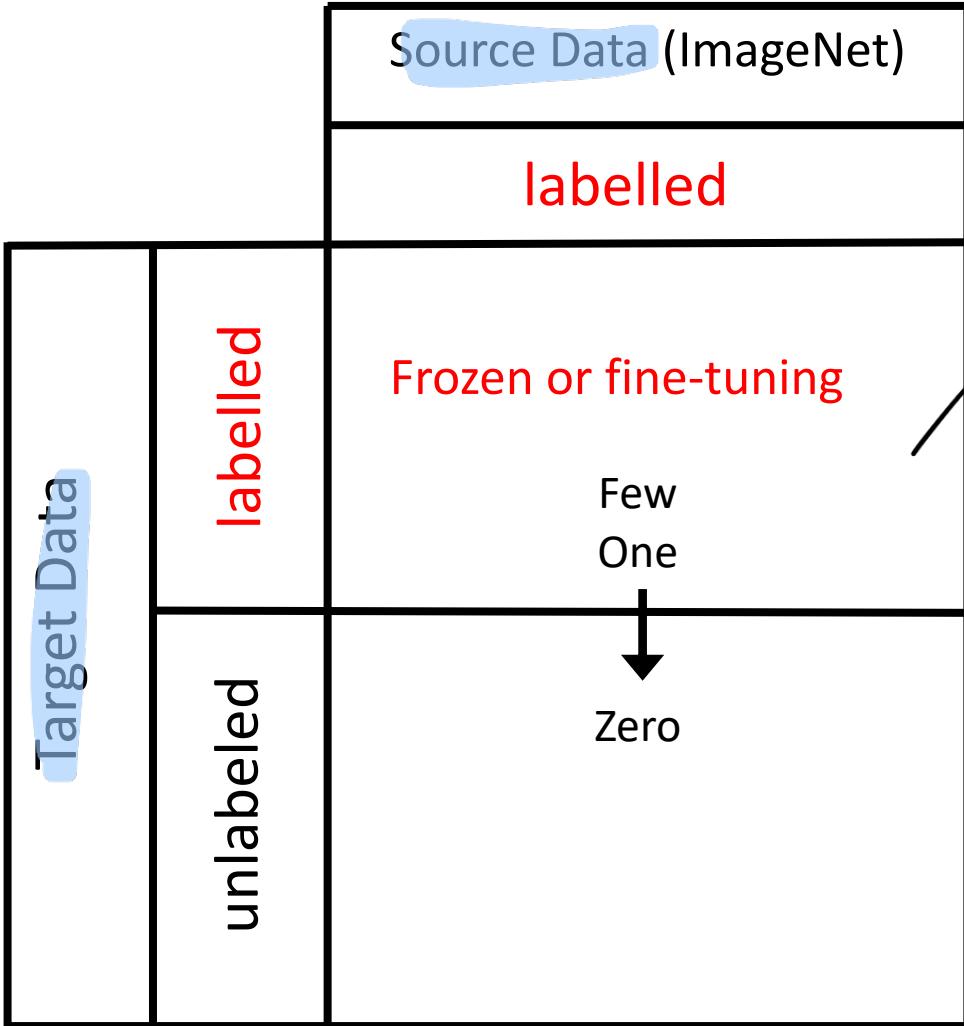
Many different contexts:

In vision: from large dataset ([ImageNet](#)) to small datasets ([VOC2007](#))

In speech: (supervised) speaker adaption

- Source data: audio data and transcriptions from many speakers
- Target data: audio data and its transcriptions of specific user

# More on transfer framework



a very few labelled data

Main purposes:  
Similar visual domain?  
Same tasks (ie class)?

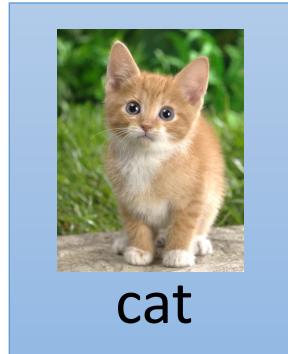
Source  $\approx$  Target similar  $\rightarrow$  easy

Source  $\neq$  Target  $\rightarrow$  hard

source = medical  
target = driving

# Similar domain: ImageNet task => Dog/Cat task

Target:  
Dog/Cat  
Classifier



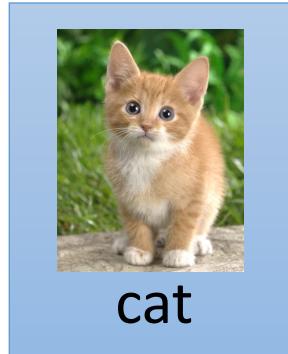
Data *not directly related to* the task considered



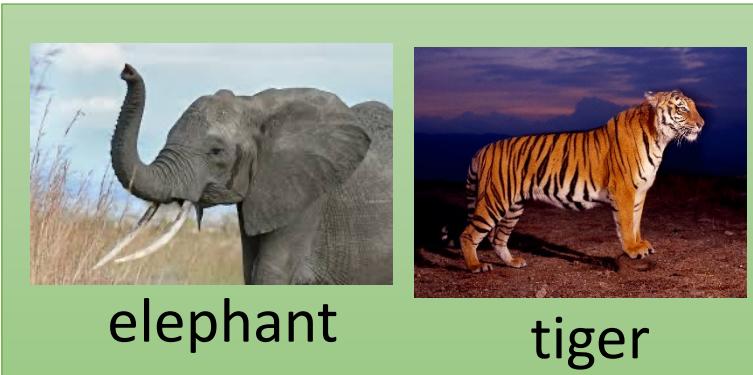
ImageNet: Similar domain,  
different task (1000 classes but NOT Dog and Cat classes)

# General Framework for Transfer Learning

Target:  
Dog/Cat  
Classifier



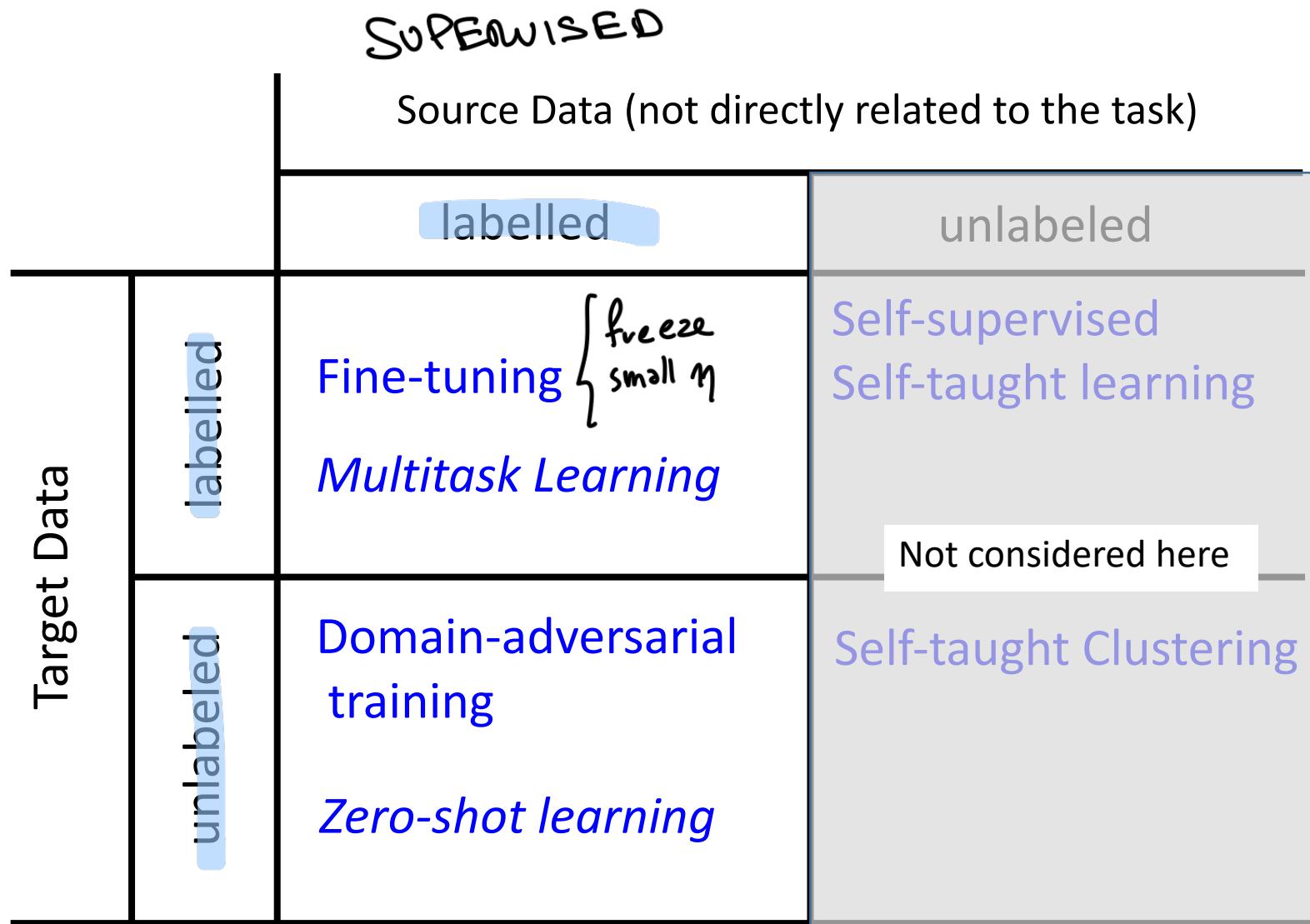
Data *not directly related to* the task considered



Similar domain, completely  
different tasks

Different domains, same task  
big problem

# General Framework for Transfer Learning

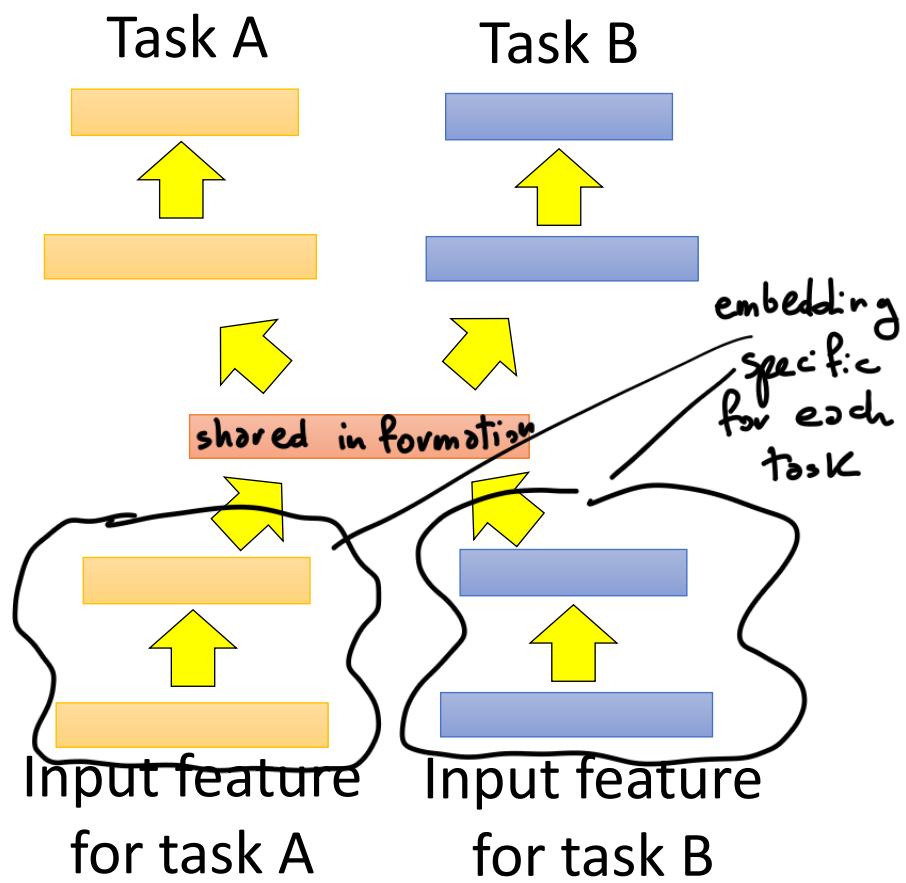
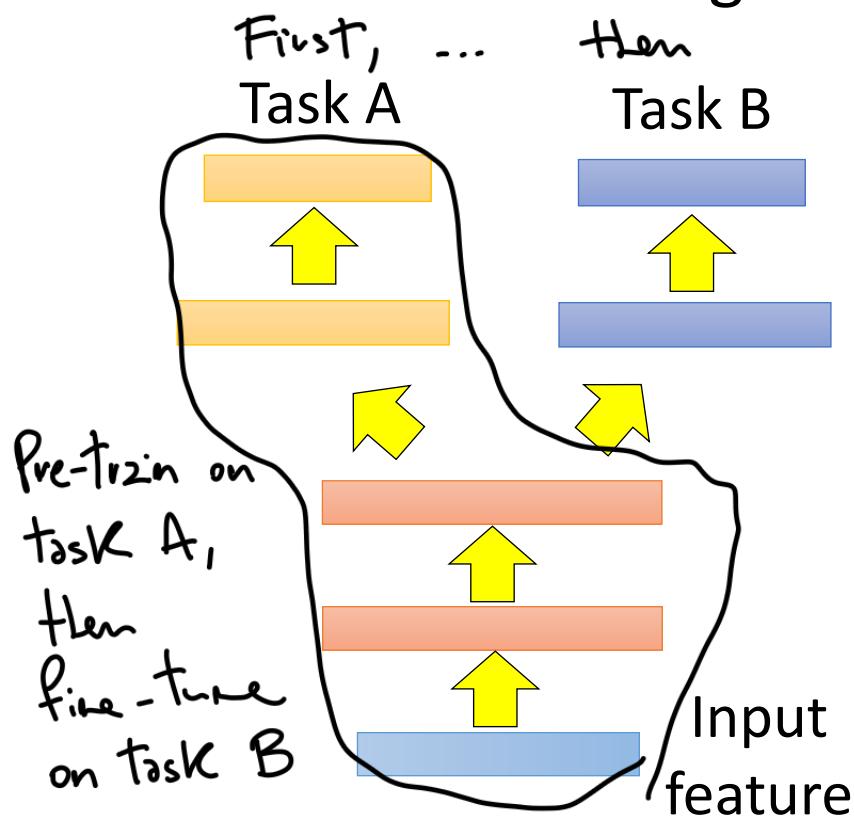


# General Framework for Transfer Learning

		Source Data (not directly related to the task)	
		labelled	unlabeled
Target Data	labelled	Fine-tuning <b><i>Multitask Learning</i></b>	Not considered here
	unlabeled		Not considered here

# Multitask Learning

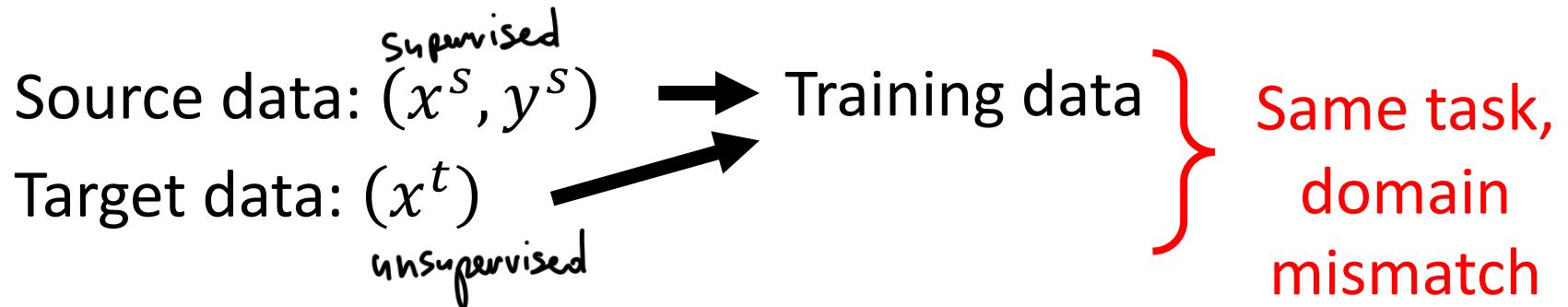
- The multi-layer structure makes NN suitable for multitask learning



