

# MAE 688: Machine Learning for Mechanical Engineers

## Lecture 8- Transfer Learning and Explainable AI



Dr. Bing Dong

*Director, Built Environment Science and Technology (BEST) Lab*

*Professor*

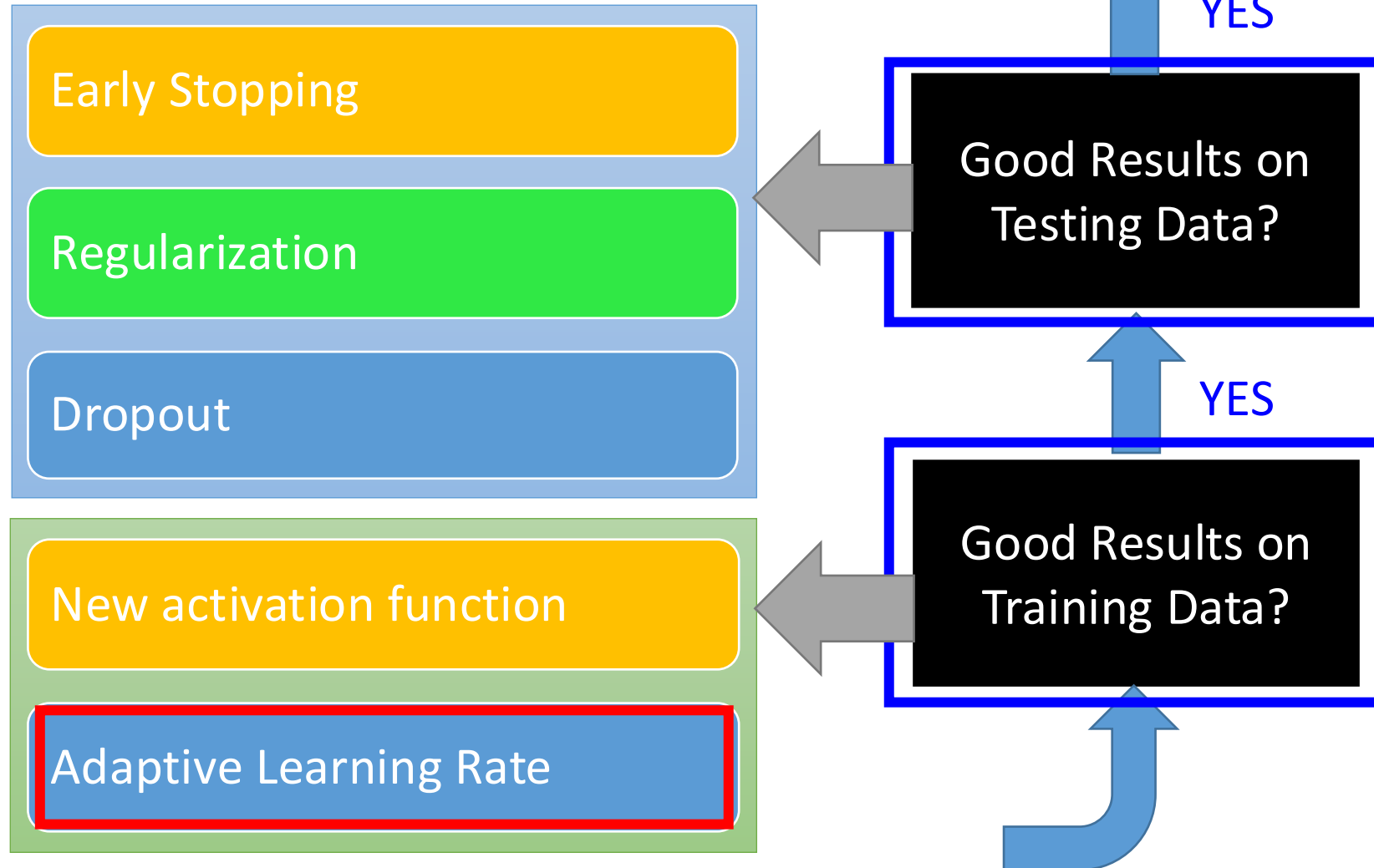
*Mechanical and Aerospace Engineering*

---

Syracuse University

Mach 28, 2023

# Review NN Training Tips:



# What is Transfer Learning

---

“You need a lot of a data if you want to train/use DNN”

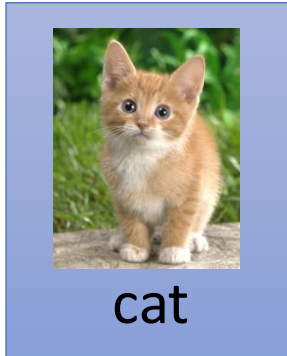
# What is Transfer Learning

---

~~“You need a lot of data if you want to train/use DNN”~~

# What is Transfer Learning

Dog/Cat  
Classifier



Data *not directly related to* the task considered

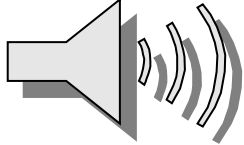

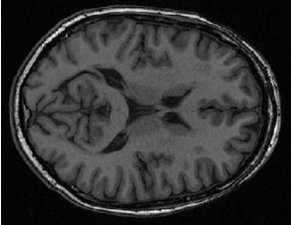


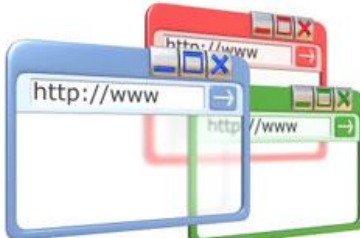


Similar domain, different tasks



Different domains, same task

# Application of Transfer Learning

Task Considered	Data not directly related
<p>Speech Recognition</p>  <p>Chinese</p>	 <p>English Japanese .....</p>
<p>Image Recognition</p>  <p>Medical Images</p>	
<p>Text Analysis</p>  <p>Specific domain</p>	 <p>Webpages</p>

# What is Transfer Learning

		Source Data (not directly related to the task)	
		labelled	unlabeled
Target Data	labelled	Model Fine-tuning	
	unlabeled		

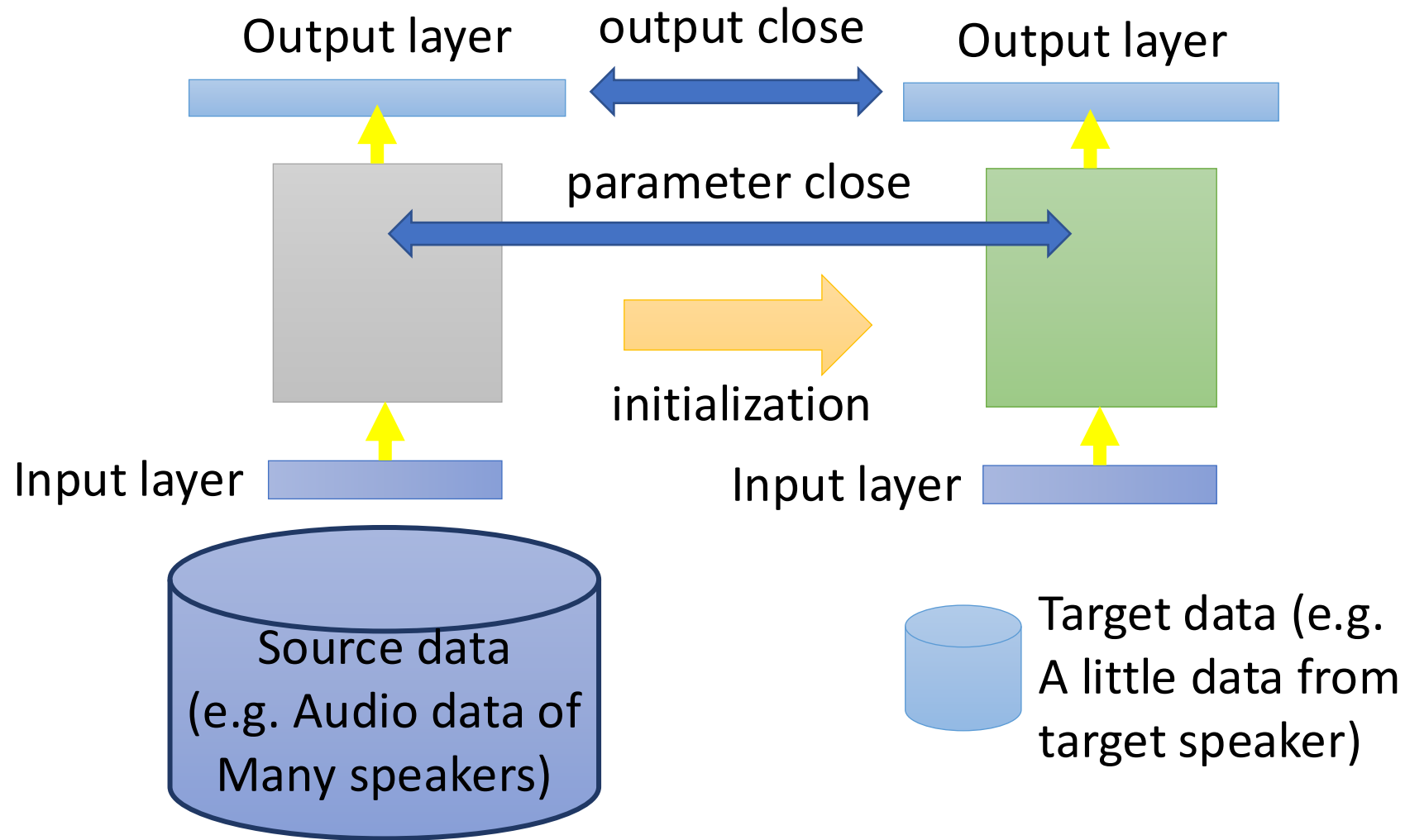
Warning: different terminology in different literature

# Model Fine-tuning

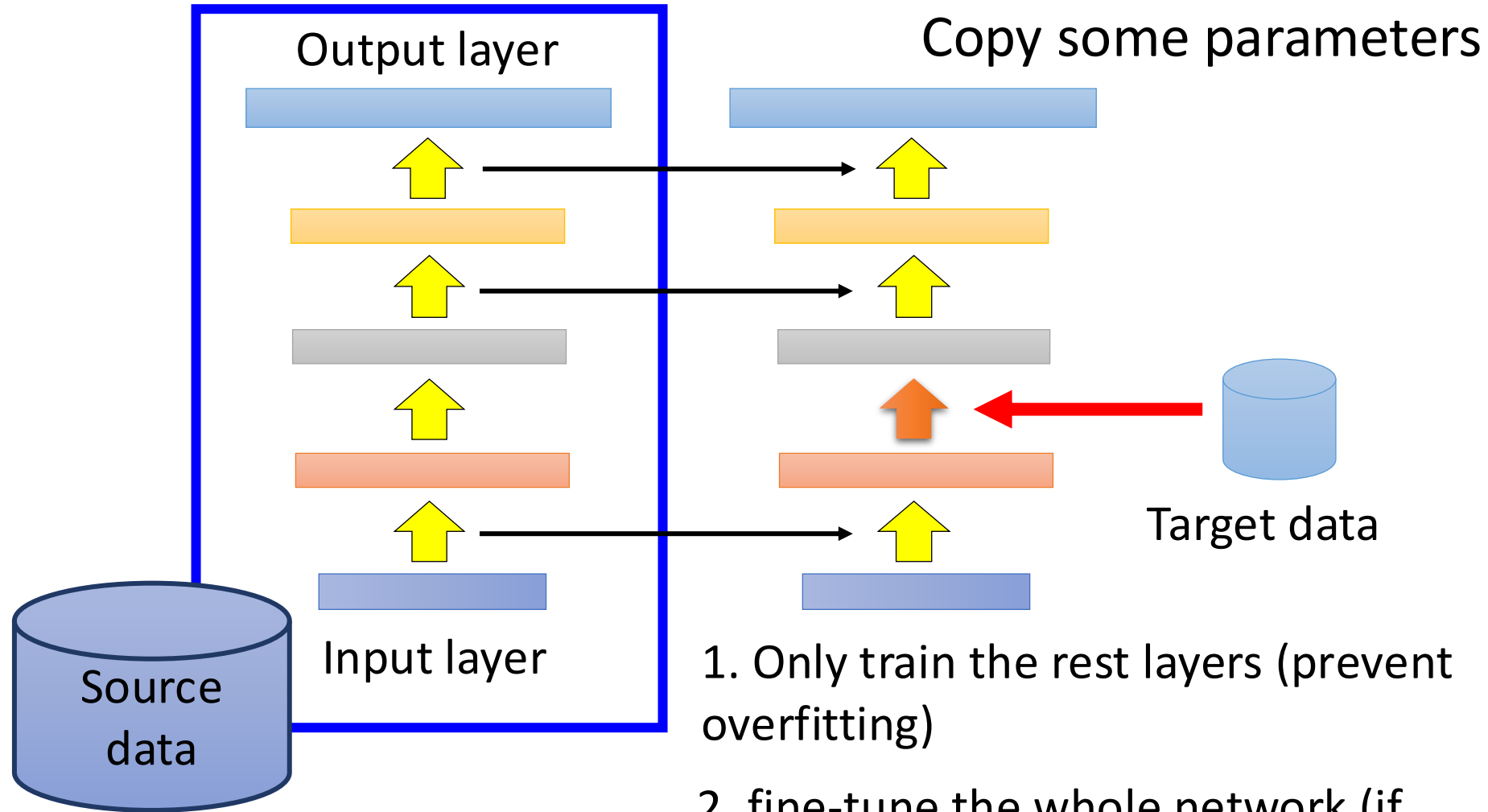
- Task description
  - Source data:  $(x^s, y^s)$  ← A large amount
  - Target data:  $(x^t, y^t)$  ← Very little
- Example: (supervised) speaker adaptation
  - Source data: audio data and transcriptions from many speakers
  - Target data: audio data and its transcriptions of specific user
- Idea: training a model by source data, then fine-tune the model by target data
  - Challenge: only limited target data, so be careful about overfitting