# Transfer Learning - Overview

		Source Data (not directly related to the task)	
		labelled	unlabeled
Target Data	labelled	Fine-tuning <i>{ together, but sequentially }</i>  <i>Multitask Learning</i> <i>parallel, different encoders</i>	Not considered here
	unlabeled	<b>Domain adaptation- adversarial training</b>	Not considered here

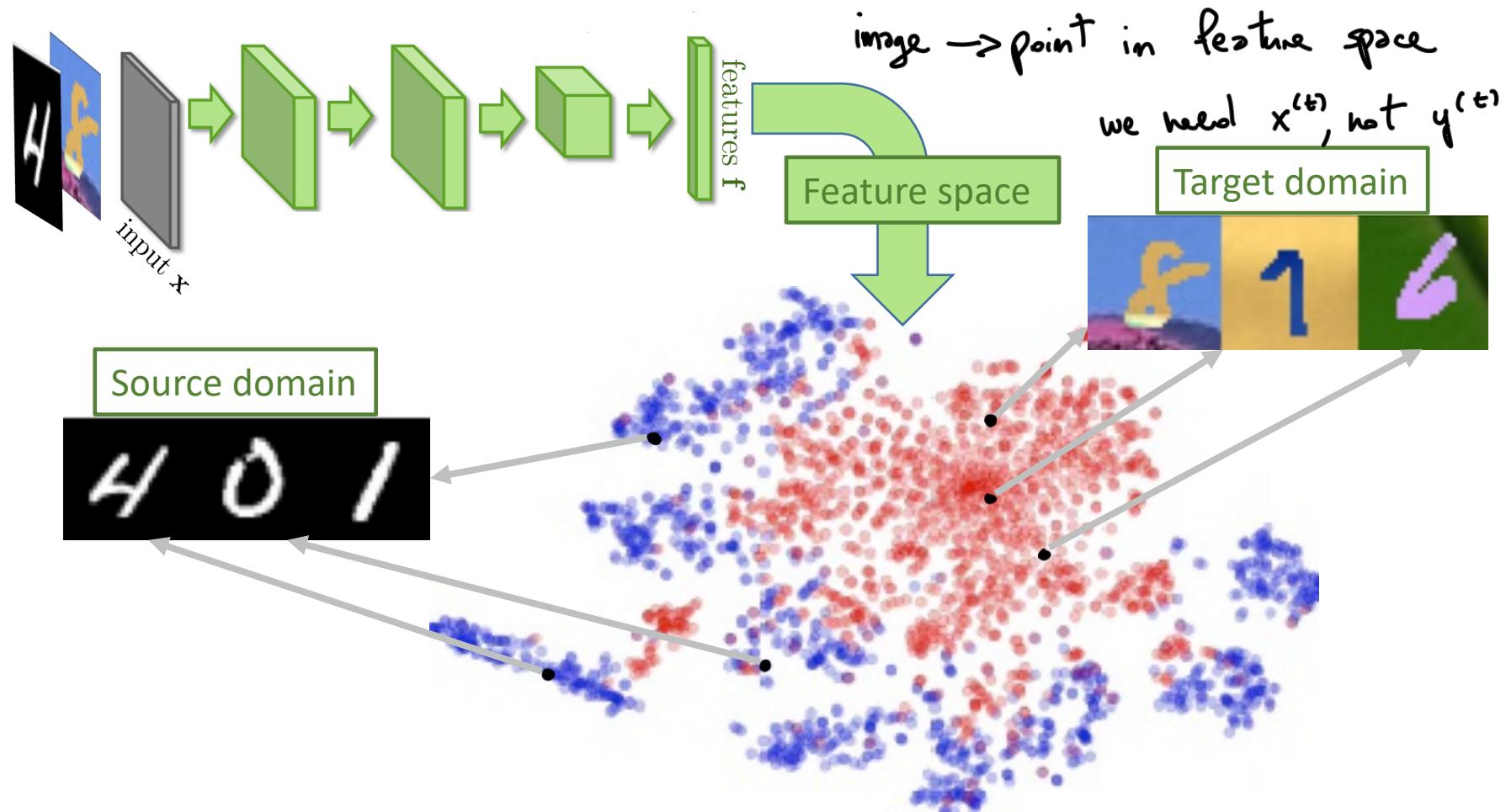
# Unsupervised Domain Adaptation (UDA)



MNIST-M

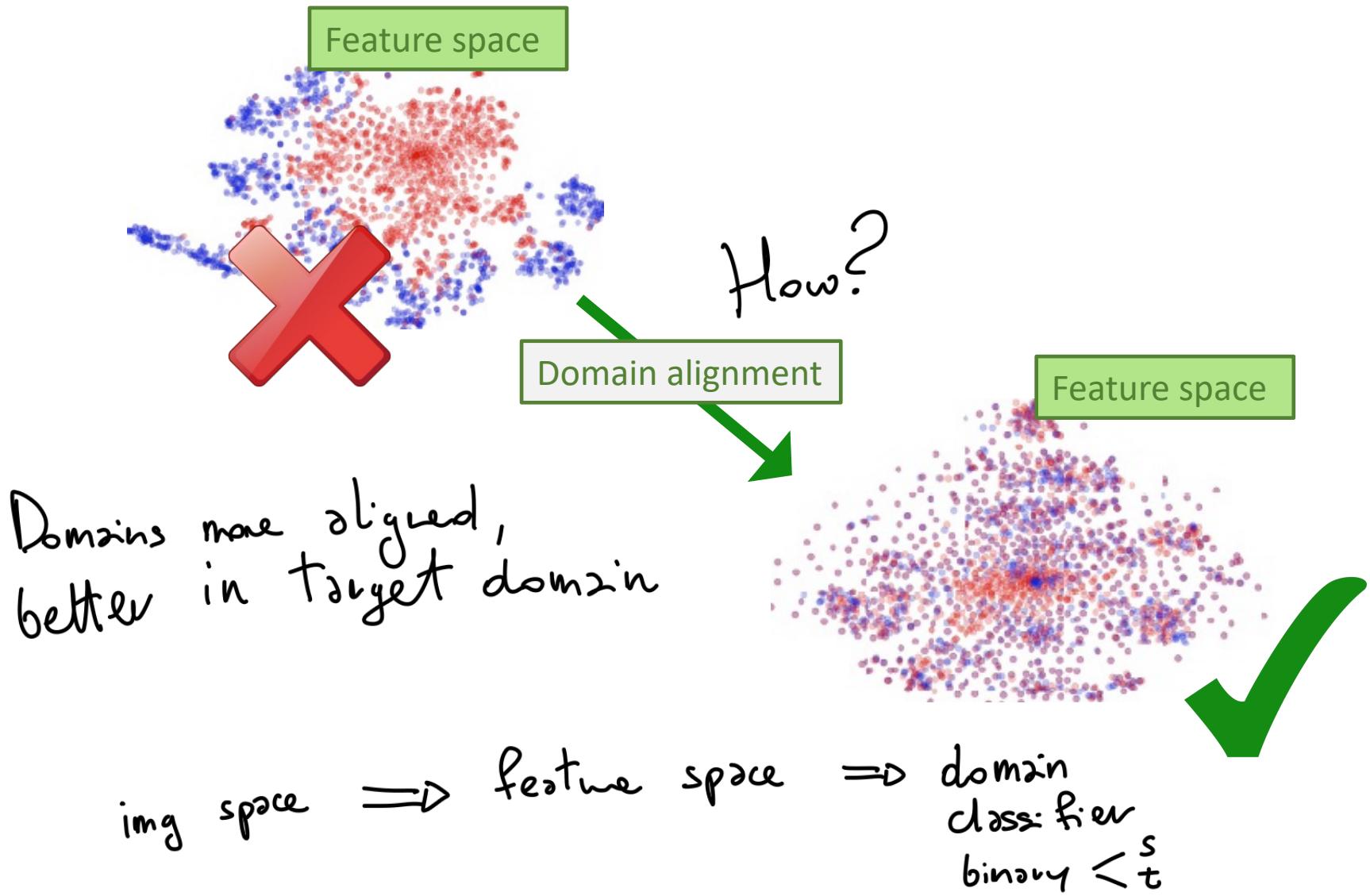
Final test on target domain! (but same task 0..9)

# Unsupervised Domain adaptation (UDA): objectives



Main principle: diminish the **domain** shift in the learned features, encourage domain confusion

# UDA strategy: align both domains



# UDA strategy: 1/ domain-adversarial training

Add to the feature generator (G) a domain classifier (discriminant D) for which labels are available!

Learn G and D:

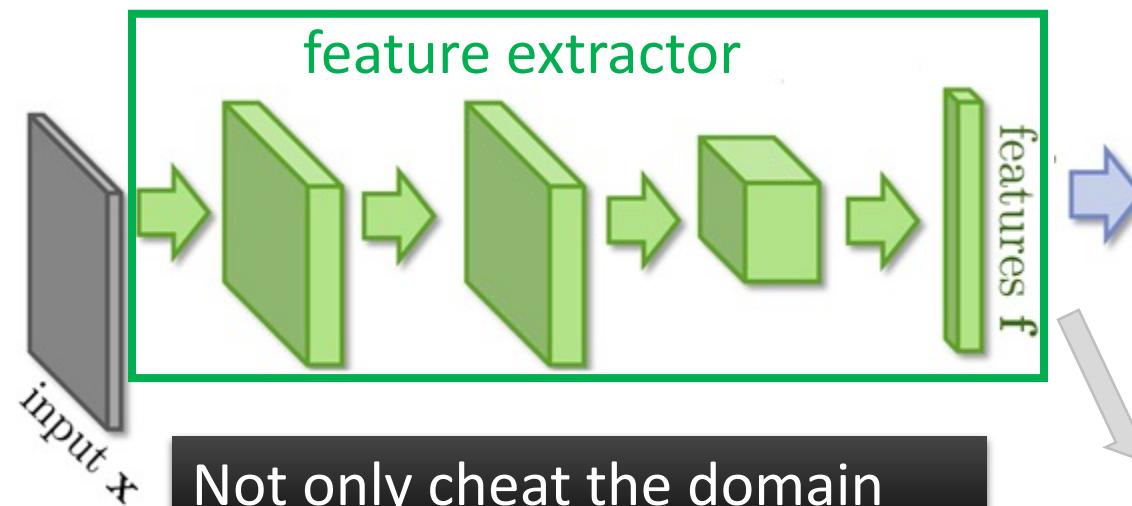
G tries to align domains

D tries to identify domains



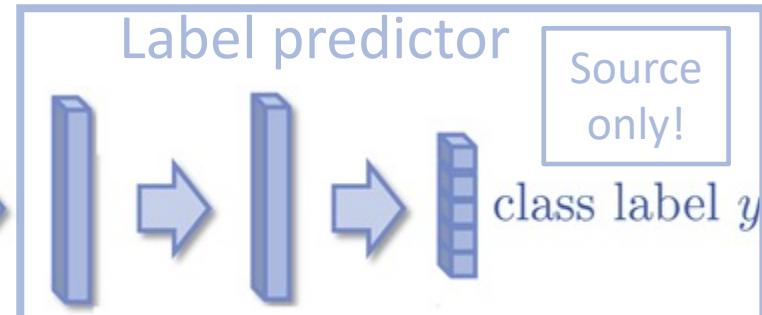
UDA strategy: 1/ domain-adversarial training  
2/ classification task (same for source and target here)

Maximize label classification accuracy + minimize domain classification accuracy



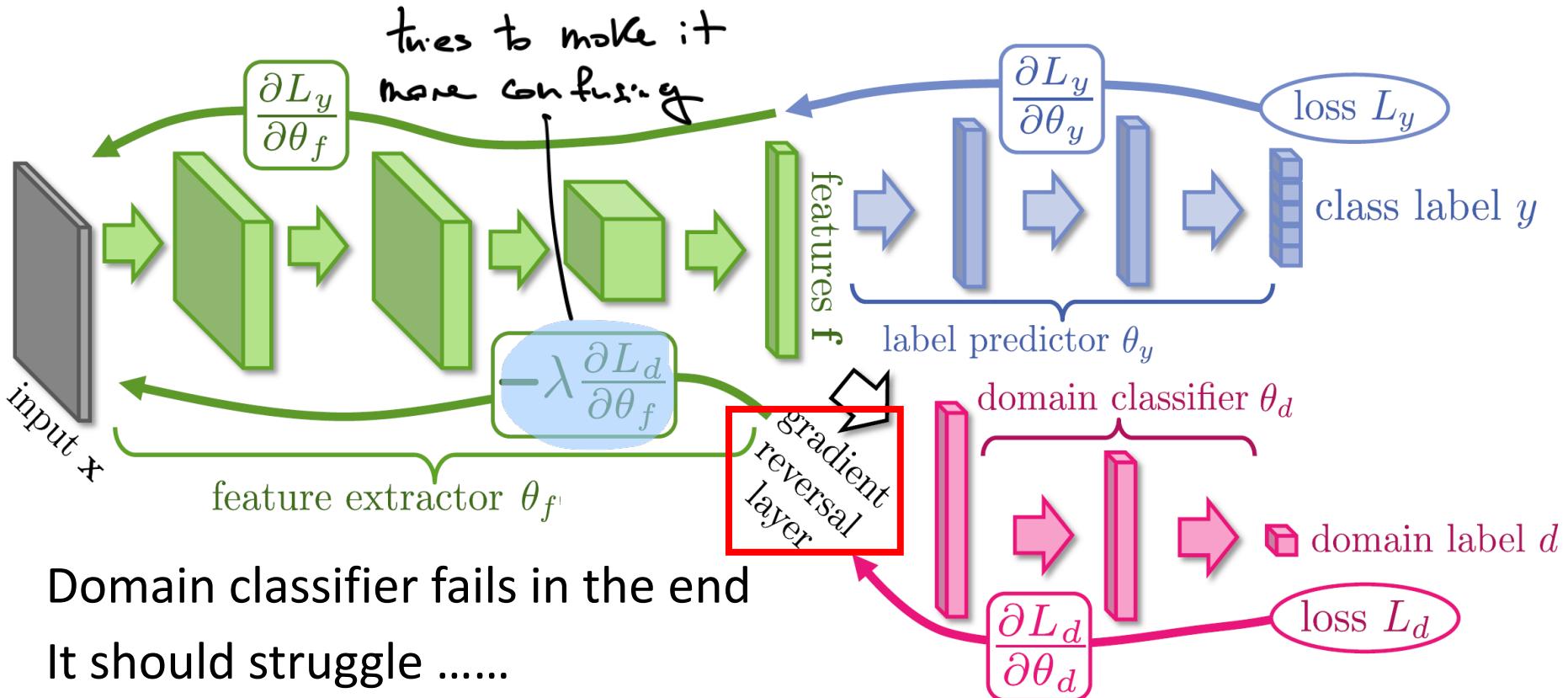
Not only cheat the domain classifier, but satisfying label classifier at the same time

Maximize label classification accuracy



Maximize domain classification accuracy

# UDA strategy: joint learning



Yaroslav Ganin, Victor Lempitsky, Unsupervised Domain Adaptation by Backpropagation, ICML, 2015

Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, Domain-Adversarial Training of Neural Networks, JMLR, 2016