One-shot learning: only a few examples in target domain

# Conservative Training



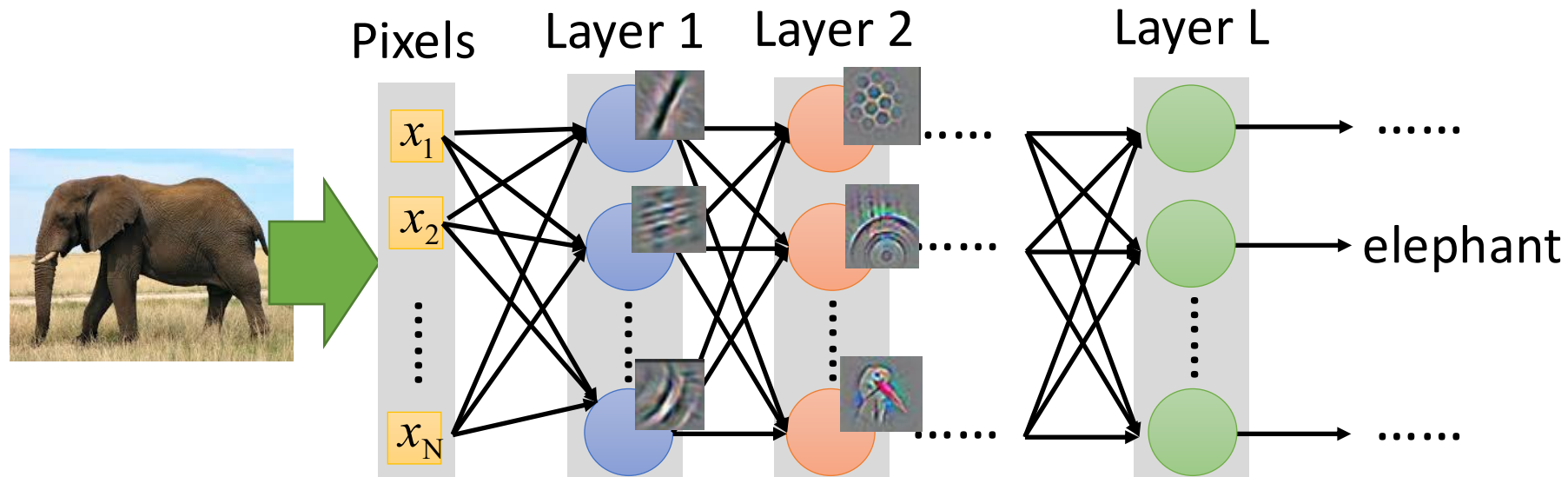
# Layer Transfer



1. Only train the rest layers (prevent overfitting)
2. fine-tune the whole network (if there is sufficient data)

# Layer Transfer

- Which layer can be transferred (copied)?
  - Speech: usually copy the last few layers
  - Image: usually copy the first few layers

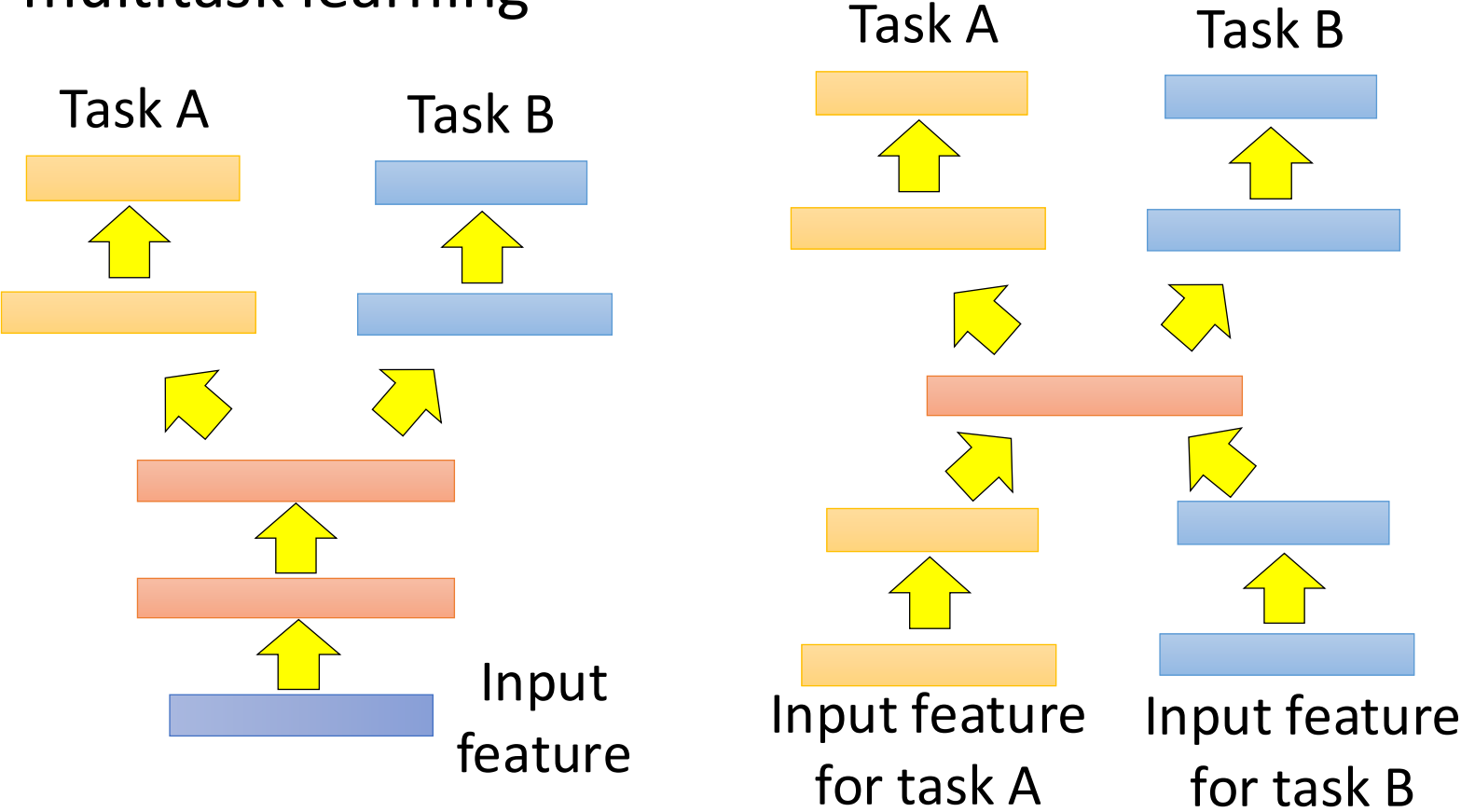


# What is Transfer Learning

		Source Data (not directly related to the task)	
		labelled	unlabeled
Target Data	labelled	Model Fine-tuning Multitask Learning	
	unlabeled		

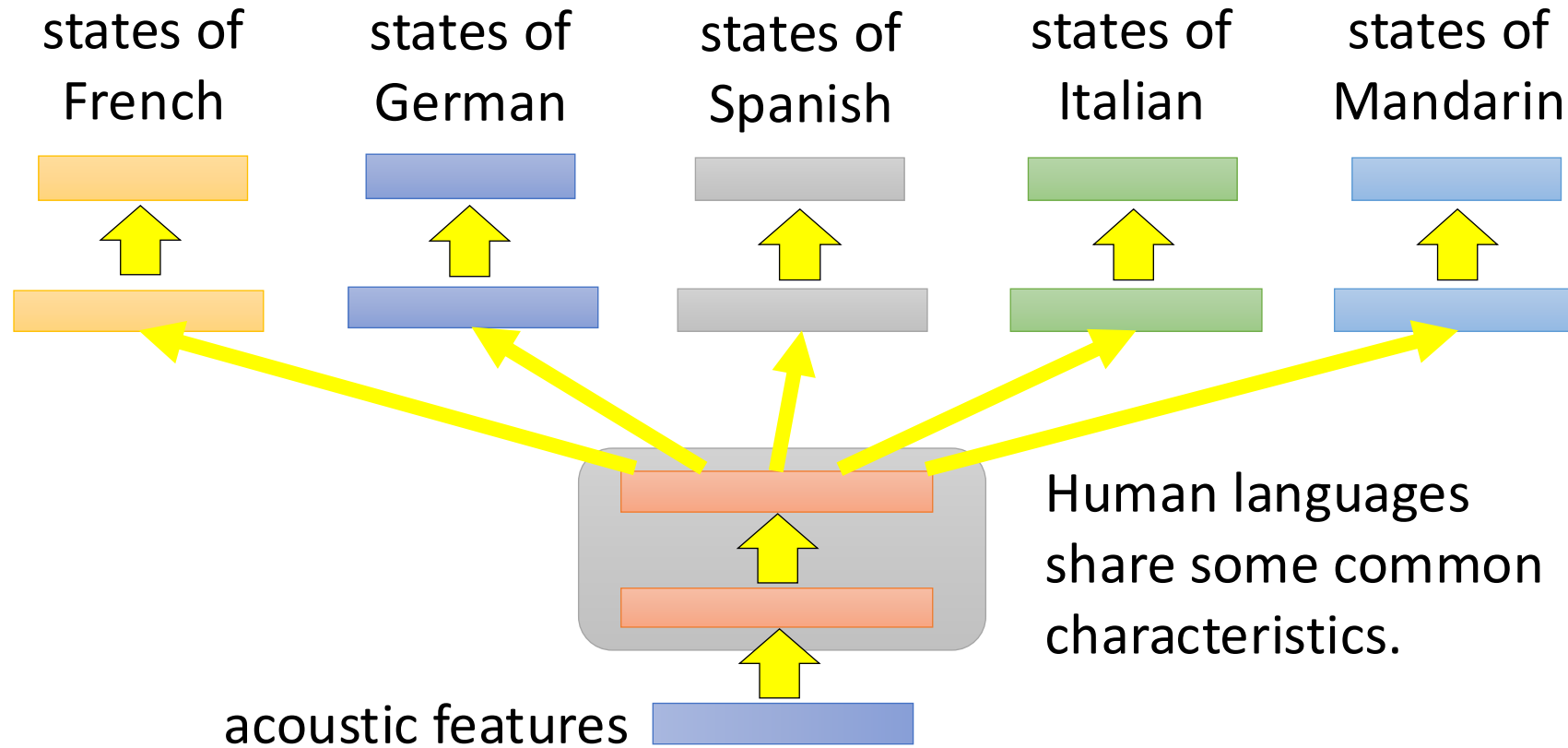
# Multitask Learning

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



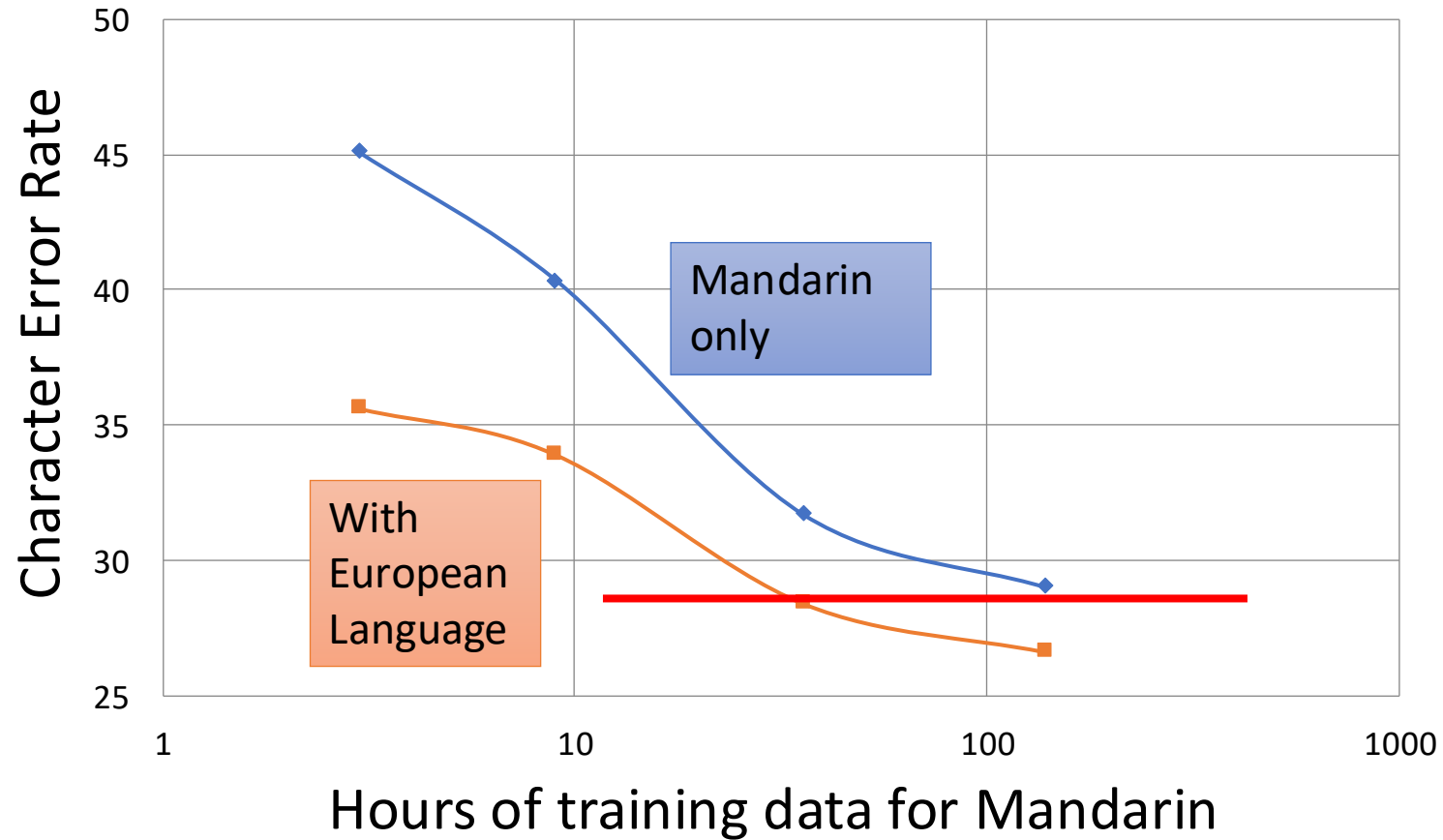
# Multitask Learning

## - Multilingual Speech Recognition



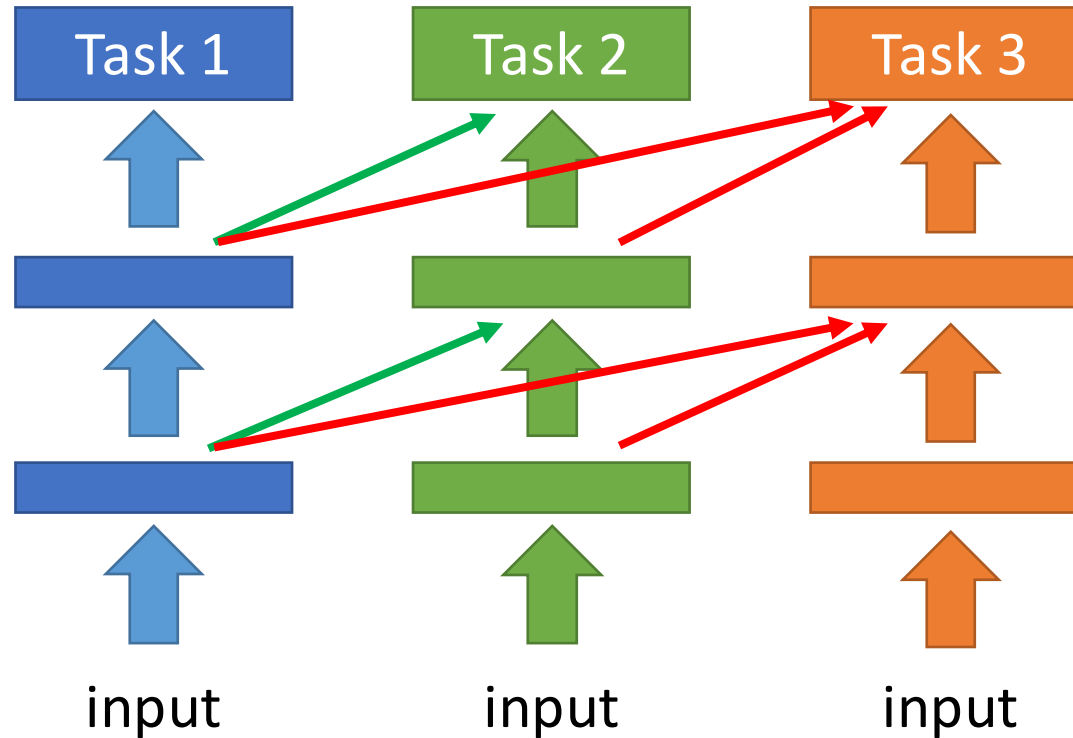
***Similar idea in translation:*** Daxiang Dong, Hua Wu, Wei He, Dianhai Yu and Haifeng Wang, "Multi-task learning for multiple language translation.", ACL 2015