# Domain-adversarial training

	MNIST	SYN NUMBERS	SVHN	SYN SIGNS	
SOURCE					
TARGET					
METHOD	SOURCE	MNIST	SYN NUMBERS	SVHN	SYN SIGNS
	TARGET	MNIST-M	SVHN	MNIST	GTSRB
SOURCE ONLY	<i>(no tuning on target)</i>	.5749	.8665	.5919	.7400
SA (FERNANDO ET AL., 2013)		.6078 (7.9%)	.8672 (1.3%)	.6157 (5.9%)	.7635 (9.1%)
PROPOSED APPROACH		<b>.8149</b> (57.9%)	<b>.9048</b> (66.1%)	<b>.7107</b> (29.3%)	<b>.8866</b> (56.7%)
TRAIN ON TARGET		.9891	.9244	.9951	.9987

Yaroslav Ganin, Victor Lempitsky, Unsupervised Domain Adaptation by Backpropagation,  
ICML, 2015

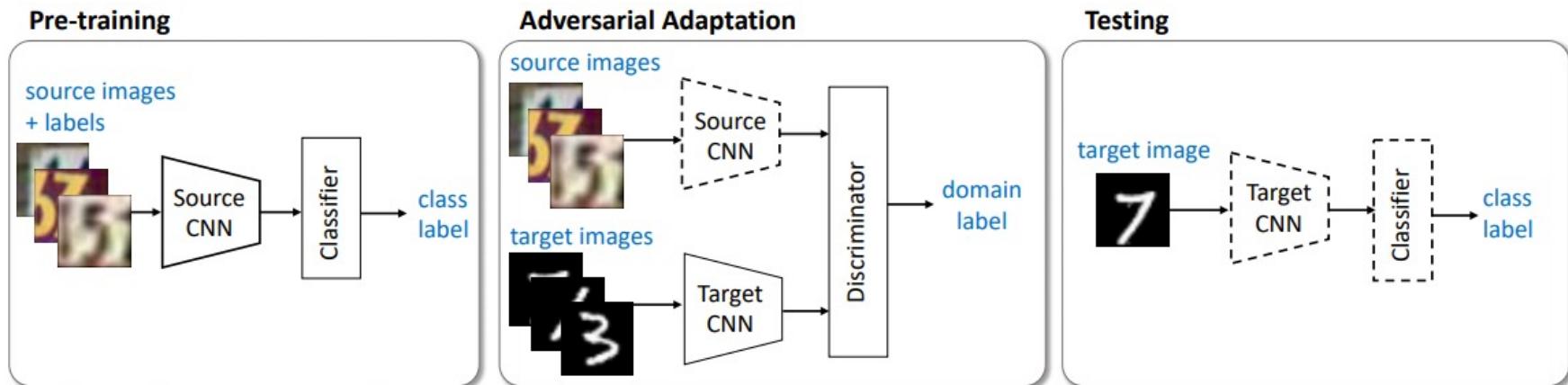
Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand,  
Domain-Adversarial Training of Neural Networks, JMLR, 2016

# Domain adaptation

skip

Main principle: diminish the domain shift in the learned features, encourage domain confusion

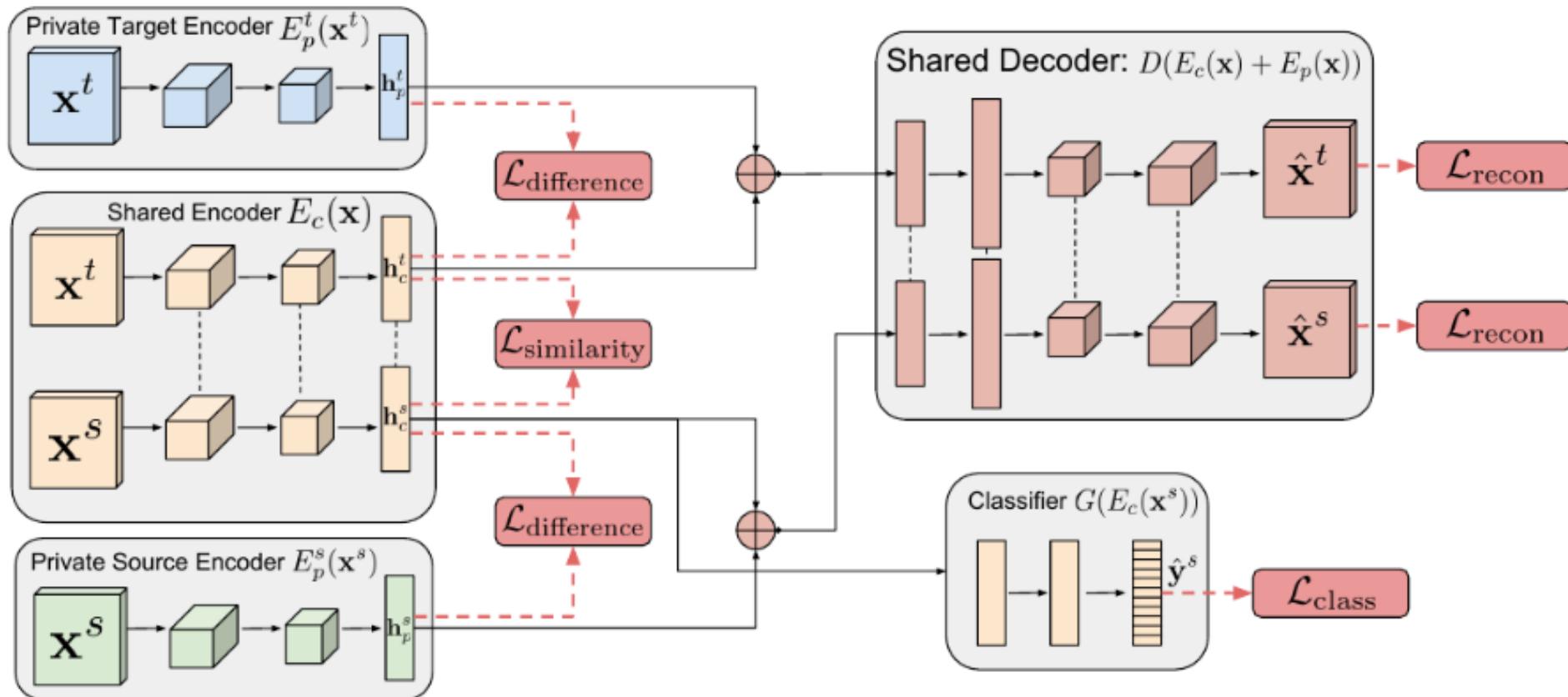
Another example: Adversarial Discriminative Domain Adaptation [Tzeng et al. 2017]



# Domain adaptation

skip

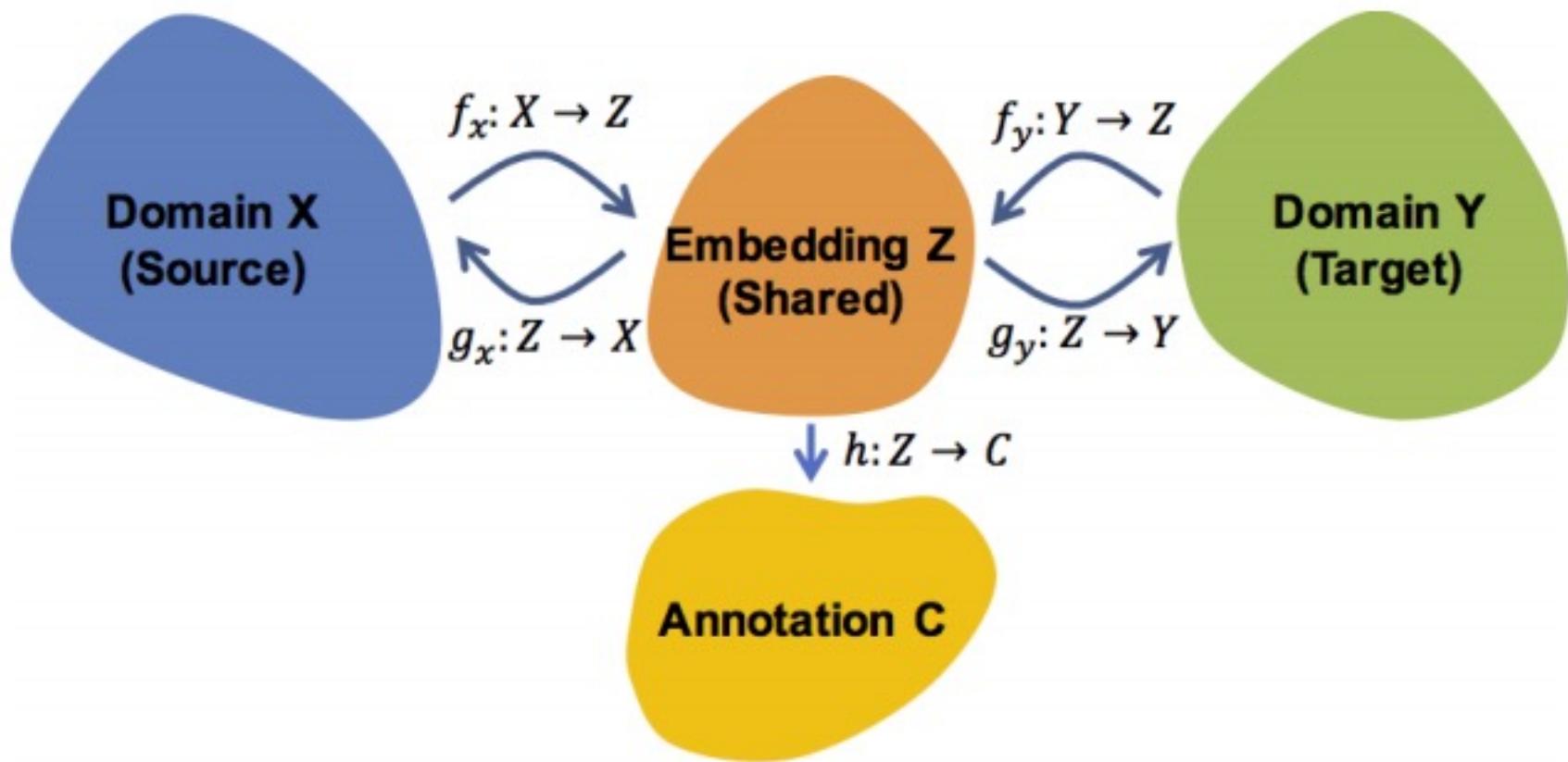
## Other architecture



# Domain adaptation

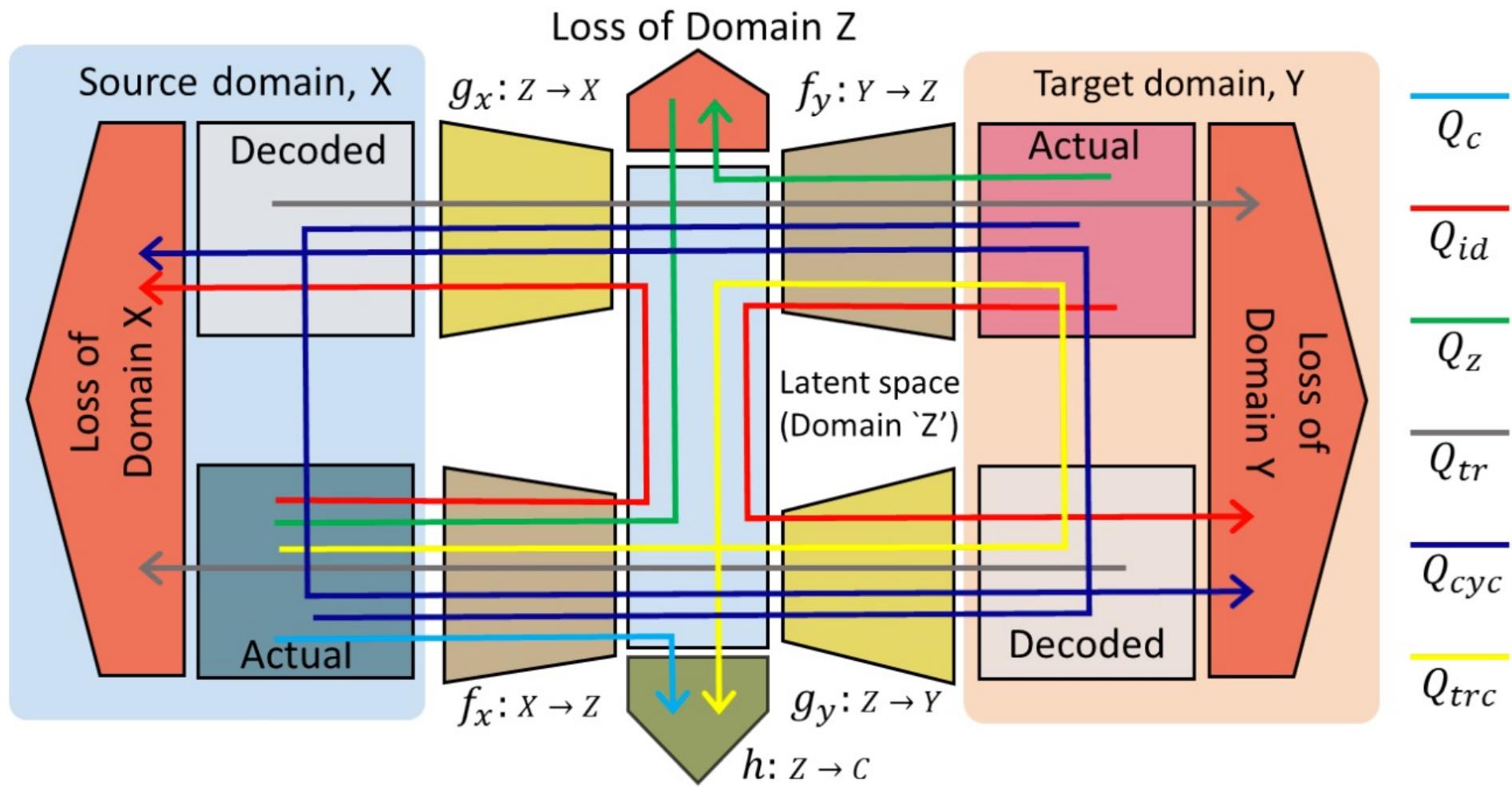
Simple design  
2 domains, 1 embedding

Other architecture: Image translation for Domain adaptation [Murez 2017]