# Multitask Learning - Multilingual



Huang, Jui-Ting, et al. "Cross-language knowledge transfer using multilingual deep neural network with shared hidden layers." *ICASSP, 2013*

# Progressive Neural Networks



Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, Raia Hadsell, "Progressive Neural Networks", arXiv preprint 2016

# What is Transfer Learning

		Source Data (not directly related to the task)	
		labelled	unlabeled
Target Data	labelled	Model Fine-tuning Multitask Learning	
	unlabeled	Domain-adversarial training	

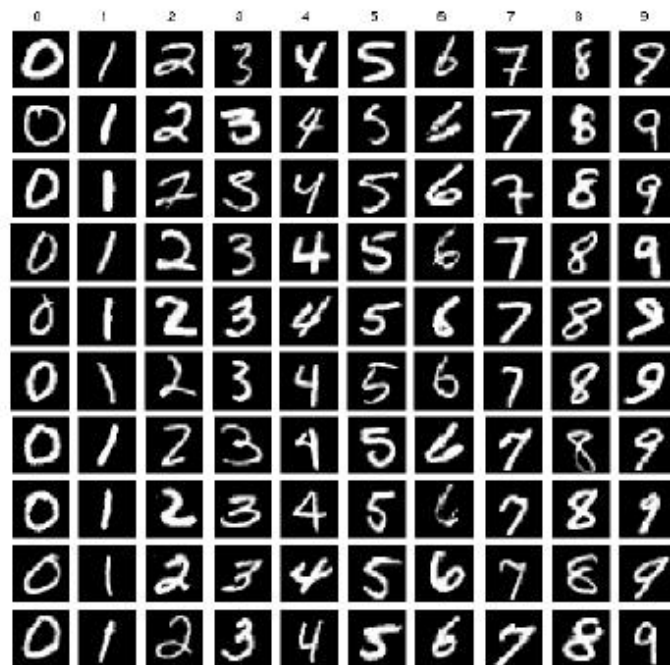
# Task Description

- Source data:  $(x^s, y^s)$   $\longrightarrow$  Training data
  - Target data:  $(x^t)$   $\longrightarrow$  Testing data
- } Same task, mismatch



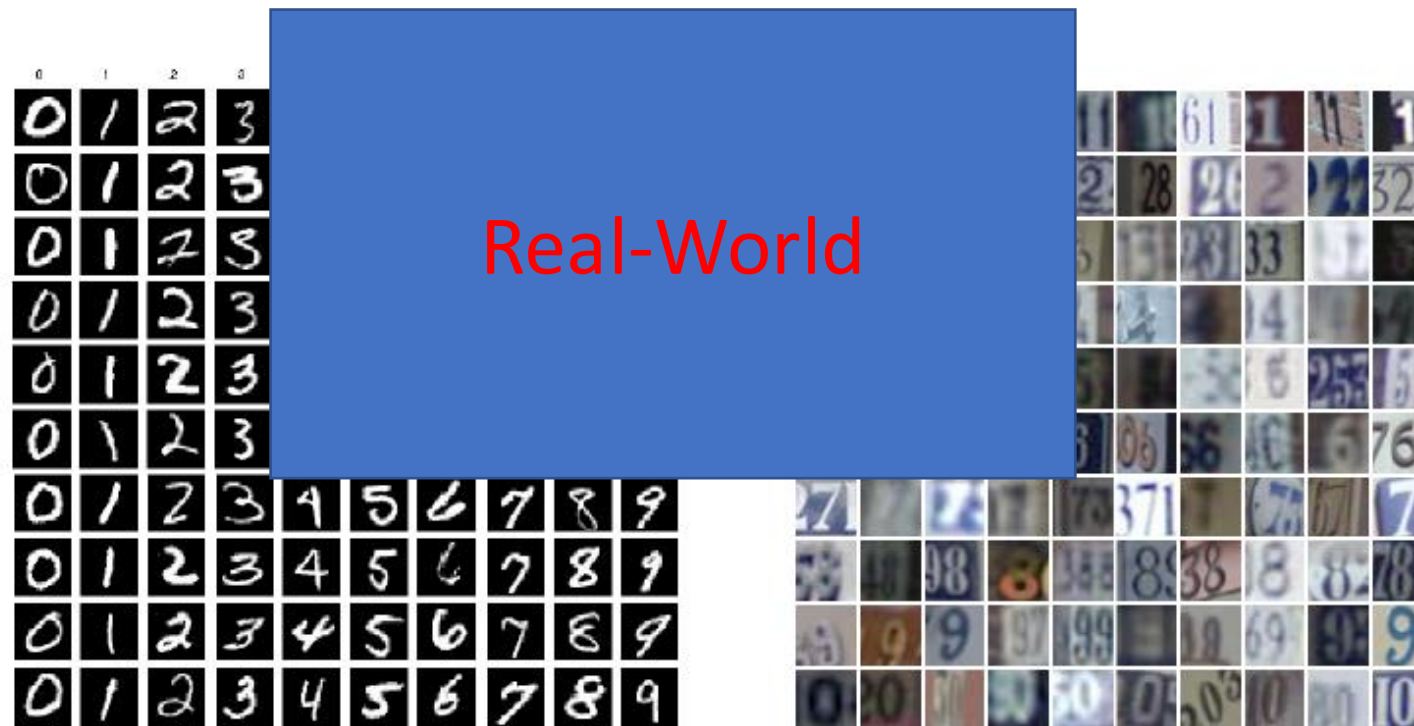
# Domain Adaption

- ★ Often times, the **data of the actual task** we want to tackle is **largely different from the training data** we use.
- ★ We need to think of ways to alleviate this problem.



# Domain Adaption

- ★ Often times, the **data of the actual task** we want to tackle is **largely different from the training data** we use.
- ★ We need to think of ways to alleviate this problem.



# Definition: Domain Adaption

## Domain $D$



# Definition: Domain Adaption

---

Task  $T$

Classification

“Biker”



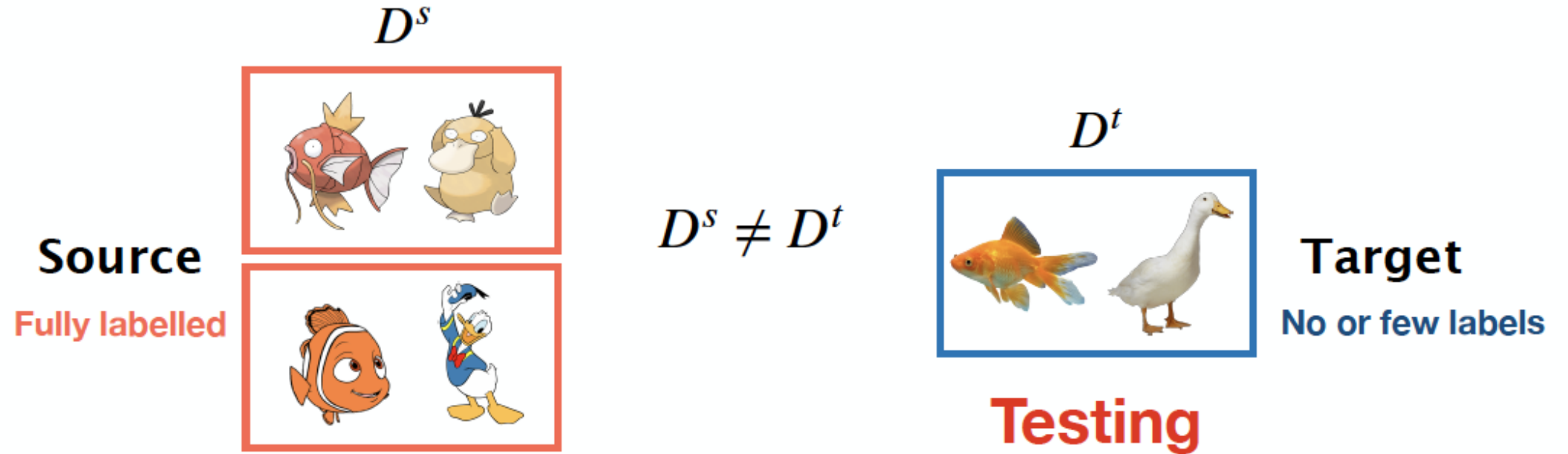
Segmentation



Bounding box detection

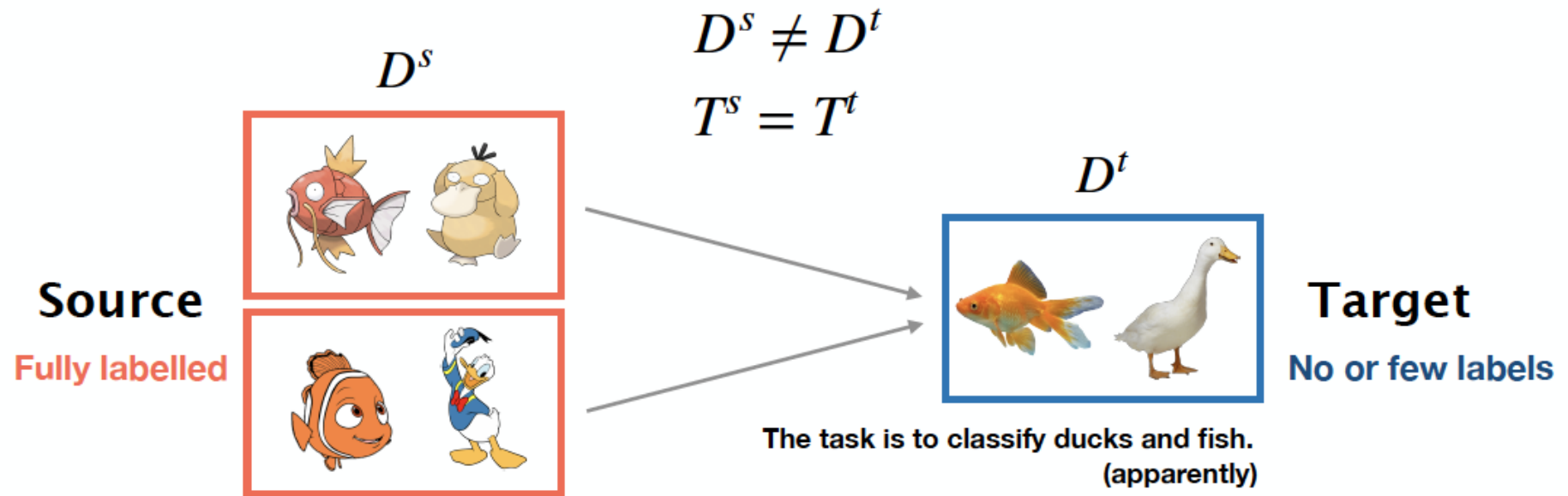


# Definition: Domain Adaption



# Definition: Domain Adaption

Tasks are the same



# Why: Domain Adaption

When it comes to **autonomous driving** ...



We are able to easily amass unlimited amounts of street view data from **Grand Theft Auto**.



Not the case for arbitrary real-life scenarios.



# Why: Domain Adaption

When it comes to **object classification** ...



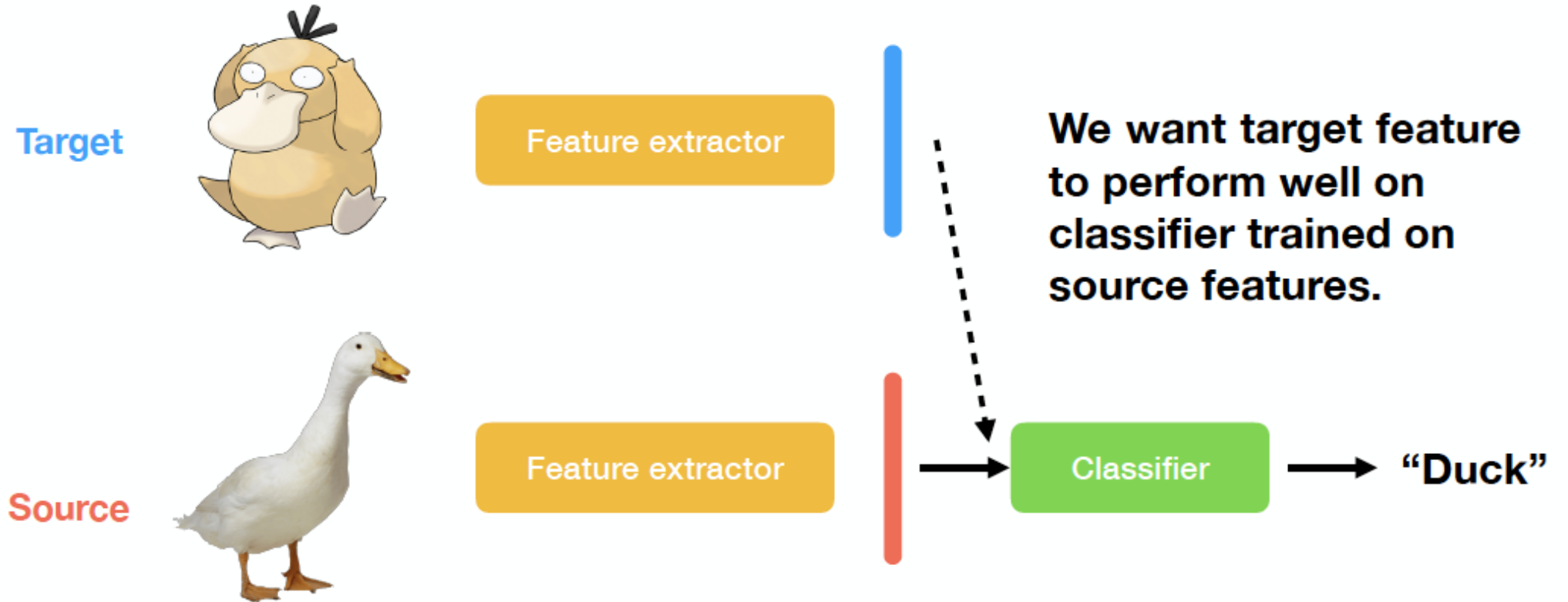
Easily scrape well-annotated images from ebay, Amazon for your model ...



Generalizes badly for natural images due to lighting, background, and etc ...