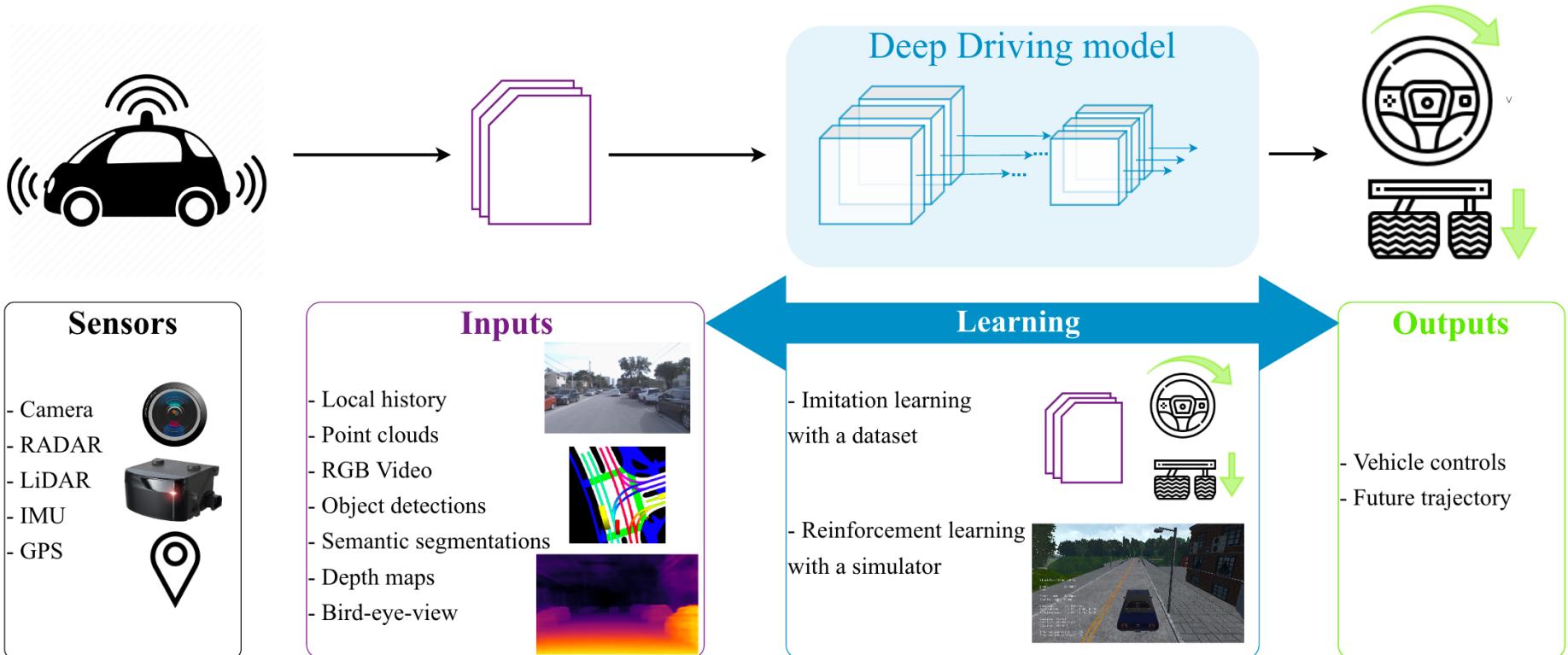
# Domain adaptation

Other architecture: Image translation for Domain adaptation [Murez 2017]



Use-Case: Domain adaptation for  
Autonomous driving

# Context: Neural network-based autonomous driving system framework



# Challenges for perception

## Multi-sensor perception

- Sensor fusion; Camera, radar and Lidar

## 3D dynamic understanding

- 3D object detection; Motion forecast; Intention prediction

## Frugal learning

- Training with limited data or supervision; Domain adaptation

## Reliability

- Robustness; Uncertainty estimation; Failure prediction

## Explainability

- Decision interpretation; Post-hoc or by-design

# Domain gap

Different, though *related* input data distributions

Source domain → Target domain



huge number of examples  
like this → large gap

- Different weather, light, location, sensor's spec/setup

# Domain gap

Different, though *related* input data distributions

Source domain → Target domain



- Different weather, light, location, sensor's spec/setup

# Domain gap

Different, though *related* input data distributions

Source domain → Target domain



- Different weather, light, location, sensor's spec/setup

# Domain gap

Different, though *related* input data distributions

Source domain → Target domain



- Different weather, light, location, sensor's spec/setup

# Domain gap

Different, though *related* input data distributions

Source domain → Target domain



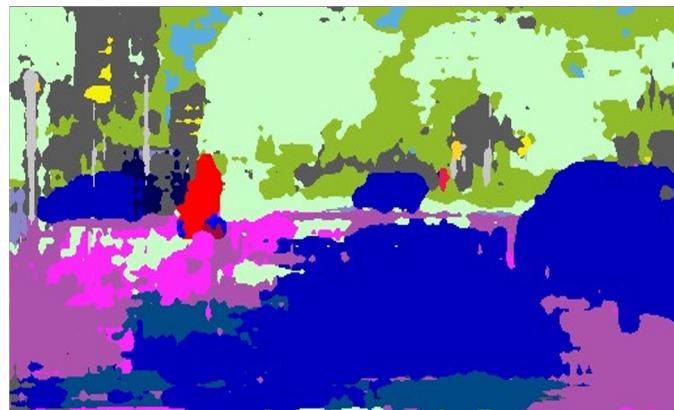
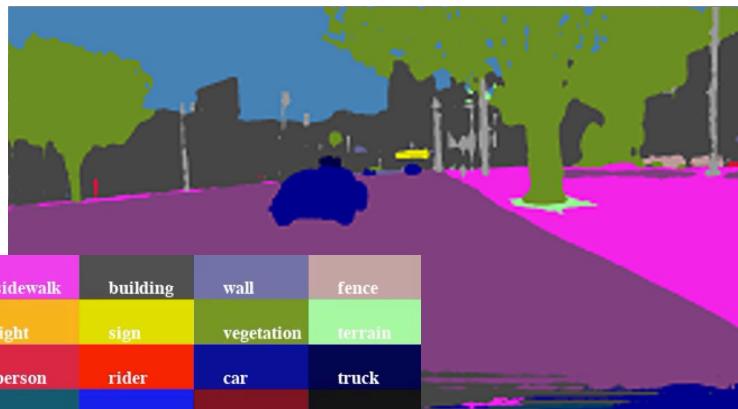
train on a simulator → infer in real world

- Synthetic vs. real

# Domain gap

Different, though *related* input data distributions

Source domain → Target domain



- Synthetic vs. real

# Unsupervised Domain Adaptation (UDA)

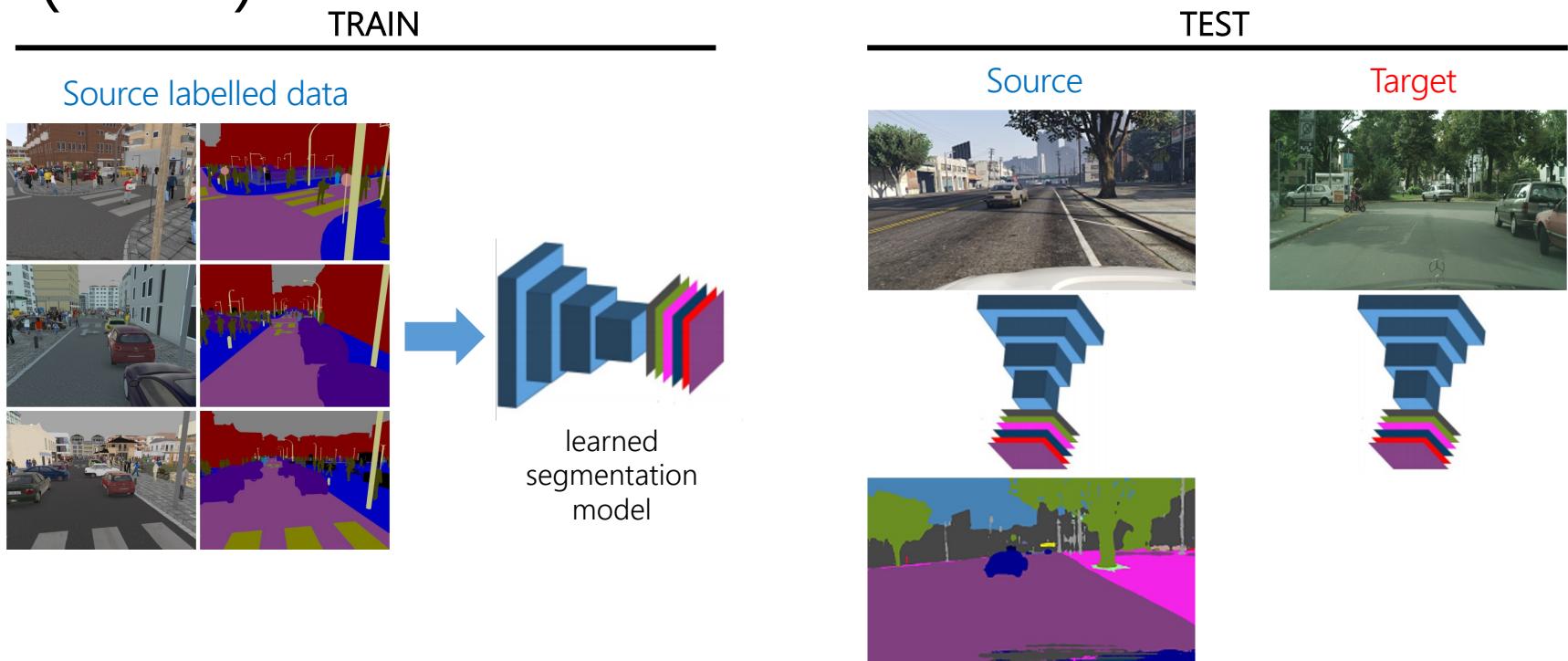
Labelled source domain data



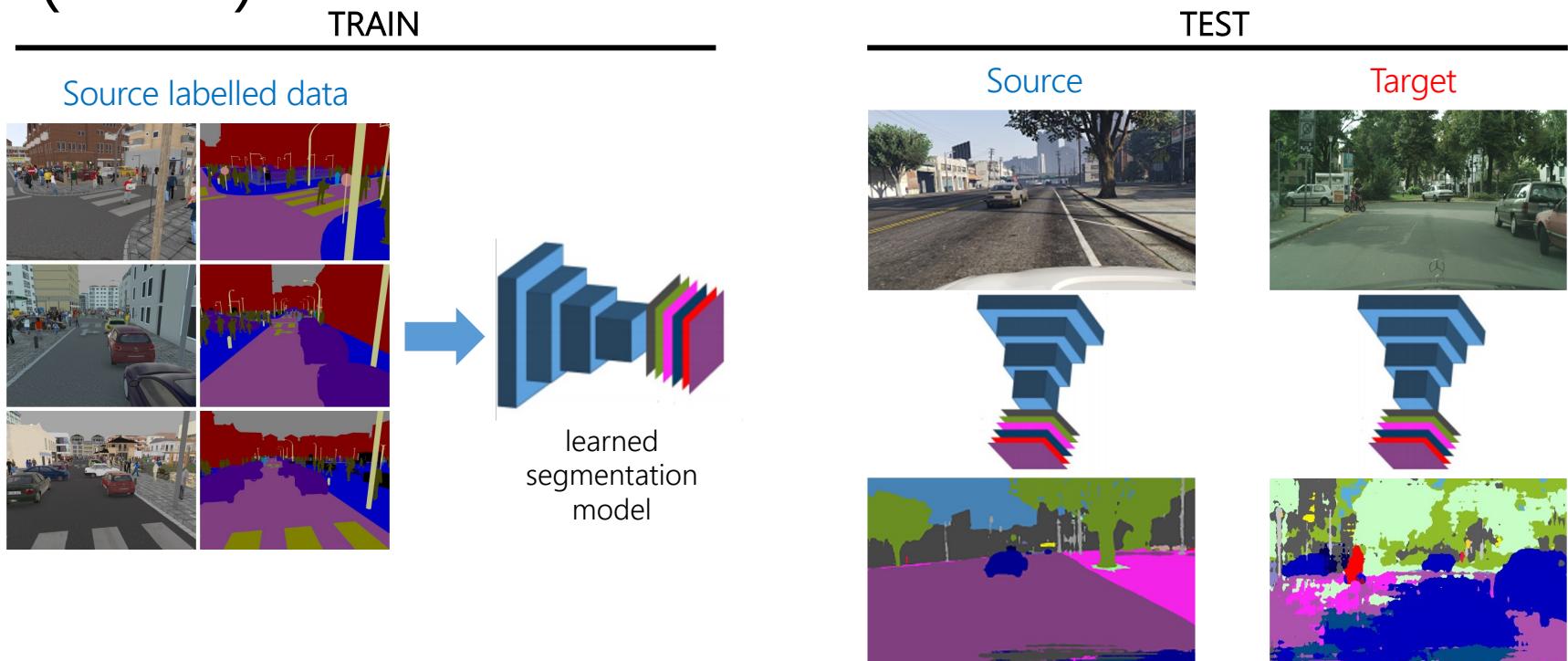
Unlabelled target domain data



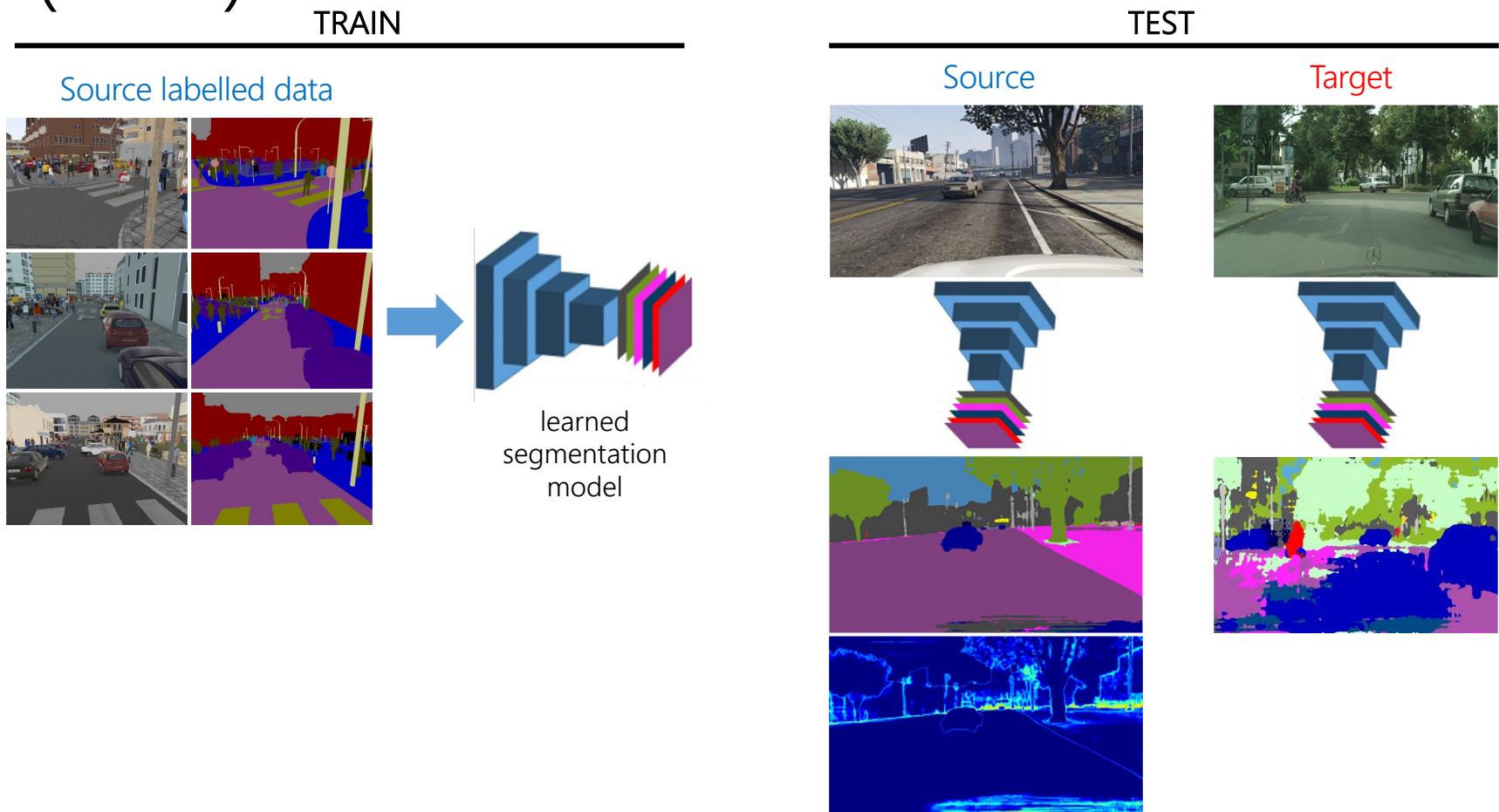
# Unsupervised Domain Adaptation (UDA)