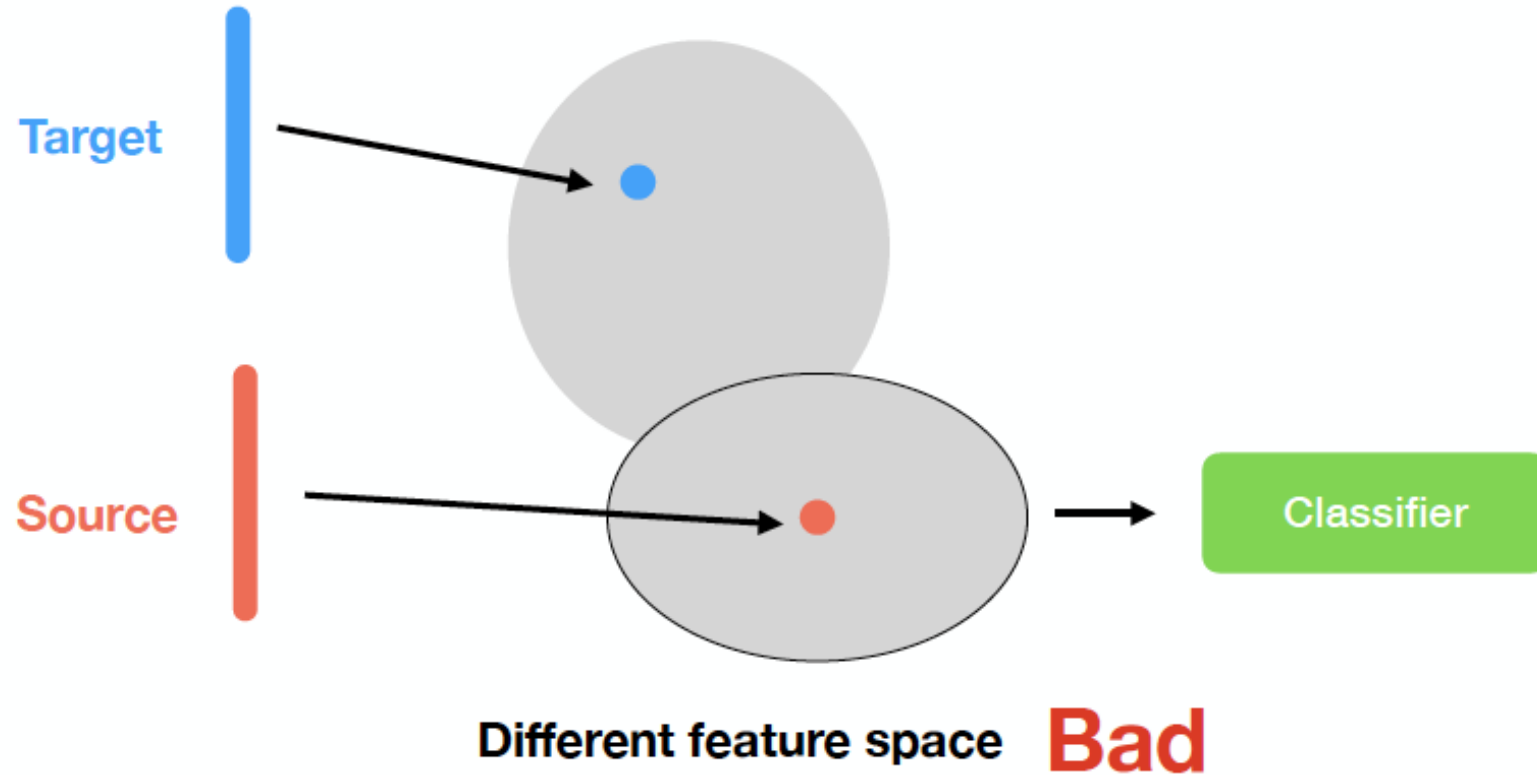


# Feature: Domain Adaption

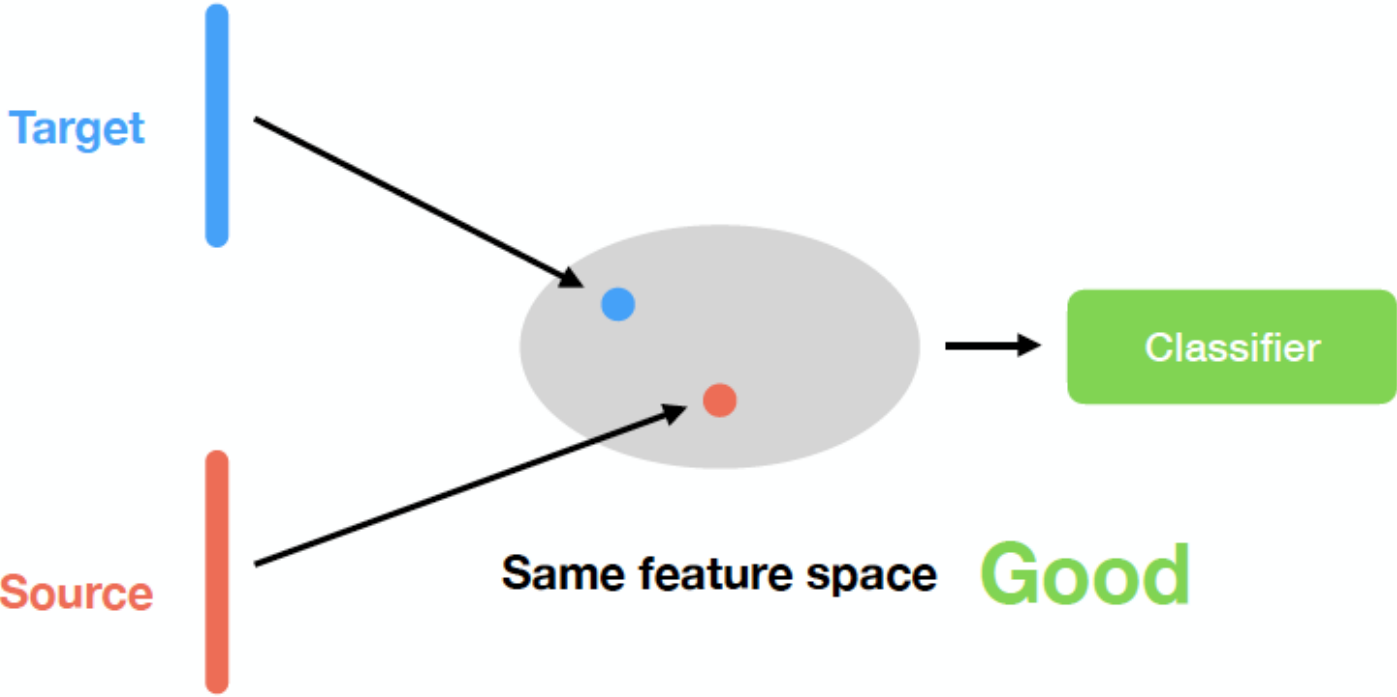


# Feature: Domain Adaption

---



# Feature: Domain Adaption

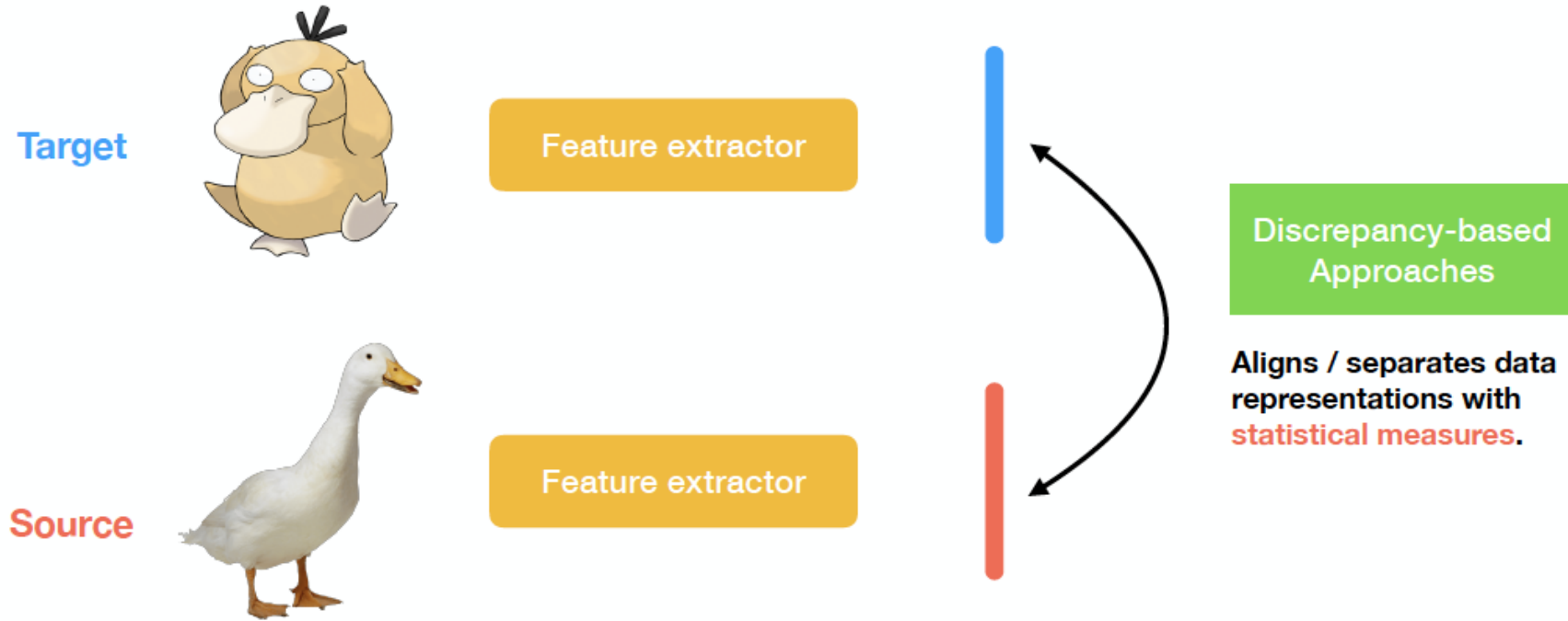


# Methods for Domain Adaption

---

- 1. Discrepancy-based methods**
- 2. Adversarial-based methods**
- 3. Reconstruction-based methods**

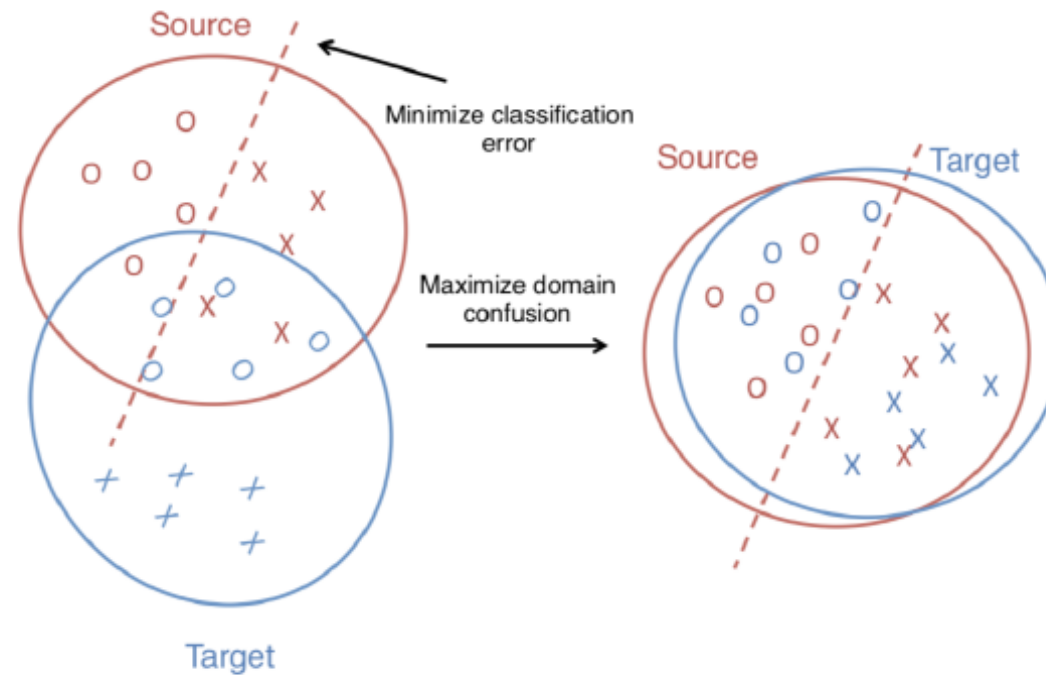
# Discrepancy-based



# Deep Domain Confusion

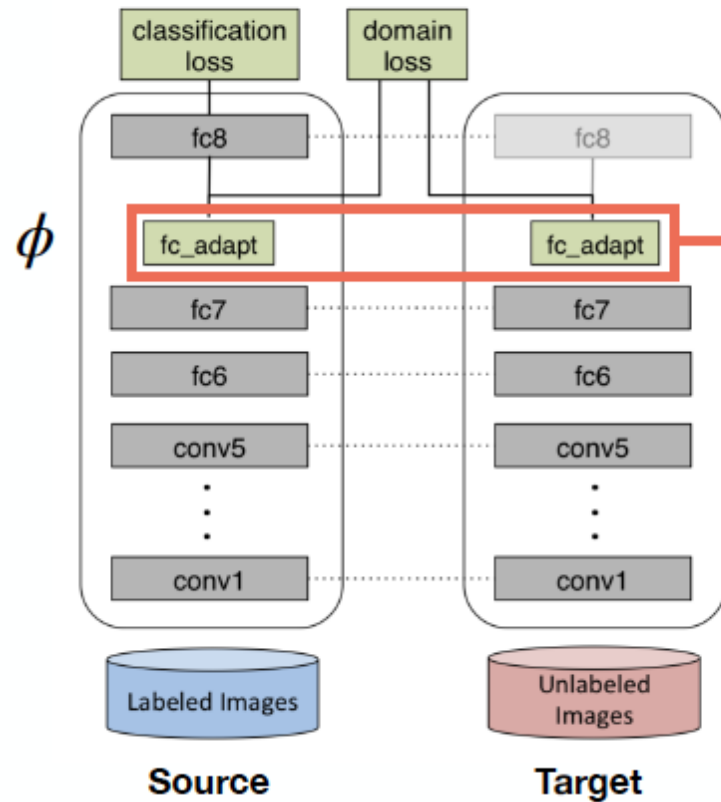
## Deep Domain Confusion: Maximizing for Domain Invariance

Tzeng et al, arXiv 1412.3474, 2014



# Deep Domain Confusion

[https://en.wikipedia.org/wiki/Kernel\\_embedding\\_of\\_distributions](https://en.wikipedia.org/wiki/Kernel_embedding_of_distributions)  
[https://en.wikipedia.org/wiki/Kernel\\_method](https://en.wikipedia.org/wiki/Kernel_method)