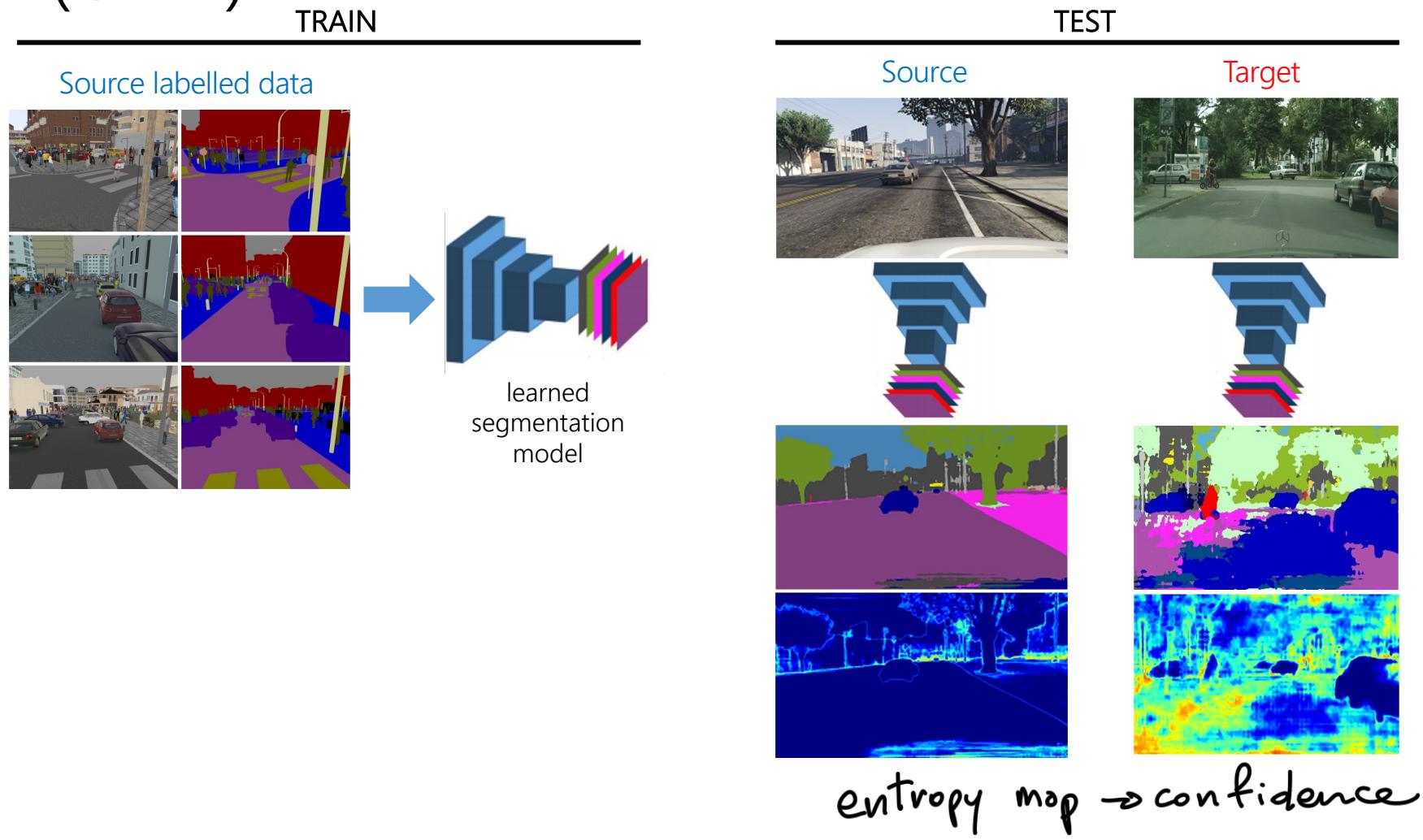
# Unsupervised Domain Adaptation (UDA)



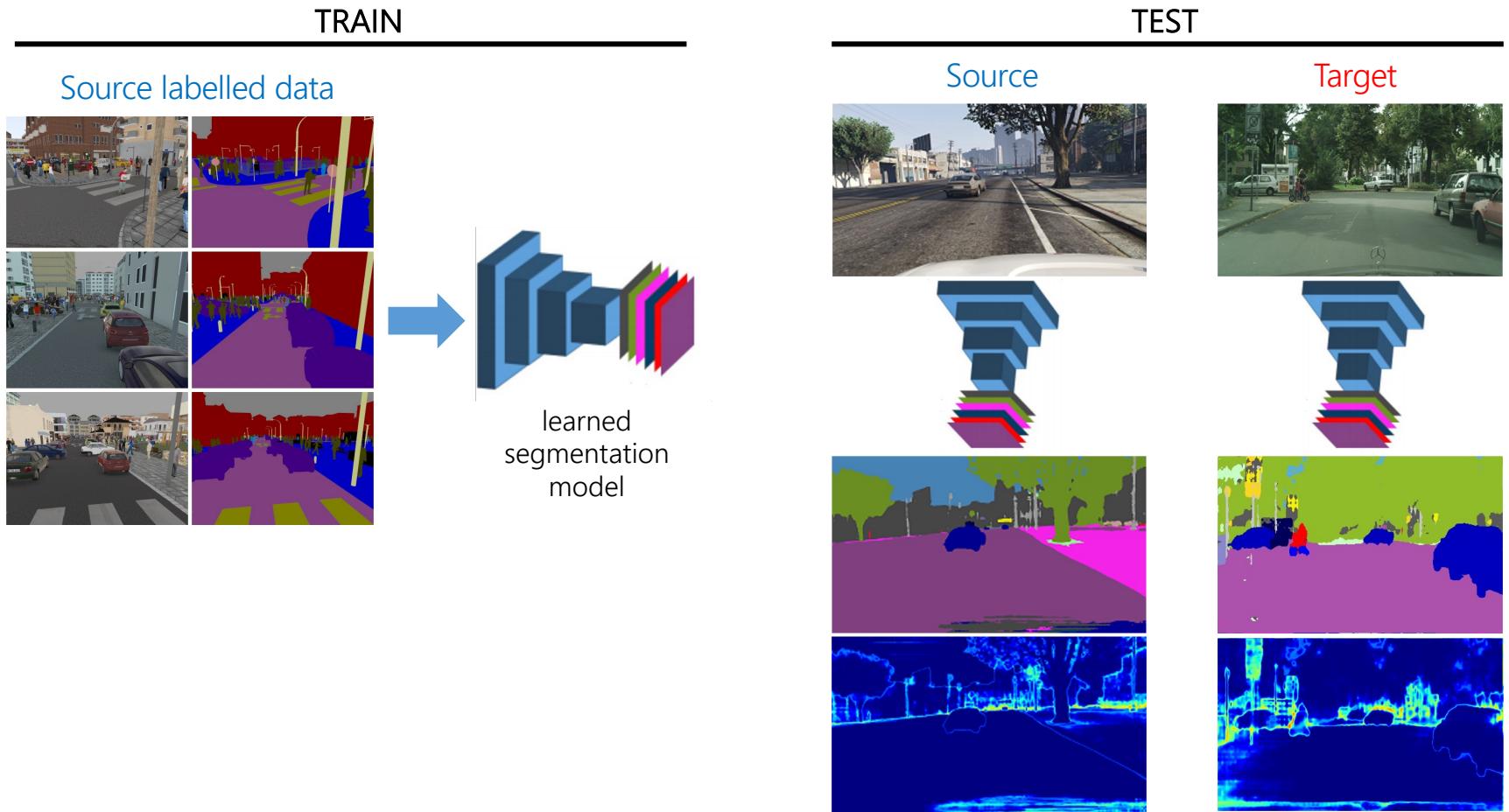
# Unsupervised Domain Adaptation (UDA)



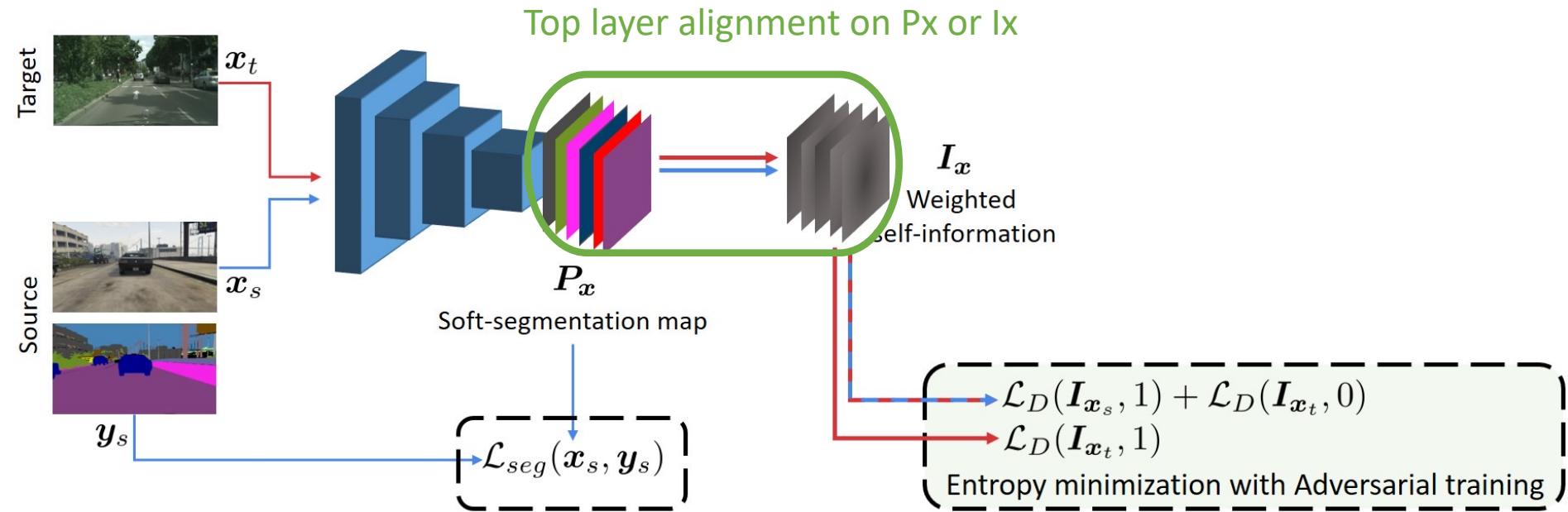
# Unsupervised Domain Adaptation (UDA)



# Expected results with UDA training



# Unsupervised Domain Adaptation (UDA)

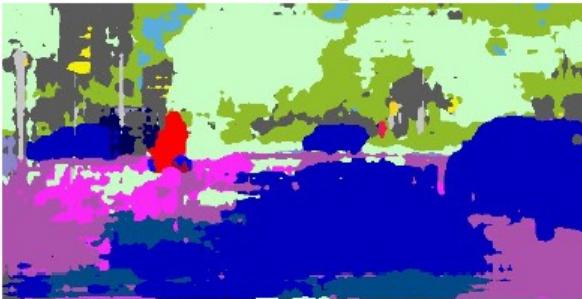


# Qualitative results

input image



without UDA

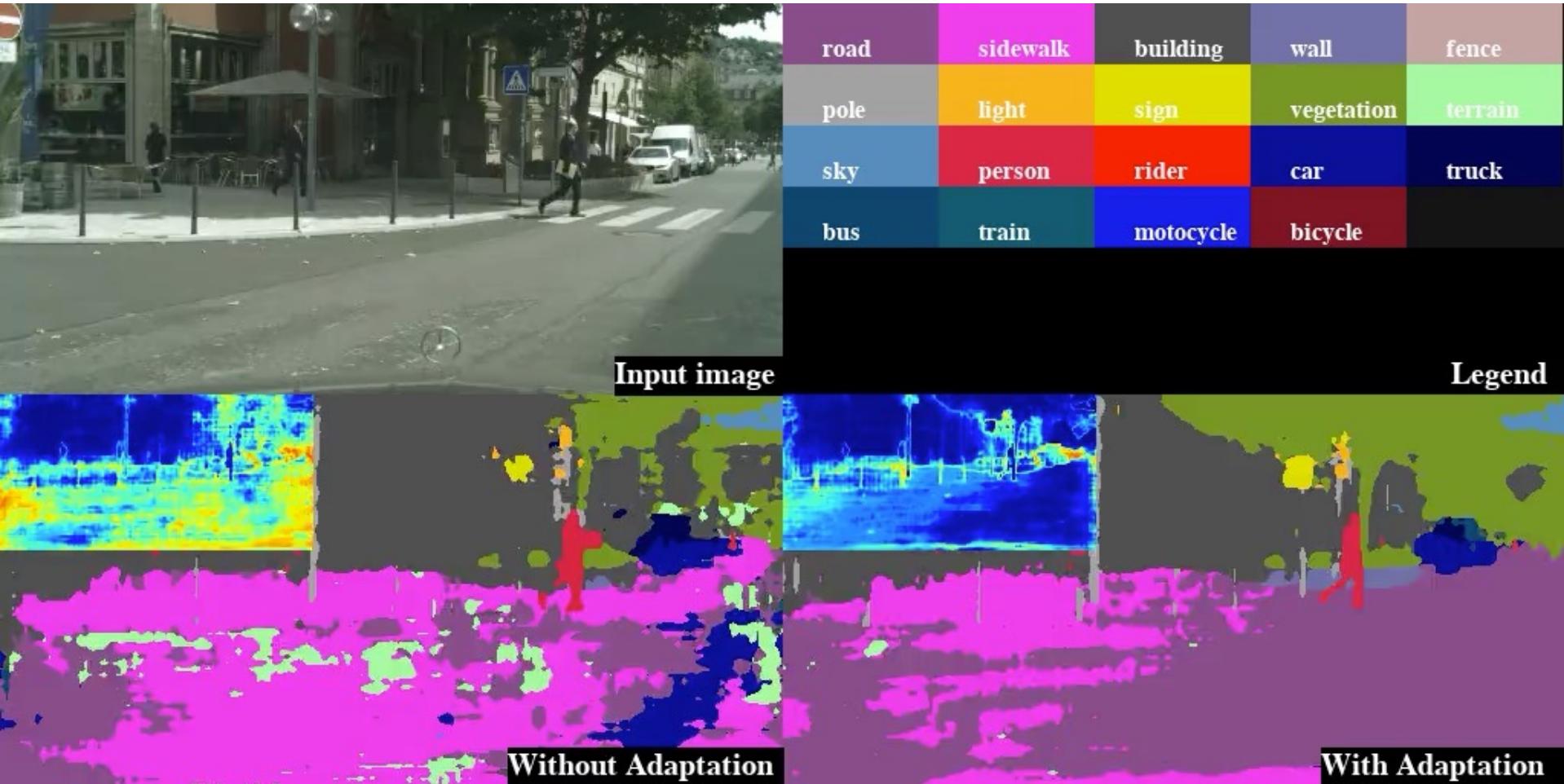


with UDA



road	sidewalk	building	wall	fence
pole	light	sign	vegetation	terrain
sky	person	rider	car	truck
bus	train	motorcycle	bicycle	

# UDA Results (with Adversarial Entropy)



# Extension: Zero shot + Domain adaptation



Private target classes: **tuk-tuk**, **animal**. Some shared classes: **truck**, **road**, **side walk**, **car**, **person**, **motorbike**, **tree**, **building**.

# Transfer Learning - Overview

		Source Data (not directly related to the task)	
		labelled	unlabeled
Target Data	labelled	Fine-tuning <i>Multitask Learning</i>	Not considered here
	unlabeled	Domain adaptation- adversarial training  <i>Zero-shot learning</i>	Not considered here

# Zero-shot Learning

SAME DOMAIN,  
DIFF TASKS

- Source data:  $(x^s, y^s)$  → Training data
- Target data:  $(\emptyset)$  usually same domain

Different tasks

*Training time :*  $x^s:$



$y^s:$  cat



...  
...