**Maximum Mean Discrepancy**  
Measuring distance between distributions

$$\text{MMD}(X_S, X_T) =$$

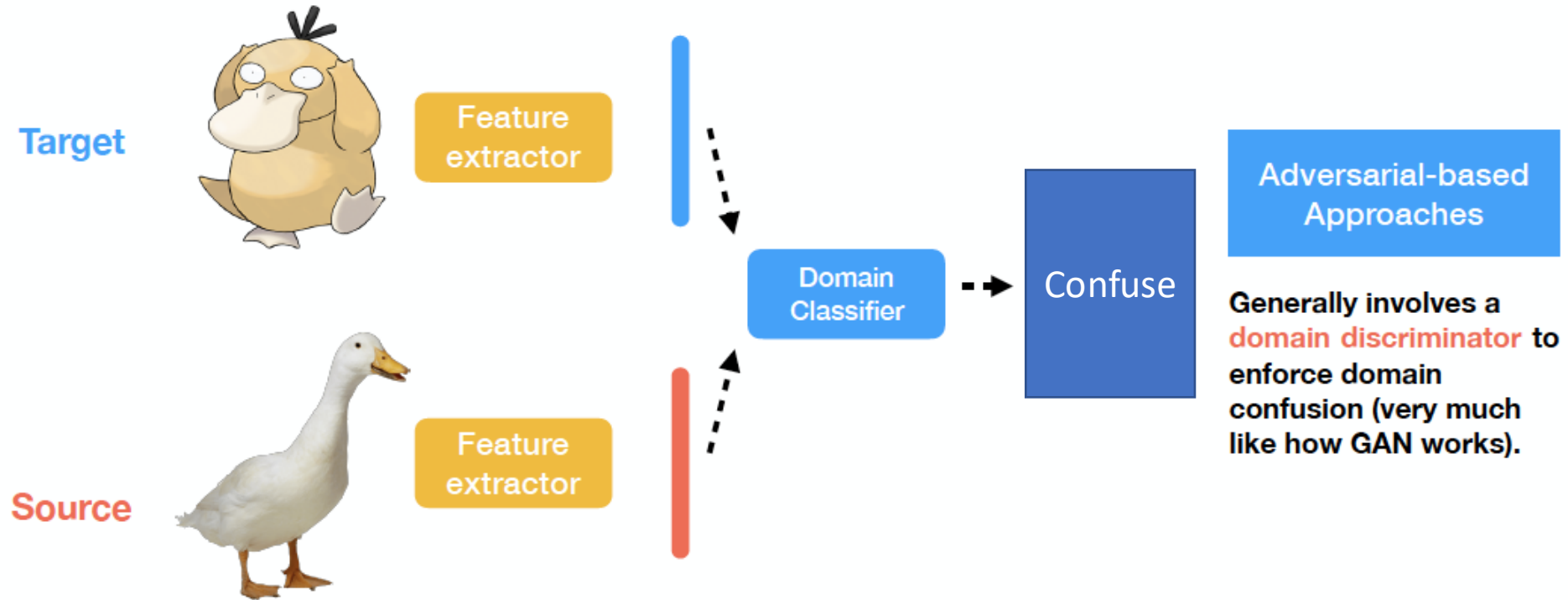
$$\left\| \frac{1}{|X_S|} \sum_{x_s \in X_S} \phi(x_s) - \frac{1}{|X_T|} \sum_{x_t \in X_T} \phi(x_t) \right\|$$

**Objective Function**  
Classification + MMD

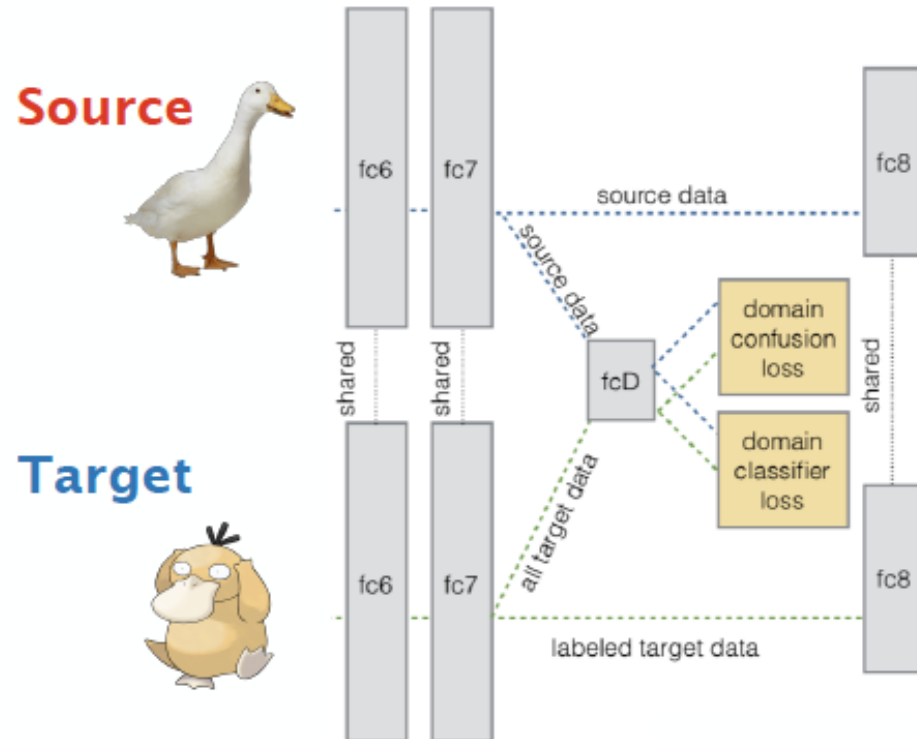
$$\mathcal{L} = \mathcal{L}_C(X_L, y) + \lambda \text{MMD}^2(X_S, X_T)$$

How do you choose  $\phi$ ?

# Adversarial based



# Domain-adversarial training



**Iteration t**  
**Learn a good domain classifier**  
**fix feature extractor**

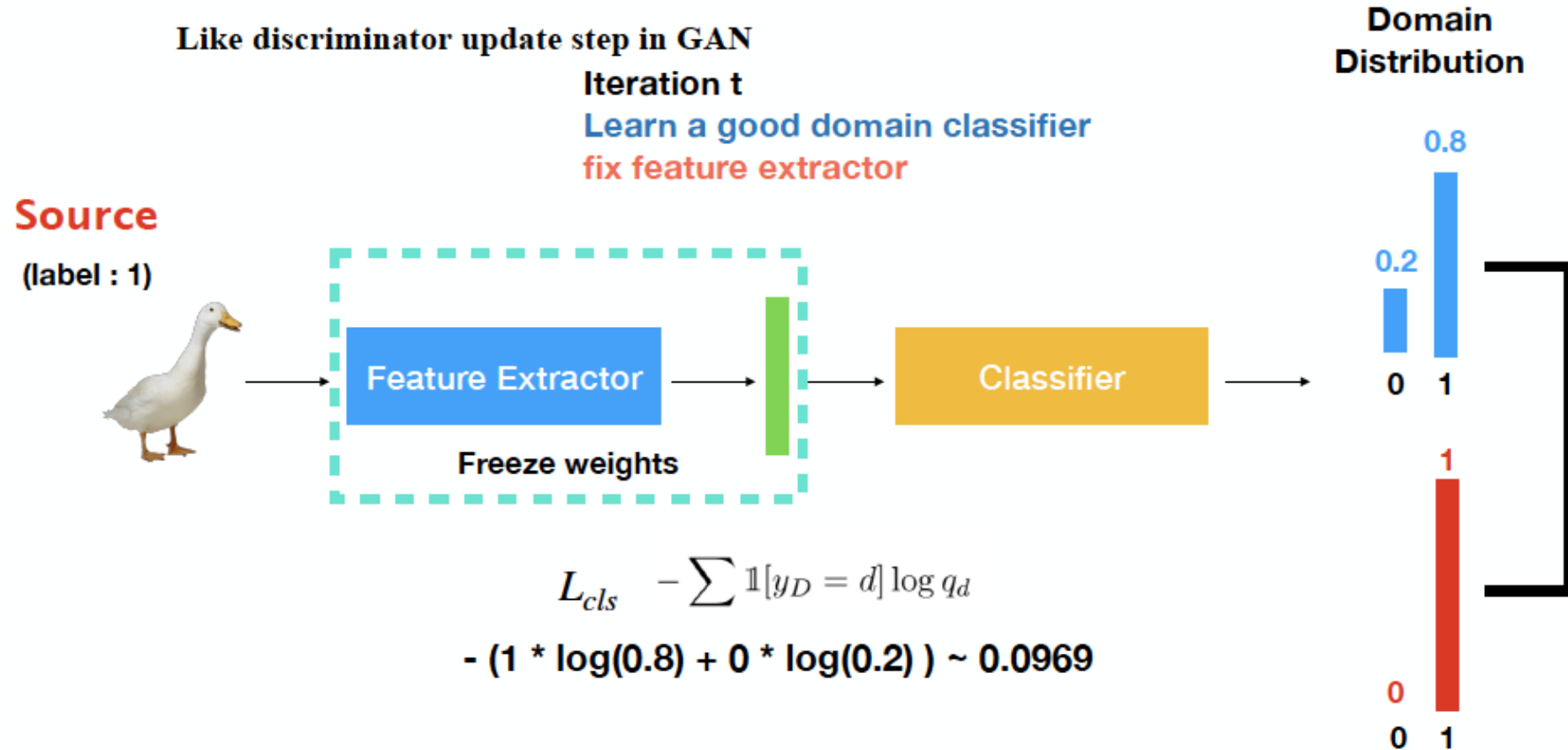
$$L_{cls} = - \sum \mathbb{1}[y_D = d] \log q_d$$

**Iteration t + 1**  
**Learn a representation that confuses the classifier (uniform distribution)**  
**fix domain classifier**

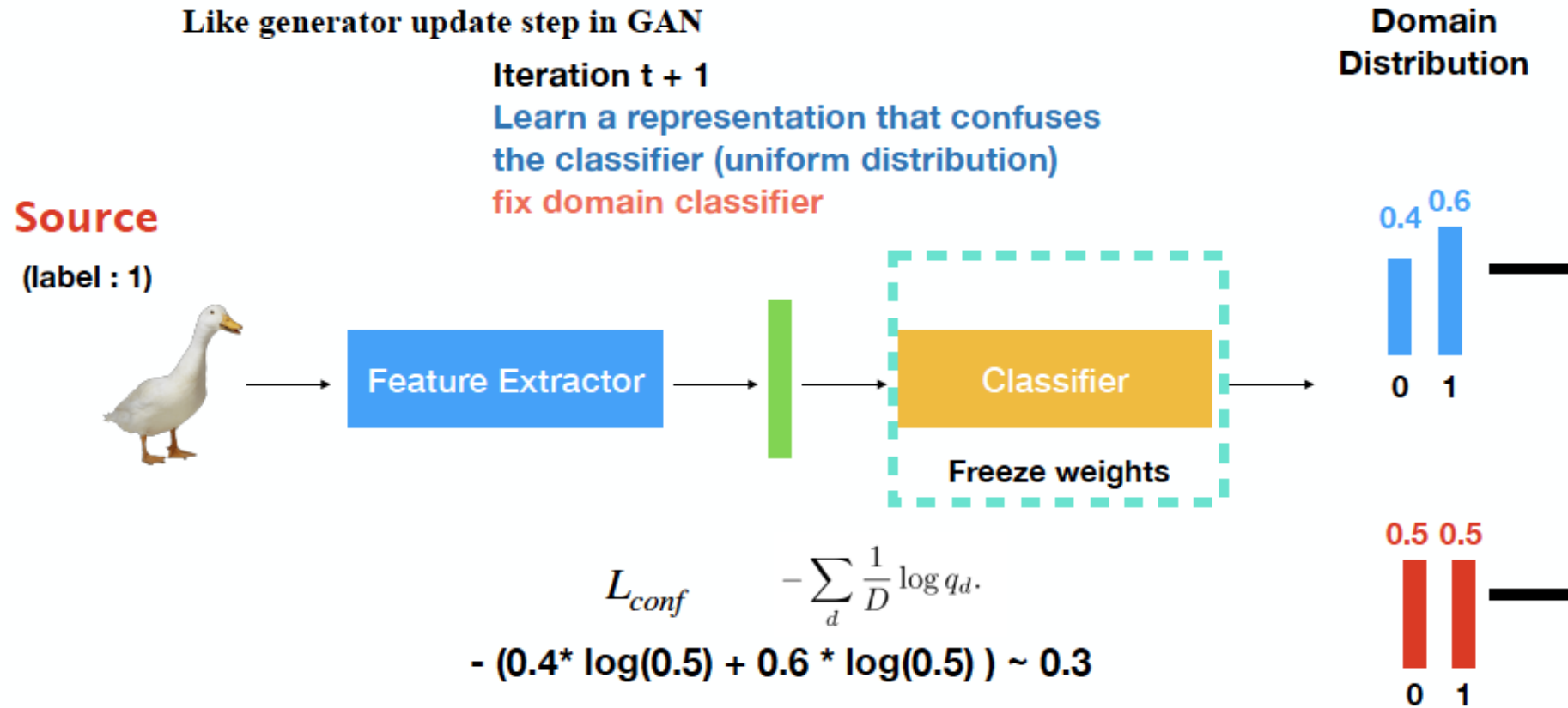
$$L_{conf} = - \sum_d \frac{1}{D} \log q_d$$

**Both  
 Cross  
 Entropy !**

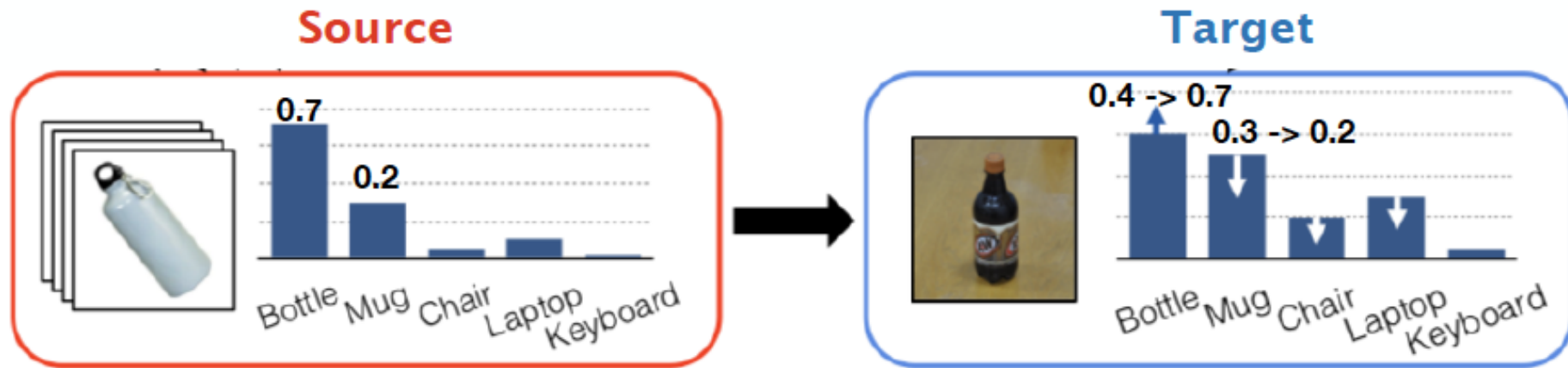
# Domain-adversarial training



# Domain-adversarial training

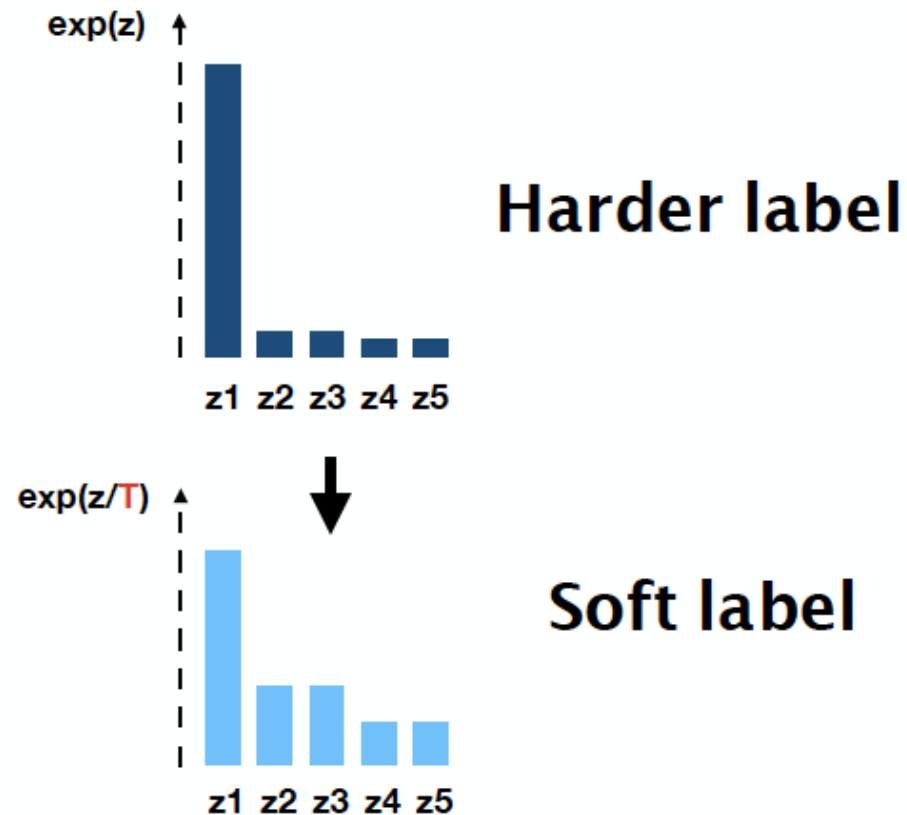


# Domain-adversarial training-Label Correlation



Idea: **“bottle”** should be similar to **“mug”** but less similar to **“laptop”**  
Class distributions of source and target should be similar

# Domain-adversarial training-Label Correlation



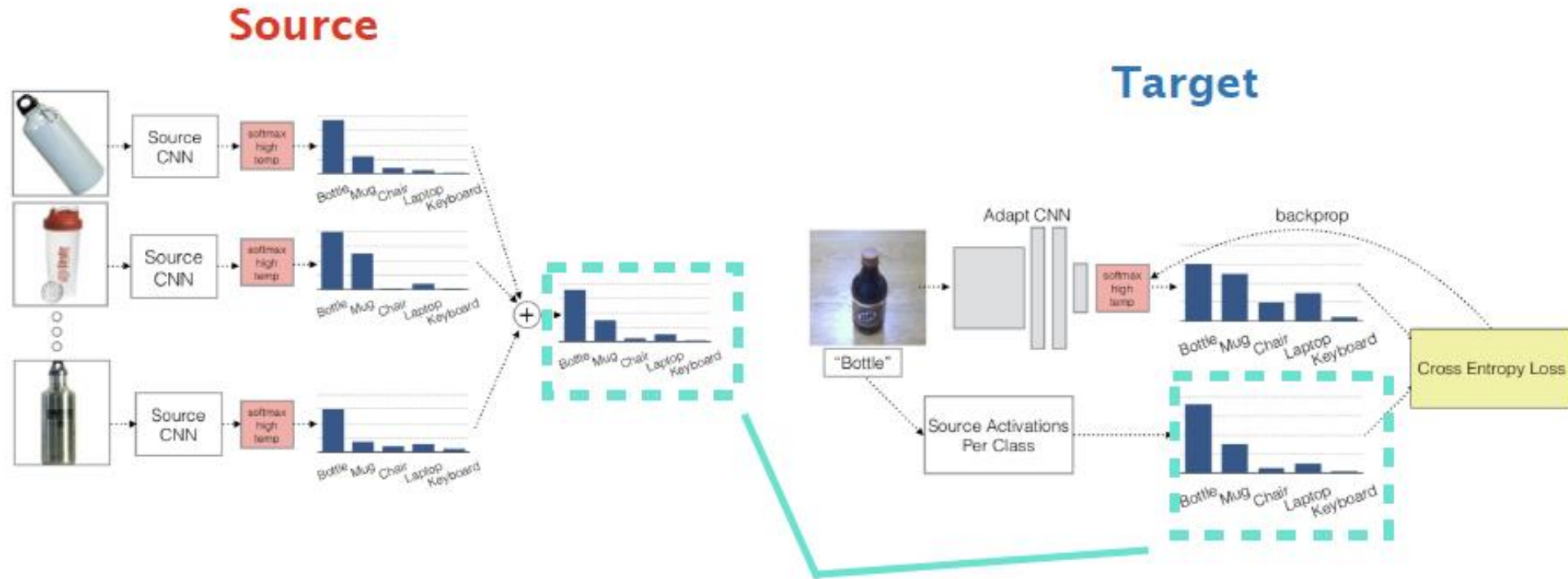
## Softmax with Temperature

Enables a smoother class probability distribution to address inter-class correlations

$$\frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}$$

“Scale down probabilities”

# Domain-adversarial training-Label Correlation

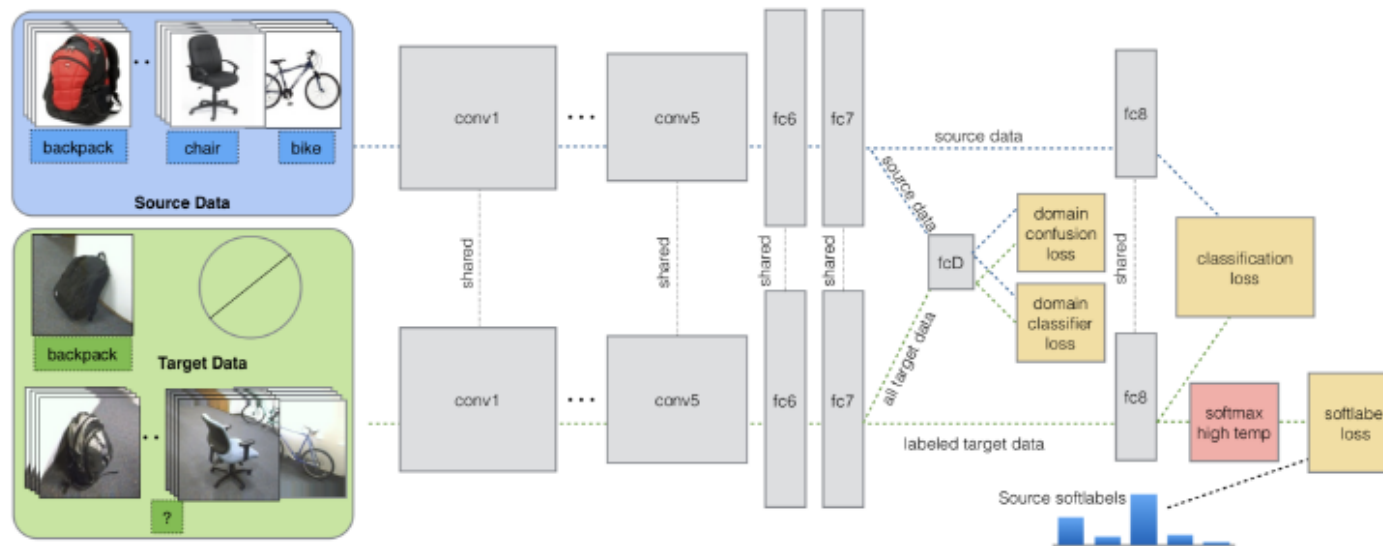


# Domain-adversarial training-Final

## Simultaneous Deep Transfer Across Domains and Tasks

Tzeng et al, ICCV 2015

$$L = L_{cls} + L_{conf} + L_{soft}$$



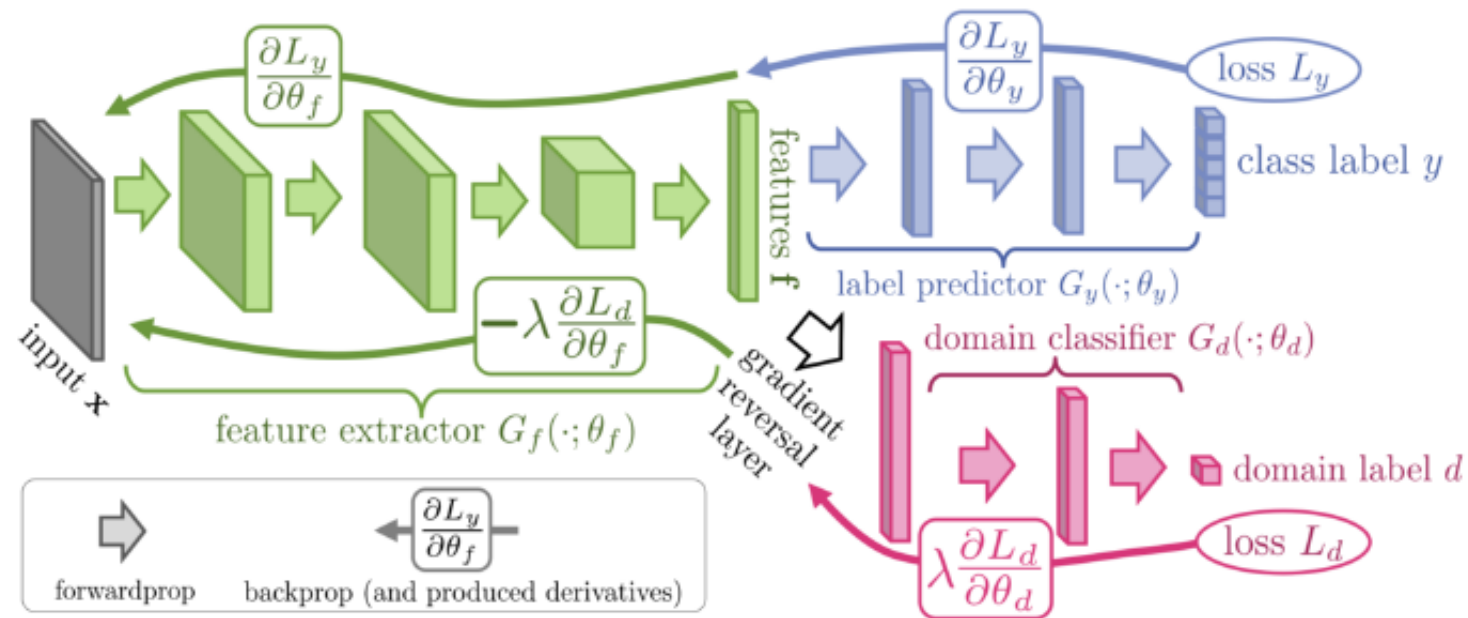
# Domain-adversarial training

## DANN

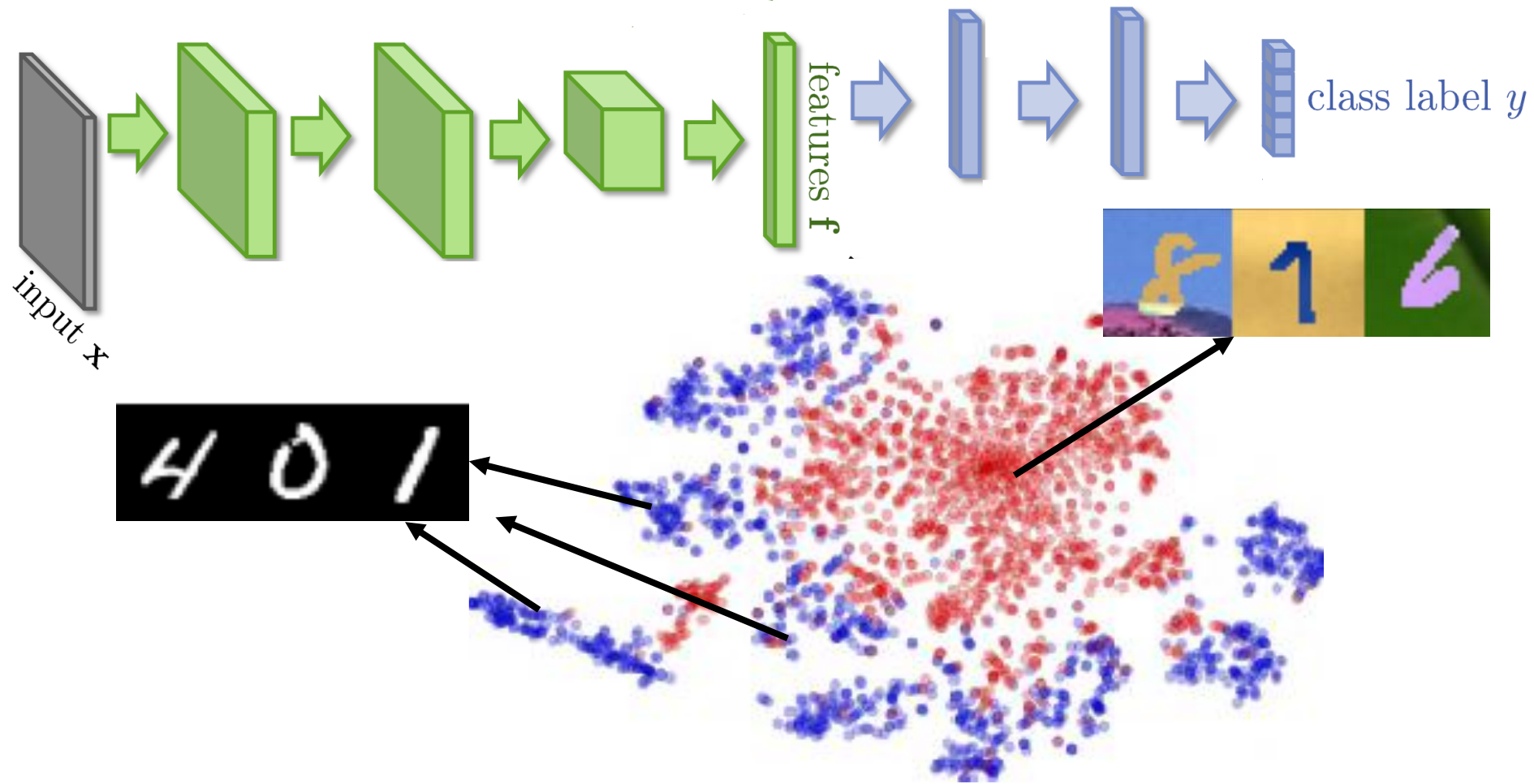
### Domain Adversarial Training of Neural Networks