+ Class Information

*Test time  $x^t:$*



=> Fish class!

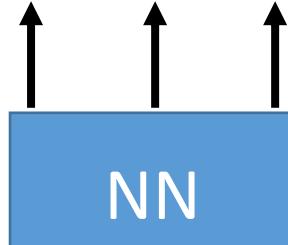
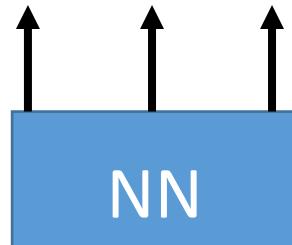
# Zero-shot Learning

- Representing each class by its attributes

## Training

1	0	0
furry	4 legs	tail

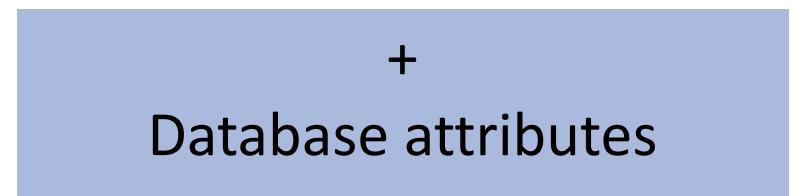
1	1	1
furry	4 legs	tail



class



...



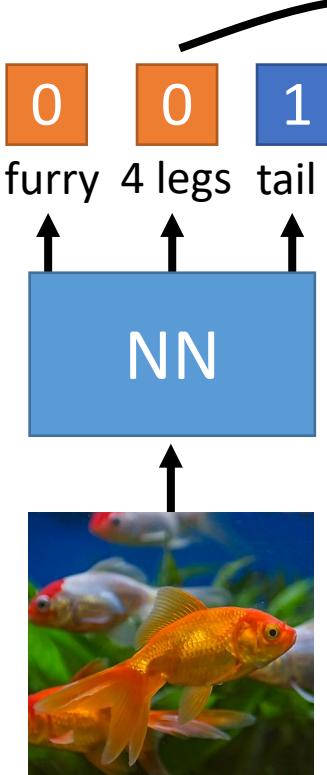
	furry	4 legs	tail	...
Dog	O	O	O	
Fish	X	X	O	
Chimp	O	X	X	
...				

sufficient attributes for one  
to one mapping

# Zero-shot Learning

- Representing each class by its attributes

## Testing



Find the class with the most similar attributes

	furry	4 legs	tail	...
Dog	o	o	o	
Fish	x	x	o	
Chimp	o	x	x	
...				

sufficient attributes for one to one mapping

*zero examples  
of my target (fish)*

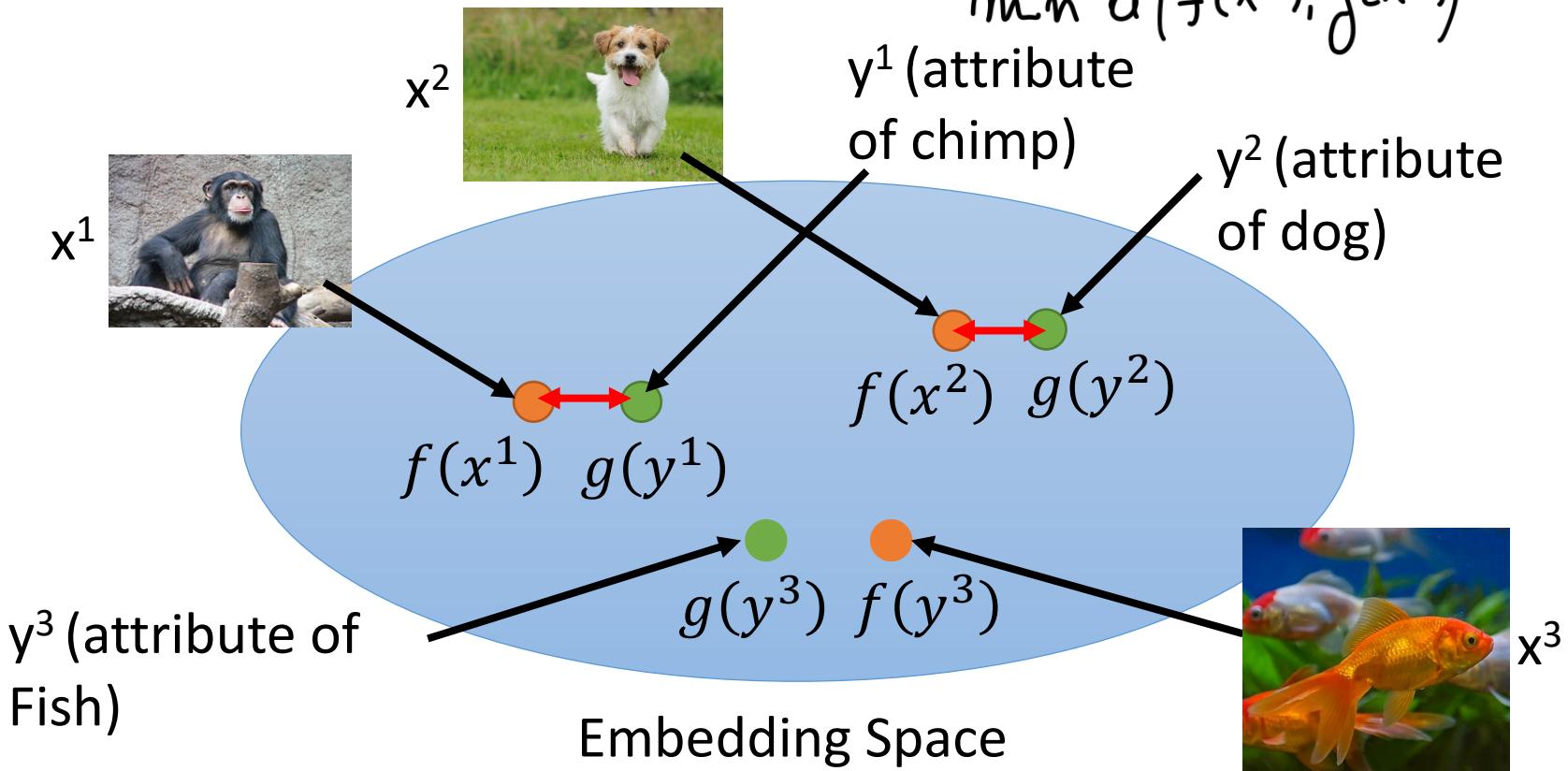
# Zero-shot Learning

What if we don't have attribute database

- Attribute embedding + class (word name) embedding

# Zero-shot Learning

- Attribute embedding



$f(*)$  and  $g(*)$  can be NN.  
Training target:

$f(x^n)$  and  $g(y^n)$  as close as possible

$$\min d(f(x^n), g(x^n))$$

$y^1$  (attribute of chimp)

$y^2$  (attribute of dog)

$$f(x^2) \quad g(y^2)$$

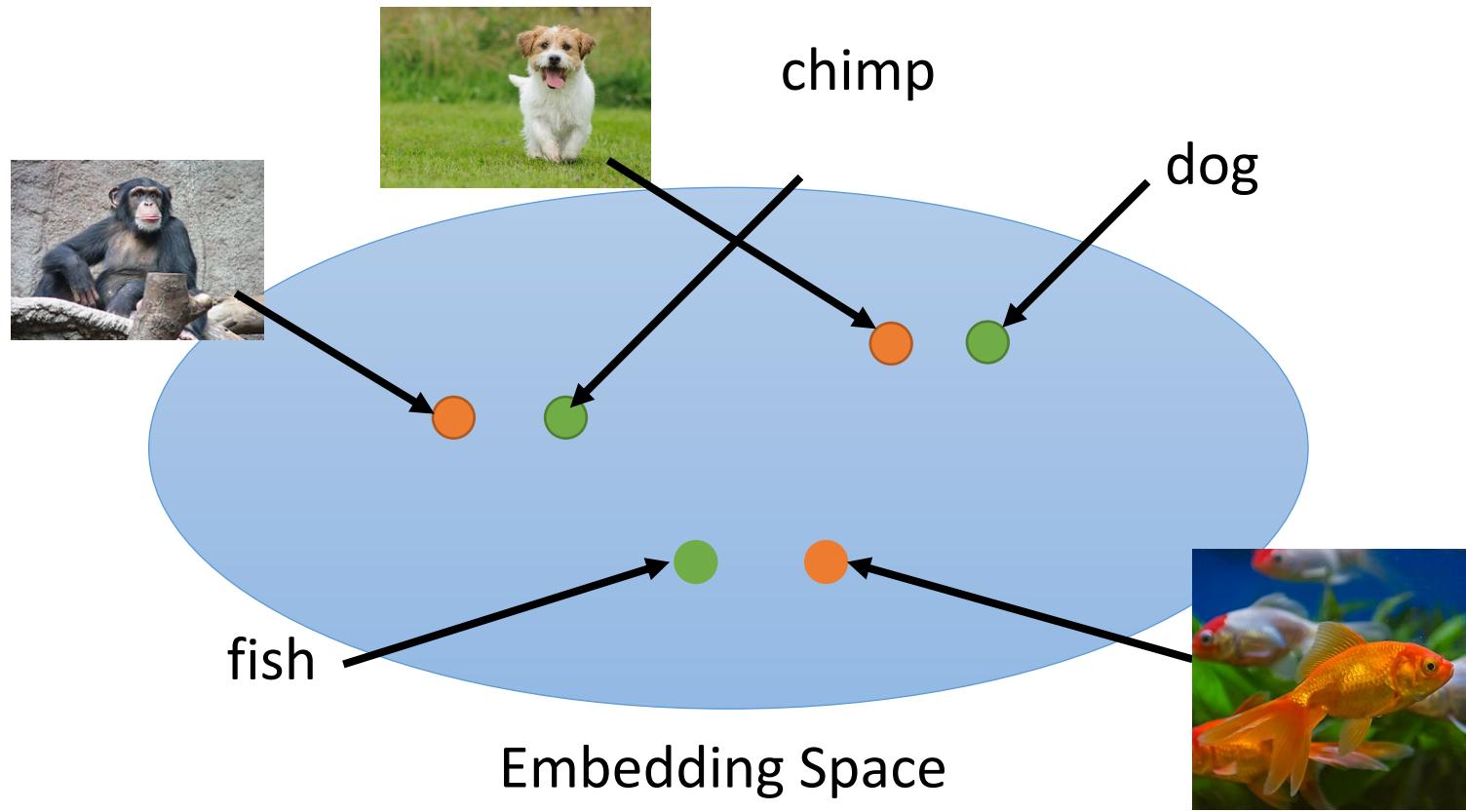
$$f(x^1) \quad g(y^1)$$

$y^3$  (attribute of Fish)

Embedding Space

$y^i$  are linked together by a class relationship (e.g. class name embedding as W2v)