Ganin et al, NIPS 2016

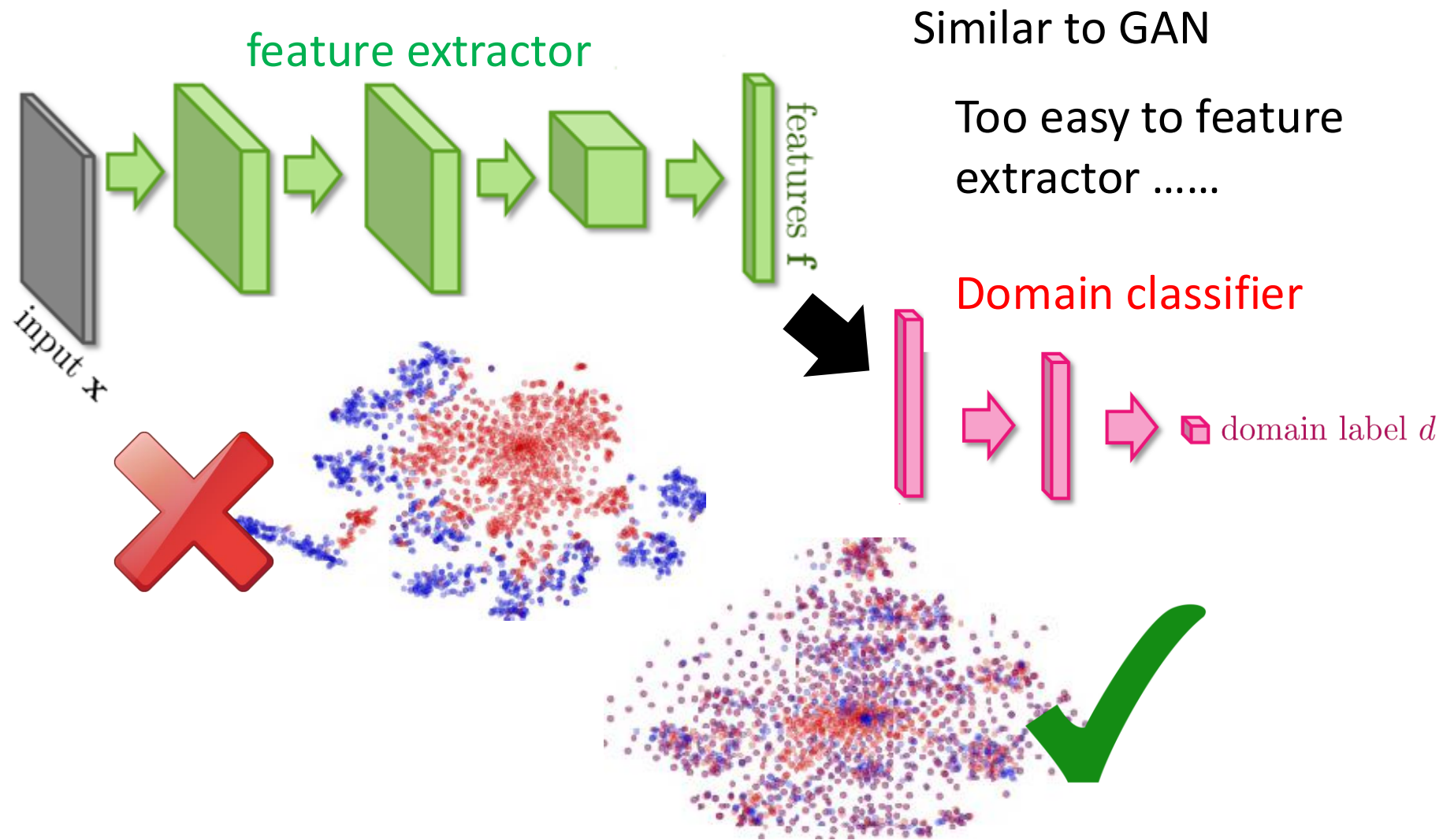
- **Spotlight : Domain Classifier + Gradient Reversal**



# Domain-adversarial training

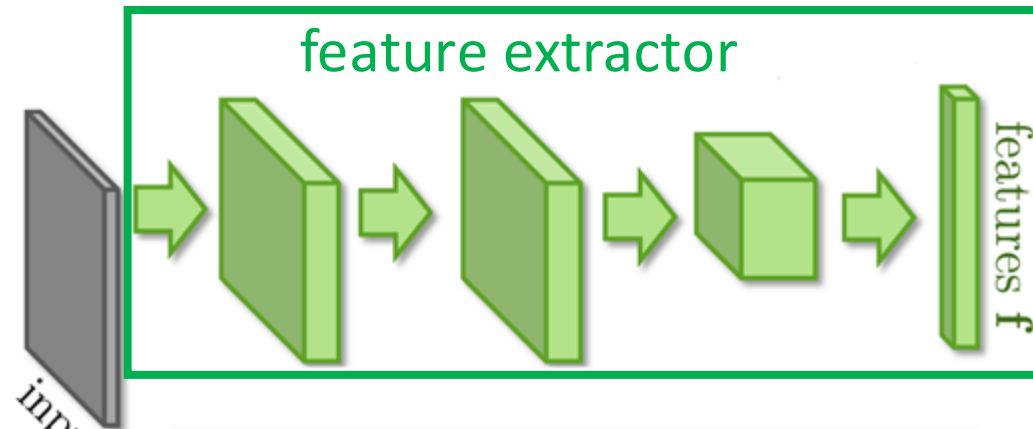


# Domain-adversarial training



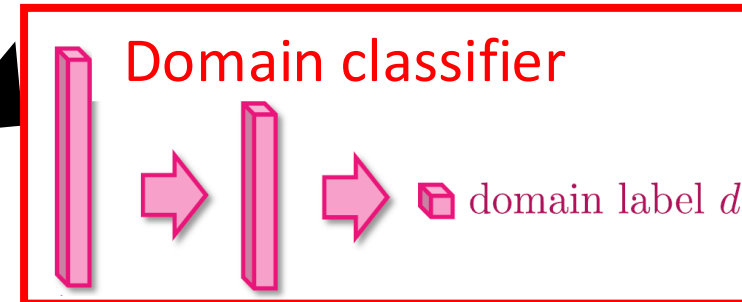
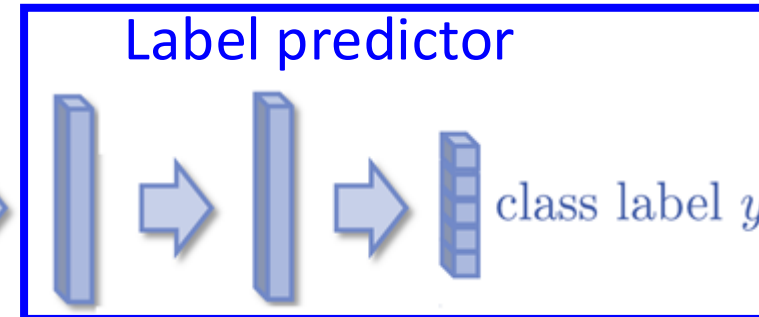
# Domain-adversarial training

Maximize label classification accuracy + minimize domain classification accuracy



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

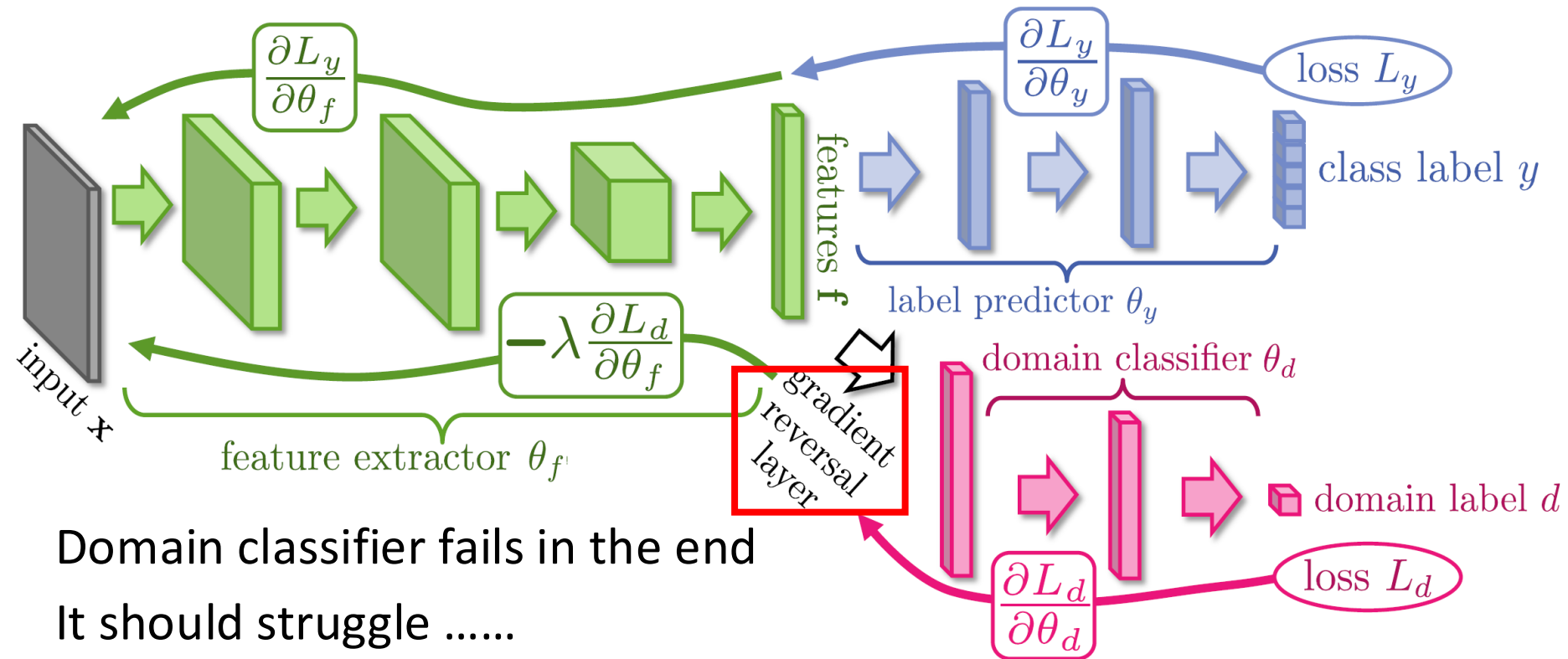
Maximize label classification accuracy



Maximize domain classification accuracy

This is a big network, but different parts have different goals.

# Domain-adversarial training



Domain classifier fails in the end  
It should struggle .....

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



METHOD	SOURCE	MNIST	SYN NUMBERS	SVHN	SYN SIGNS
	TARGET	MNIST-M	SVHN	MNIST	GTSRB
SOURCE ONLY		.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 (57.9%)</b>	<b>.9048 (66.1%)</b>	<b>.7107 (29.3%)</b>	<b>.8866 (56.7%)</b>
TRAIN ON TARGET		.9891	.9244	.9951	.9987

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

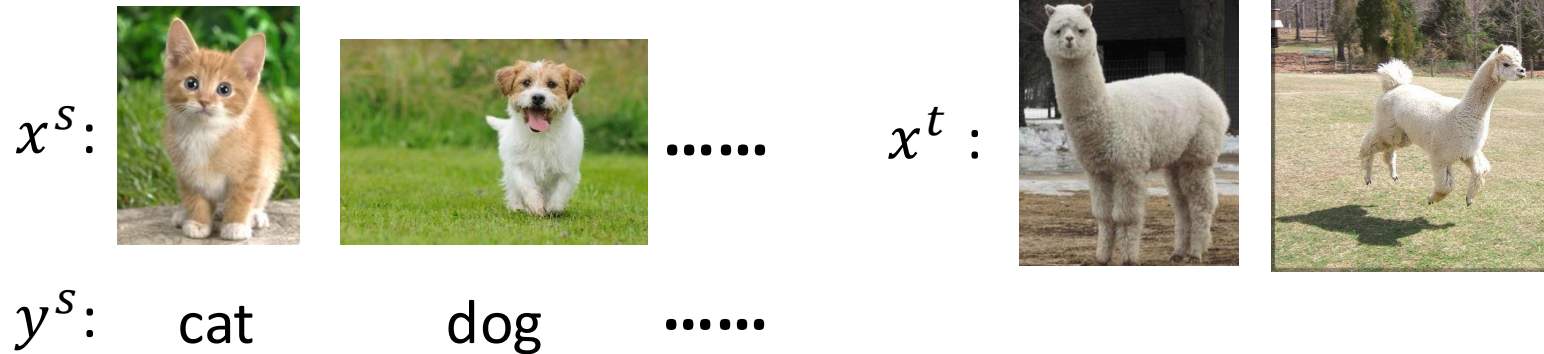
Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, Domain-Adversarial Training of Neural Networks, JMLR, 2016

# What is Transfer Learning

		Source Data (not directly related to the task)	
		labelled	unlabeled
Target Data	labelled	Model Fine-tuning Multitask Learning	
	unlabeled	Domain-adversarial training Zero-shot learning	

# Zero-shot Learning

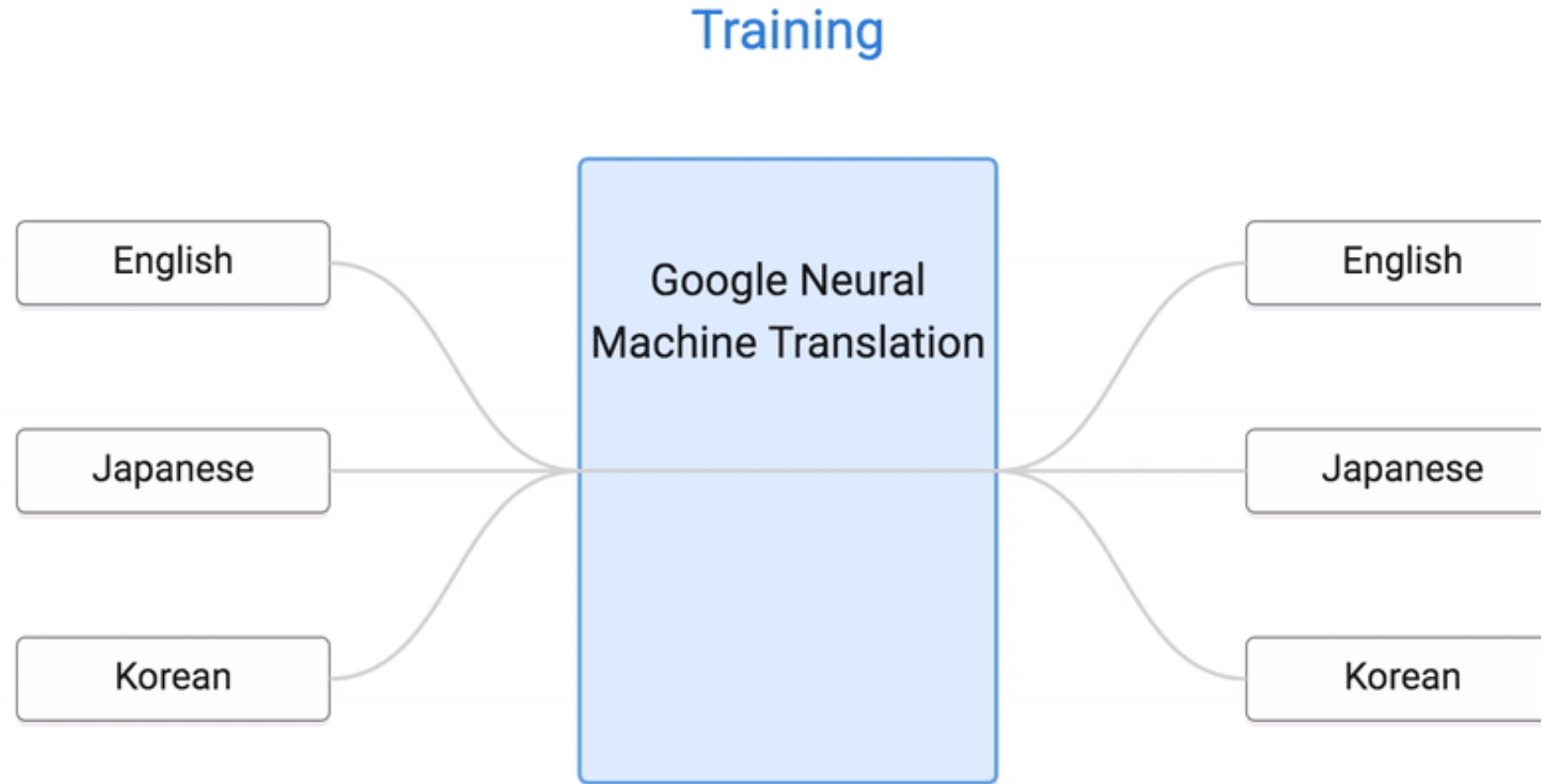
- Source data:  $(x^s, y^s)$   $\longrightarrow$  Training data
  - Target data:  $(x^t)$   $\longrightarrow$  Testing data
- } Different tasks



In speech recognition, we can not have all possible words in the source (training) data.