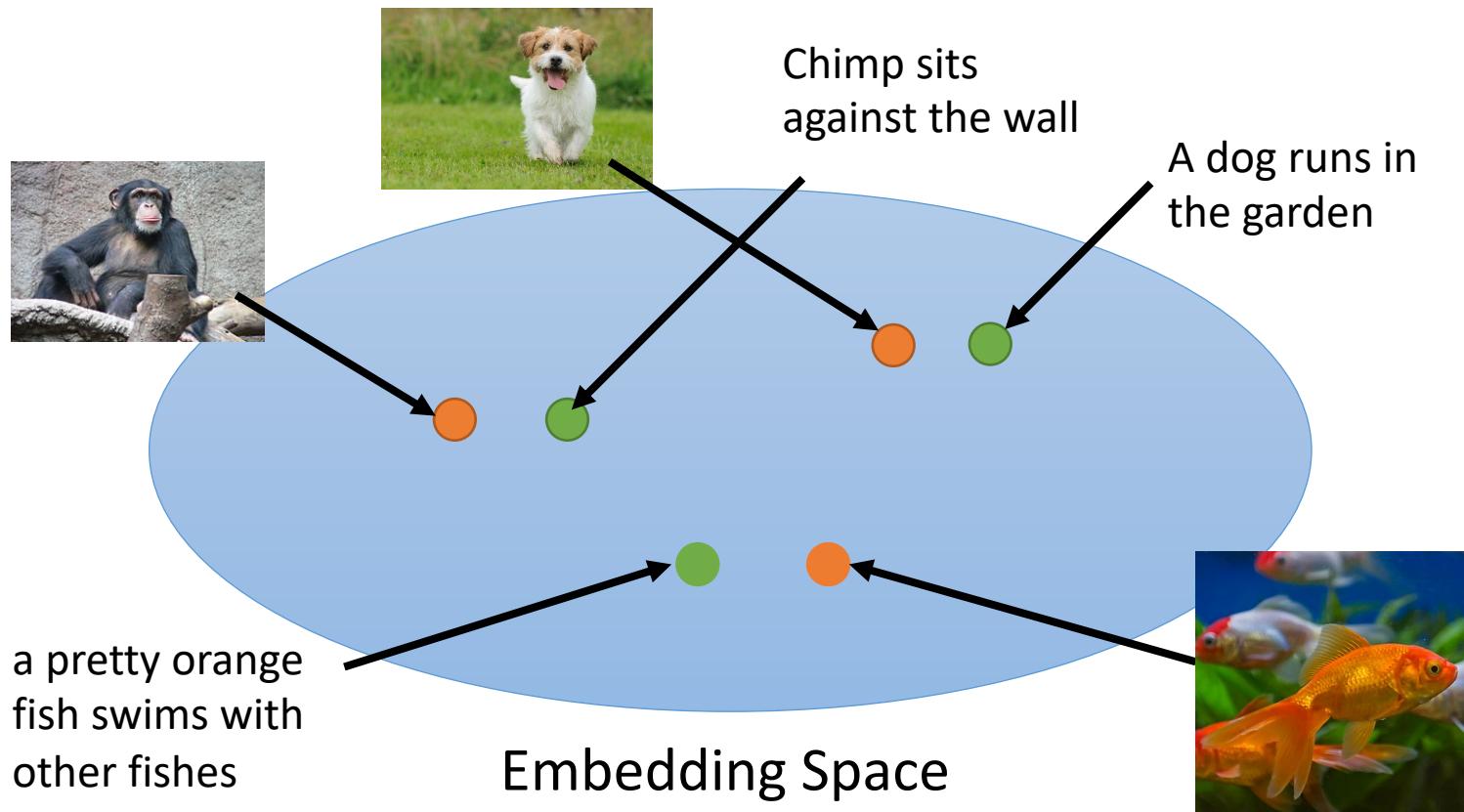
# More on Vision-Language



First time ever alignment in > latent space between imgs and classes

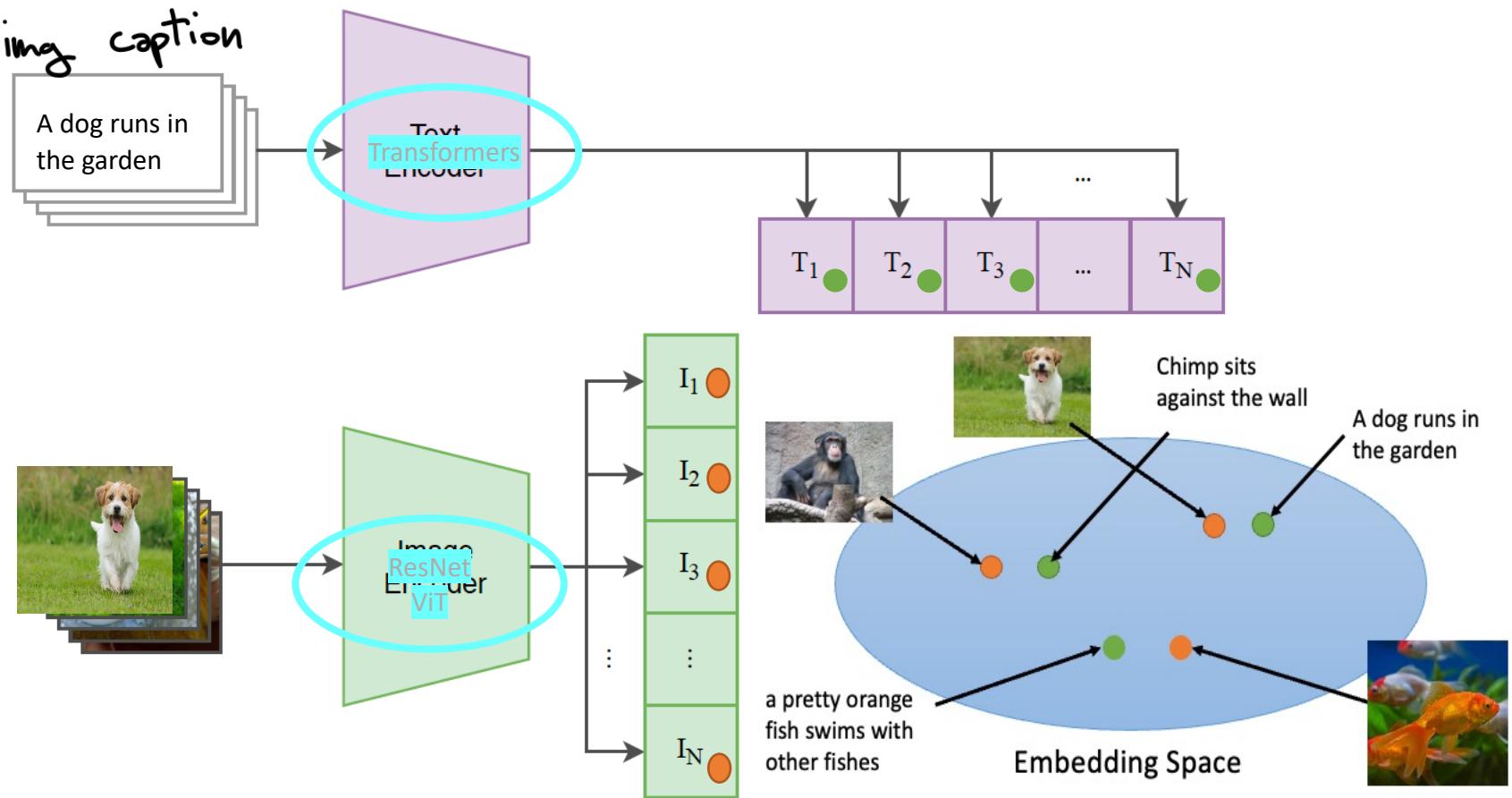
# CLIP: Vision + Language Models (VLM)

[Learning transferable visual models from natural language supervision. Radford/Sutskever ICML, 2021]



# CLIP: Vision + Language Models (VLM)

Dual architecture: Text encoder + Image encoder



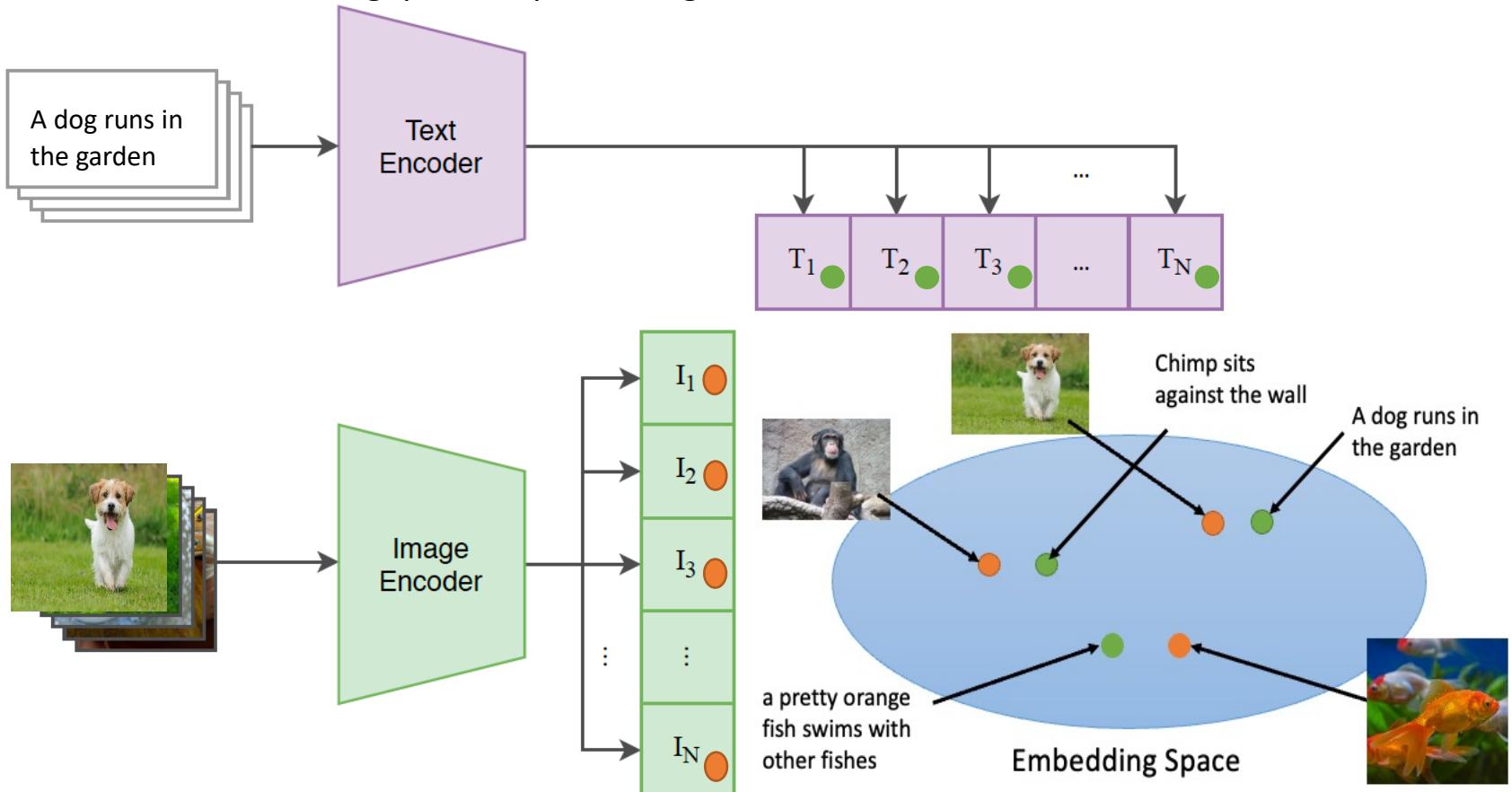
# CLIP: Vision + Language Models (VLM)

## Learning strategy

Training set:  $A = \{(\mathbf{I}_n, \mathbf{T}_n)\}_n$  of image/caption pairs (coherent!)

Massive Text+Image =**400M pairs** to train the model (from the Internet)

Contrastive loss for training: positive pair vs negative one or set



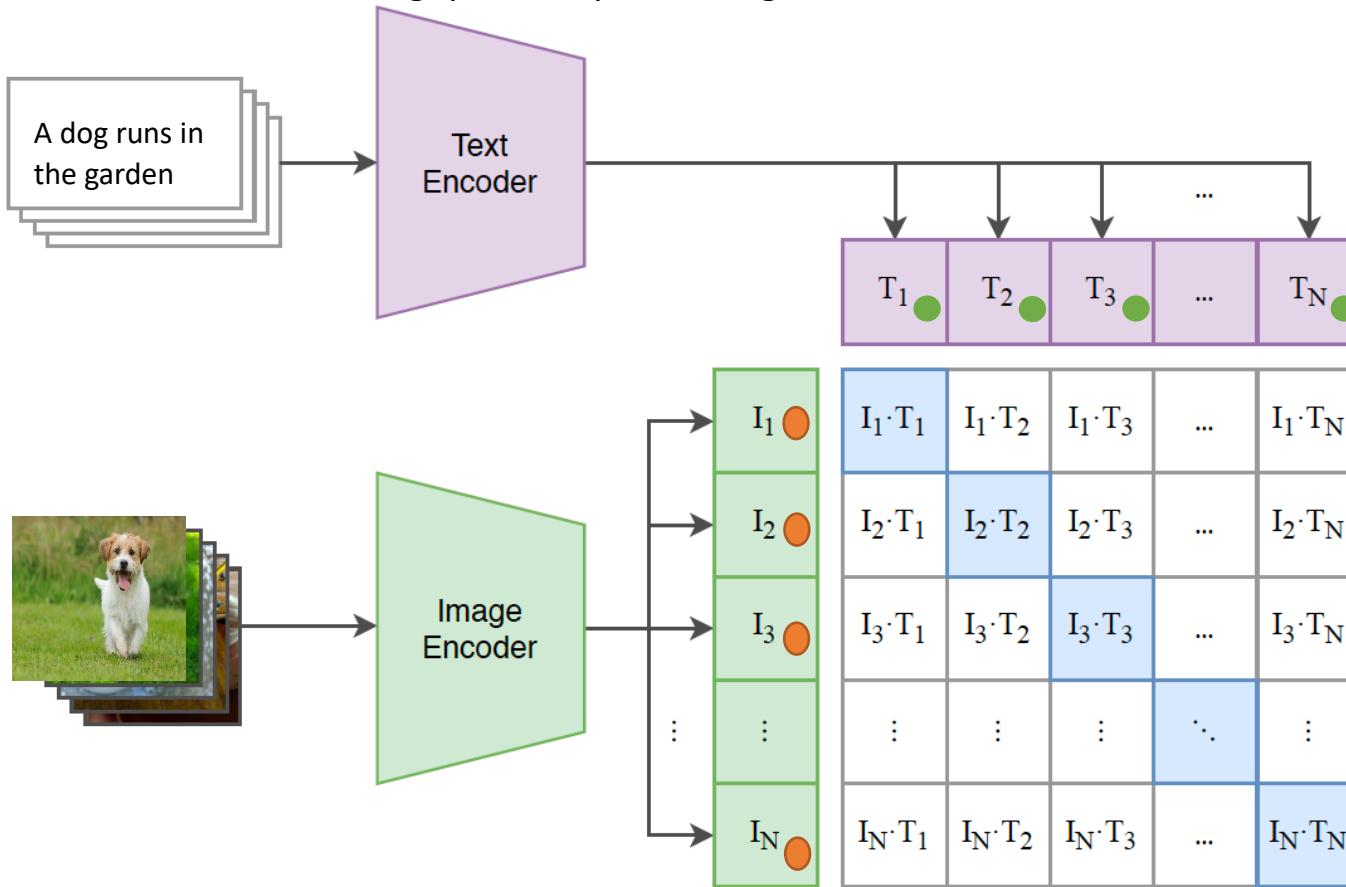
# CLIP: Vision + Language Models (VLM)

## Learning strategy

Training set:  $A = \{(\mathbf{I}_n, \mathbf{T}_n)\}_n$  of image/caption pairs (coherent!)

Massive Text+Image =**400M pairs** to train the model (from the Internet)

Contrastive loss for training: positive pair vs negative one or set



# CLIP: Vision + Language Models (VLM)

## Learning strategy

Training set:  $A = \{(\mathbf{I}_n, \mathbf{T}_n)\}_n$  of image/caption pairs (coherent!)

Massive Text+Image =**400M pairs** to train the model (from the Internet)

Contrastive loss for training: positive pair vs negative pair or more

(contrastive) Triplet loss: A variant of the standard margin based loss (SVM)

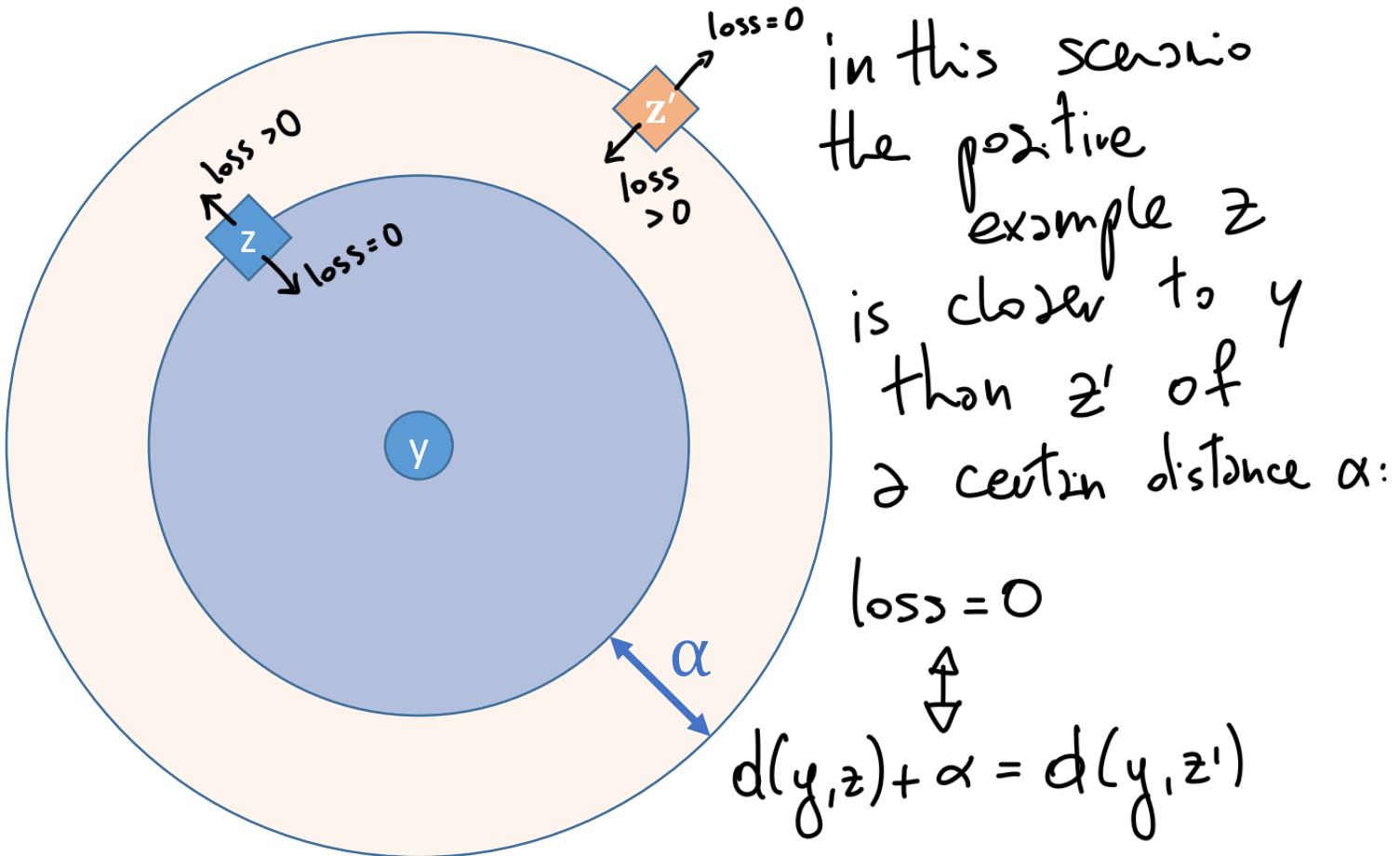
- Triplet  $(\mathbf{y}, \mathbf{z}, \mathbf{z}')$
- Anchor:  $\mathbf{y}$  (E.g image representation)
- Positive:  $\mathbf{z}$  (E.g associated caption representation)
- Negative:  $\mathbf{z}'$  (E.g contrastive caption representation) *push them far away*
- Margin parameter  $\alpha$

$$\text{loss}(\mathbf{y}, \mathbf{z}, \mathbf{z}') = \max\{0, \alpha - \langle \mathbf{y}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z}' \rangle\}$$

*push them close*      *far away*  
CONTRASTIVE APPROACH

## Learning strategy: triplet loss

$$\text{loss}(\mathbf{y}, \mathbf{z}, \mathbf{z}') = \max\{0, \alpha + d(\mathbf{y}, \mathbf{z}) - d(\mathbf{y}, \mathbf{z}')\}$$



# Learning strategy: triplet loss

**Hard negative** margin-based loss:

Loss for a **batch**  $\mathcal{B} = \{(\mathbf{I}_n, \mathbf{T}_n)\}_{n \in B}$  of image/sentence pairs:

$$\mathcal{L}(\Theta; \mathcal{B}) = \frac{1}{|B|} \sum_{n \in B} \left( \max_{m \in C_n \cap B} \text{loss}(\mathbf{x}_n, \mathbf{v}_n, \mathbf{v}_m) + \max_{m \in D_n \cap B} \text{loss}(\mathbf{v}_n, \mathbf{x}_n, \mathbf{x}_m) \right)$$

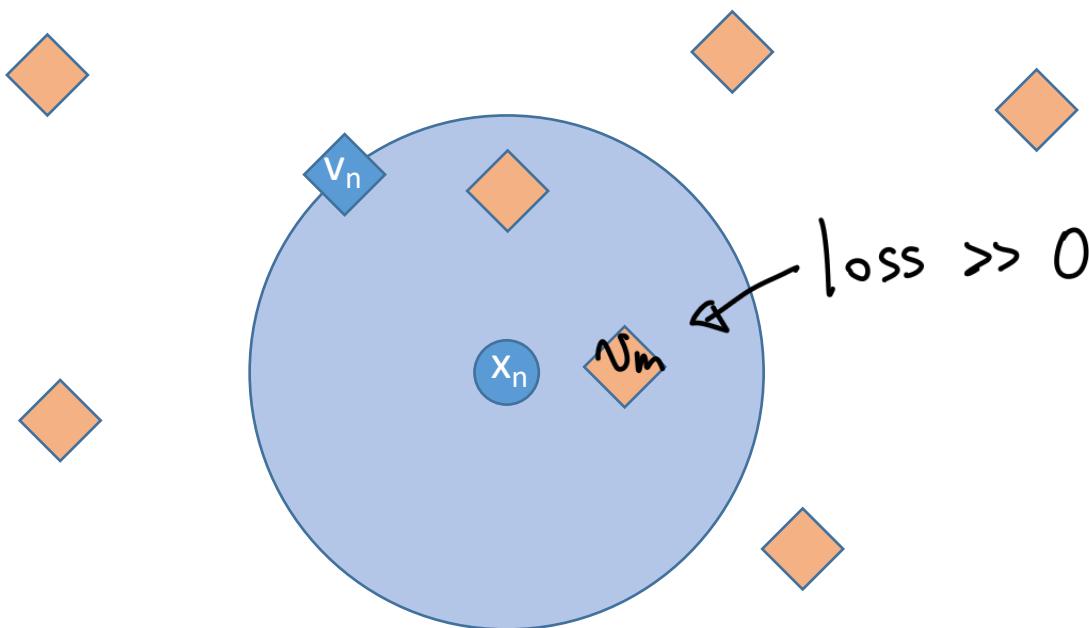
With  $C_n$  (resp.  $D_n$ ) set of indices of caption (resp. image) unrelated to  $n$ -th element

# Learning strategy: hard negative triplet loss

**Mining hard negative contrastive example:**

$$\mathcal{L}(\Theta; \mathcal{B}) = \frac{1}{|\mathcal{B}|} \sum_{n \in \mathcal{B}} \left( \max_{m \in C_n \cap \mathcal{B}} \text{loss}(\mathbf{x}_n, \mathbf{v}_n, \mathbf{v}_m) + \max_{m \in D_n \cap \mathcal{B}} \text{loss}(\mathbf{v}_n, \mathbf{x}_n, \mathbf{x}_m) \right)$$

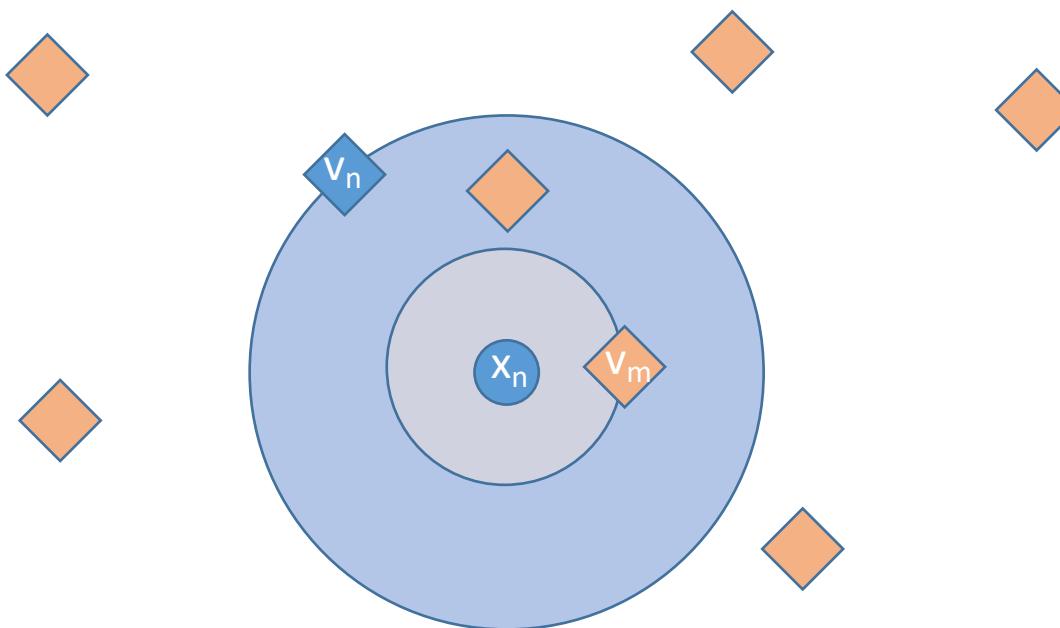
*select the most critical  $\mathbf{v}_m$  that max loss*



# Learning strategy: hard negative triplet loss

**Mining hard negative contrastive example:**

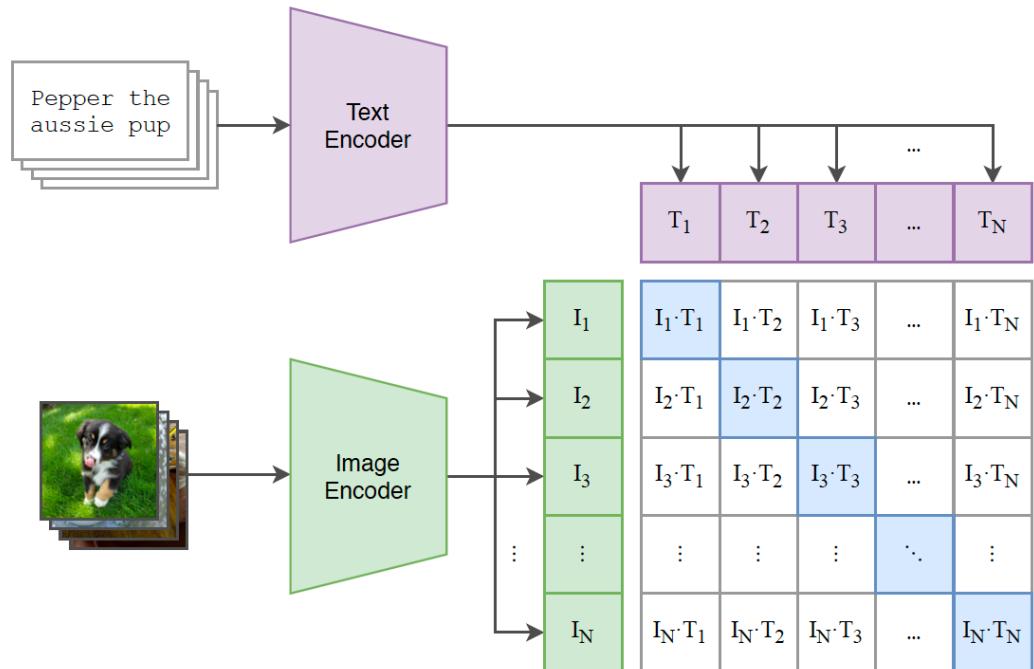
$$\mathcal{L}(\Theta; \mathcal{B}) = \frac{1}{|\mathcal{B}|} \sum_{n \in \mathcal{B}} \left( \max_{m \in C_n \cap \mathcal{B}} \text{loss}(\mathbf{x}_n, \mathbf{v}_n, \mathbf{v}_m) + \max_{m \in D_n \cap \mathcal{B}} \text{loss}(\mathbf{v}_n, \mathbf{x}_n, \mathbf{x}_m) \right)$$



# CLIP: Vision + Language Models (VLM)

Massive Text+Image =**400M**  
**pairs** to train the model (from  
the Internet)

Contrastive loss for training



$$\mathcal{L}_{InfoNCECLIP} = - \sum_i \log \left( \frac{\exp\left(\frac{\text{sim}(I_i, T_i)}{\tau}\right)}{\sum_{k=1}^N \exp\left(\frac{\text{sim}(I_i, T_k)}{\tau}\right)} \right)$$

*dot product =  $I_i \cdot T_i$*

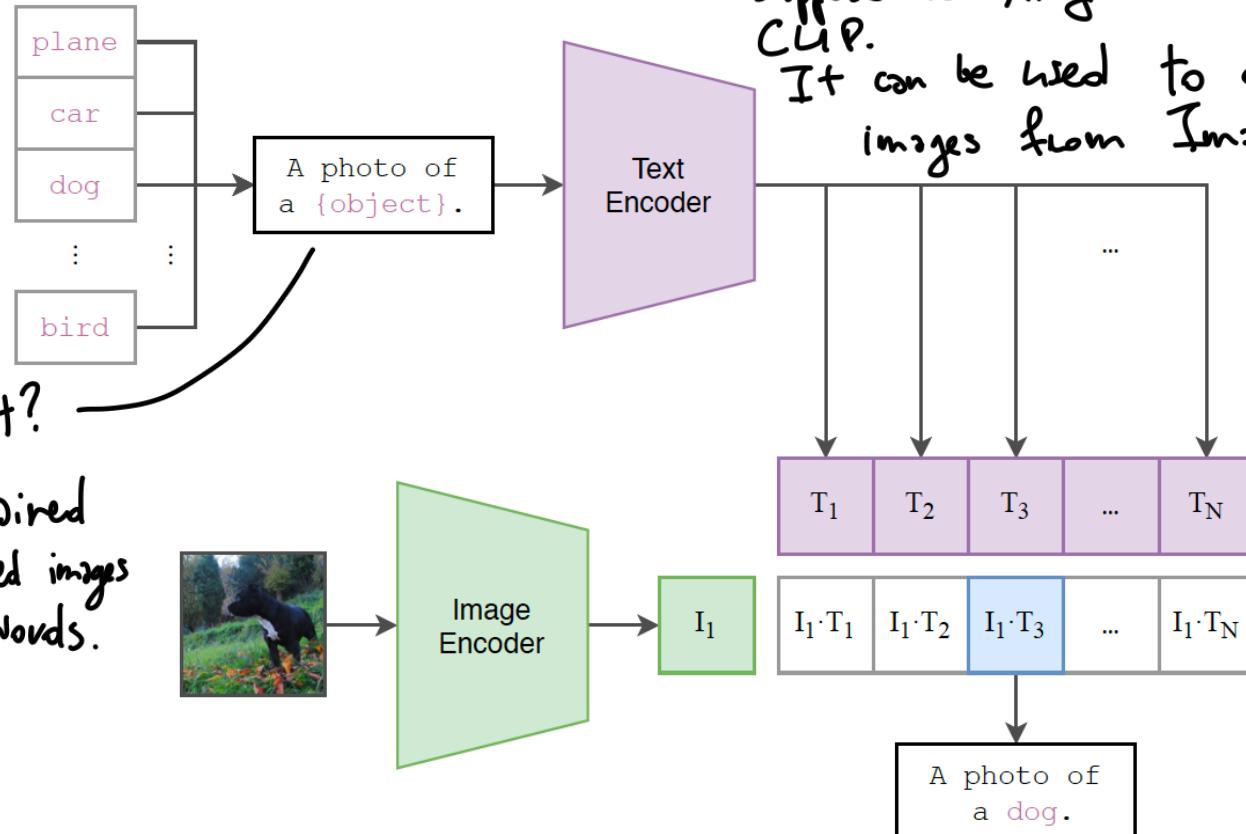
*Compared with every other example in the batch*

# CLIP: Vision + Language Models (VLM)

Pre-trained encoders = **dual encoders** (Text/Image) (RETRIEVAL)  
used for **Zero-shot classifier**, and other downstream tasks

→ Write a sentence → embed → find closest image

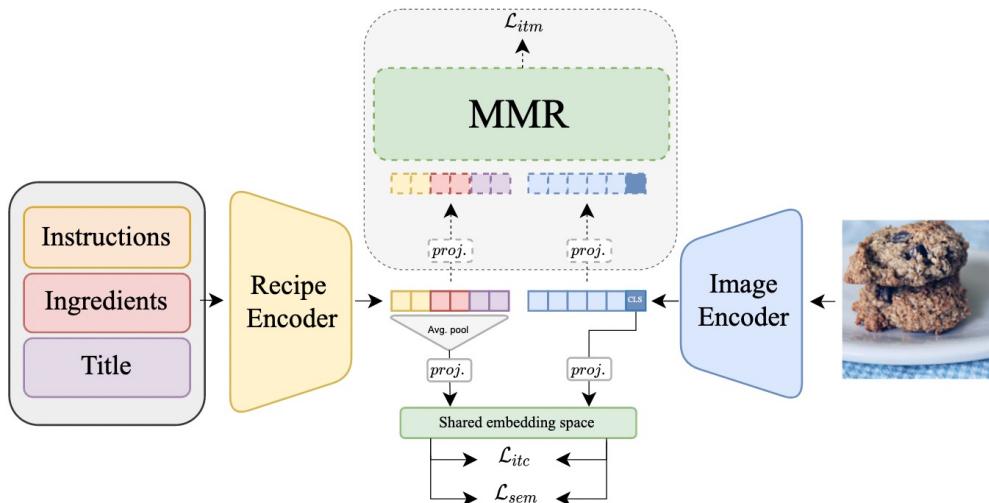
Suppose text, img encoders trained by CLIP.  
It can be used to classify images from ImageNet



why the prompt?

Because it was trained with more detailed images and not single words.

A lot of variants ....



Title query	Ingredient query	Instruction query	Top 5 retrieved images
Mint Chocolate Chip Frosting.	1 cup Unsalted Butter, 2 Tablespoons Heavy Cream, 2 drops Green Food Coloring, .. Chocolate..	Add sugar, cream, peppermint, and food coloring... ...scoop the frosting and place on top of your cupcakes Source: Chocolate Cupcakes with Mint Chocolate Chip ...	
Honey-Grilled Chicken.	1 broiler-fryer chicken, halved, 34 cup butter, melted, 14 cup honey..	Place the halved chicken in a large, shallow container... ...Combine the remaining ingredients, stirring sauce well Grill chicken, skin side up..	
The Best Kale Ever.	1/2 cup Kale, 1 teaspoon Olive Oil, 1/4 teaspoons Red Pepper Flakes..	Wash and cut kale off the stems... Heat olive oil on medium heat and add garlic Add in kale and red pepper flakes..	