How we solve this problem in speech recognition?

# Example of Zero-Shot Learning



Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat. Google's Multilingual Neural Machine Translation System: Enabling Zero-Shot Translation, arXiv preprint 2016

# What is Transfer Learning

		Source Data (not directly related to the task)	
		labelled	unlabeled
Target Data	labelled	<p>Model Fine-tuning</p> <p>Multitask Learning</p>	<p>Self-taught learning</p> <p>Rajat Raina , Alexis Battle , Honglak Lee , Benjamin Packer , Andrew Y. Ng, Self-taught learning: transfer learning from unlabeled data, ICML, 2007</p>
	unlabeled	<p>Domain-adversarial training</p> <p>Zero-shot learning</p>	<p>Different from semi-supervised learning</p> <p>Self-taught Clustering</p> <p>Wenyuan Dai, Qiang Yang, Gui-Rong Xue, Yong Yu, "Self-taught clustering", ICML 2008</p>

---

# Explainable Machine Learning

# Why we need Explainable ML?

---

- Correct answers  $\neq$  Intelligent

Clever Hans  
performing  
in 1904



# Why we need Explainable ML?

---

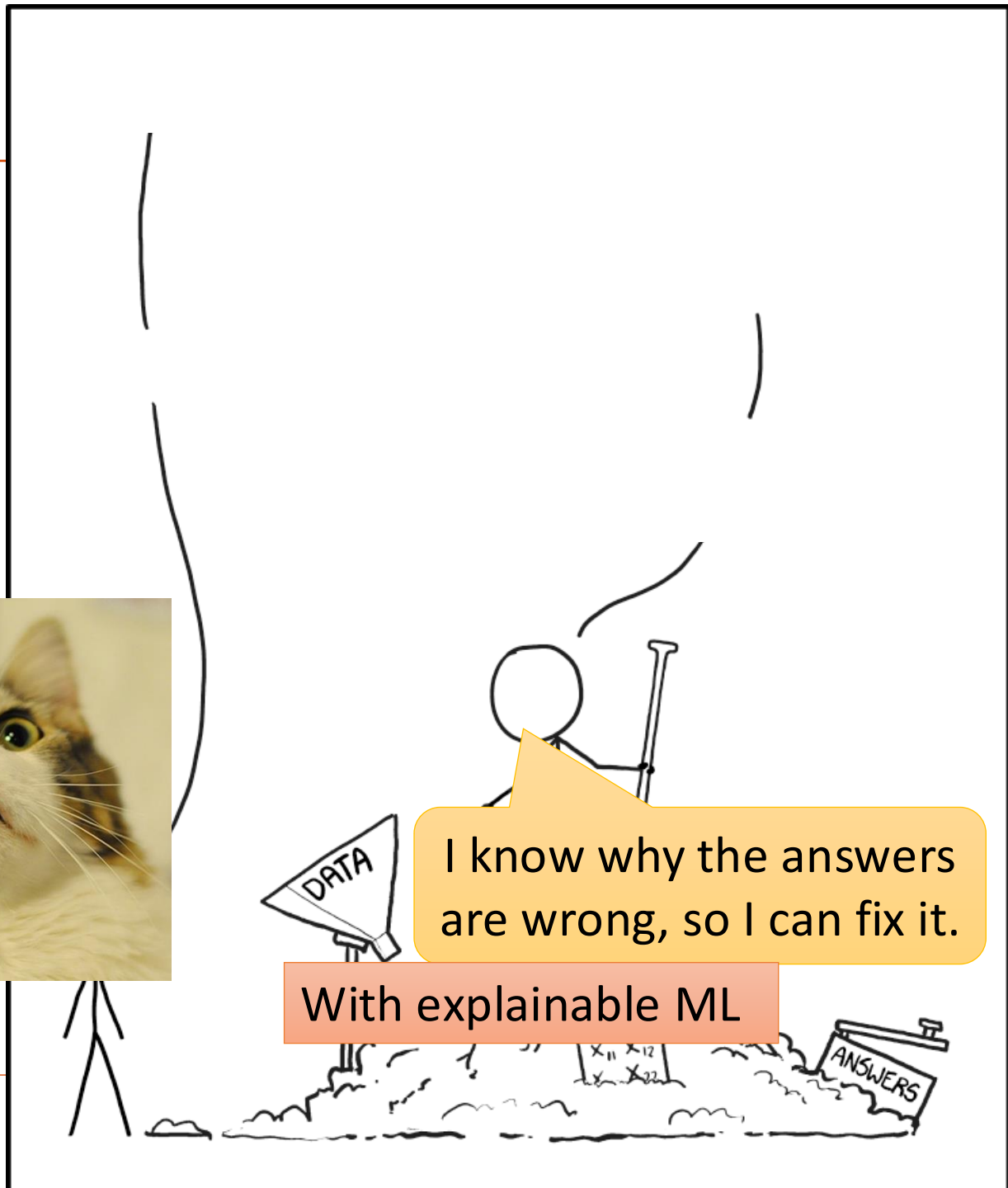
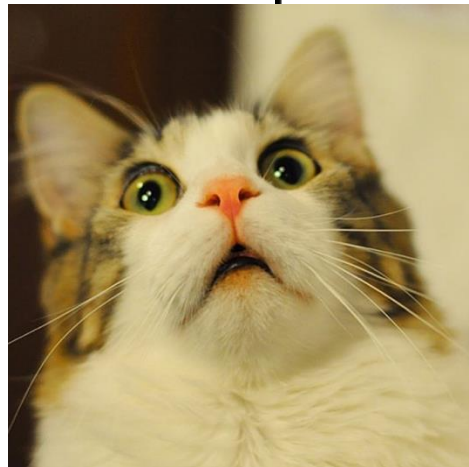
Loan issuers are required by law to explain their models.

Medical diagnosis model is responsible for human life. Can it be a black box?

If a model is used at the court, we must make sure the model behaves in a nondiscriminatory manner.

If a self-driving car suddenly acts abnormally, we need to explain why.

We can improve ML model based on explanation.



[https://www.explainkcd.com/wiki/index.php/1838:\\_Machine\\_Learning](https://www.explainkcd.com/wiki/index.php/1838:_Machine_Learning)

# Interpretable v.s. Powerful

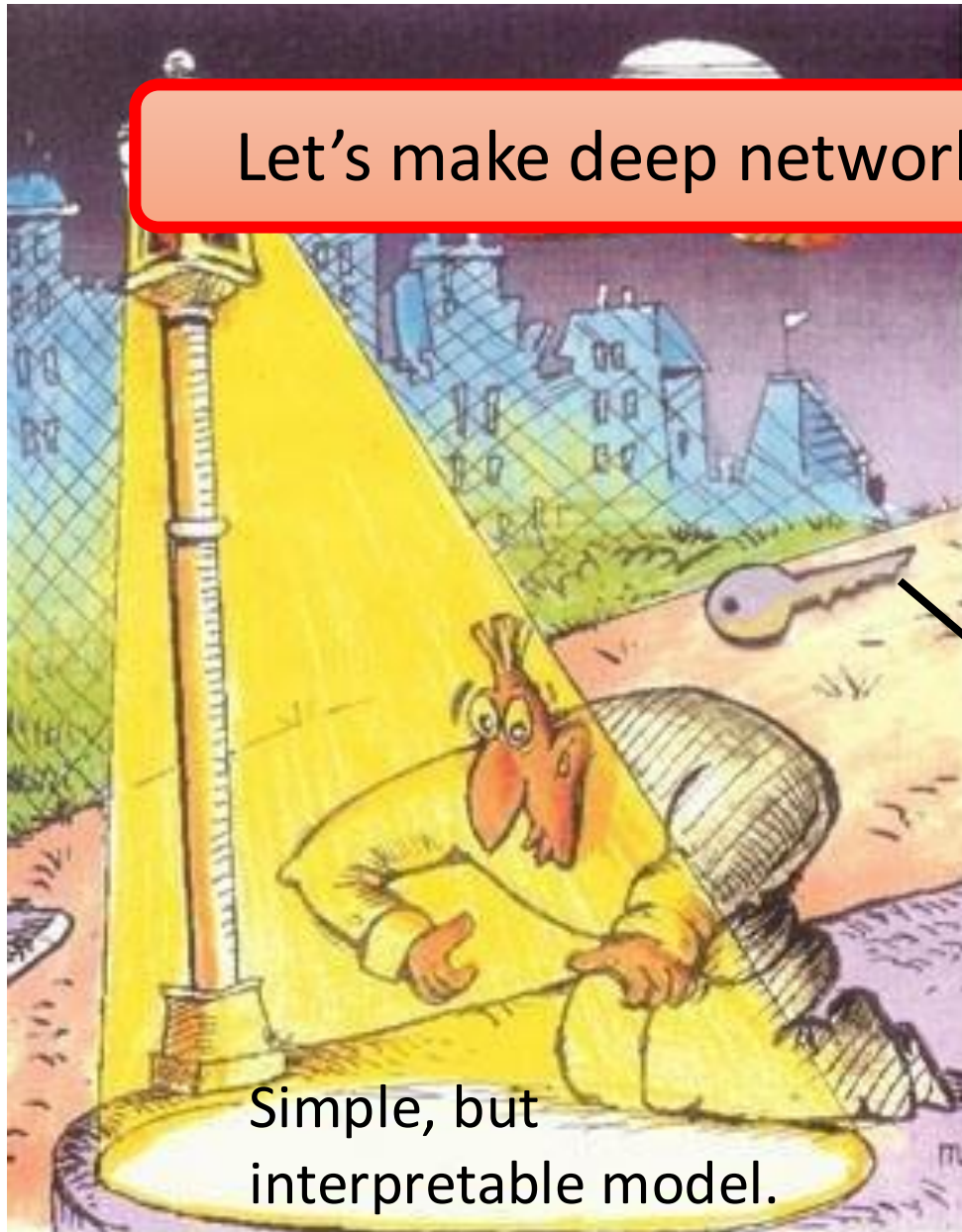
---

- Some models are intrinsically interpretable.
  - For example, linear model (from weights, you know the importance of features)
  - But not very powerful.
- Deep network is difficult to interpretable. Deep networks are black boxes ... but powerful than a linear model.

We don't want to use a more powerful model because it is a black box.

This is “cut the feet to fit the shoes.”

Let's make deep network explainable.



Simple, but interpretable model.

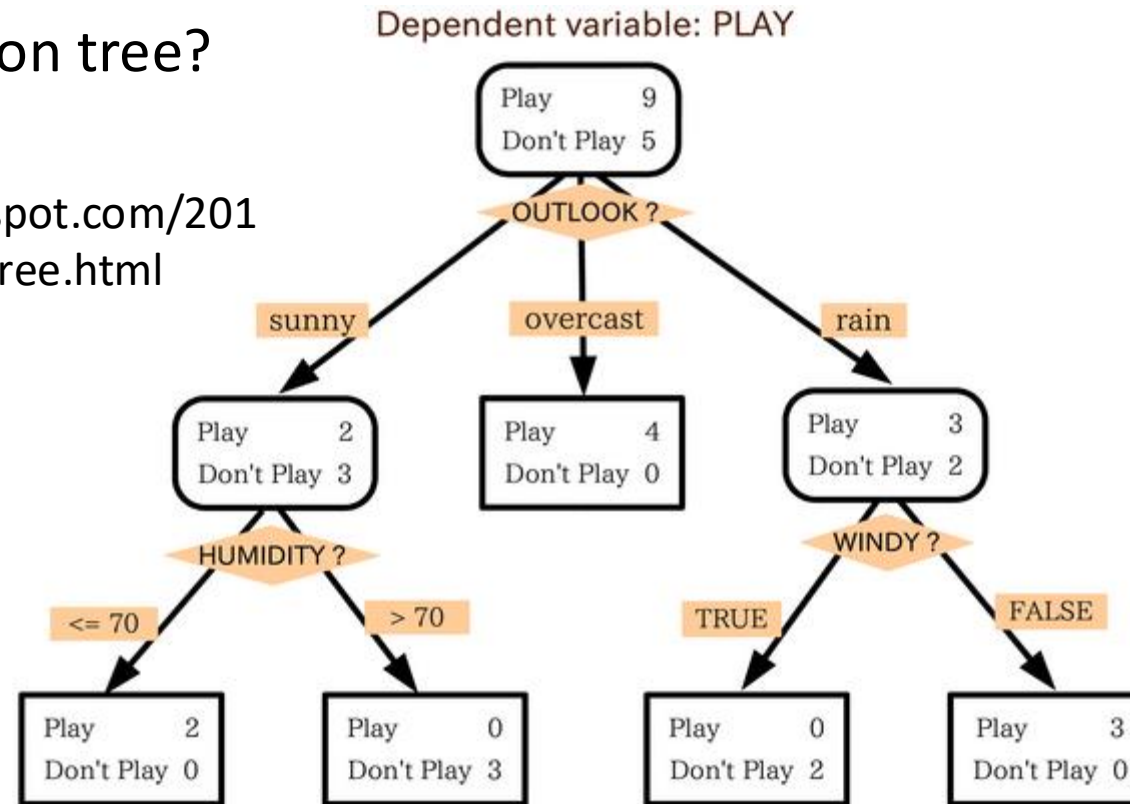
Powerful model

# Interpretable v.s. Powerful

- Are there some models interpretable and powerful at the same time?
- How about decision tree?

Source of image:

<https://mropengate.blogspot.com/2015/06/ai-ch13-2-decision-tree.html>

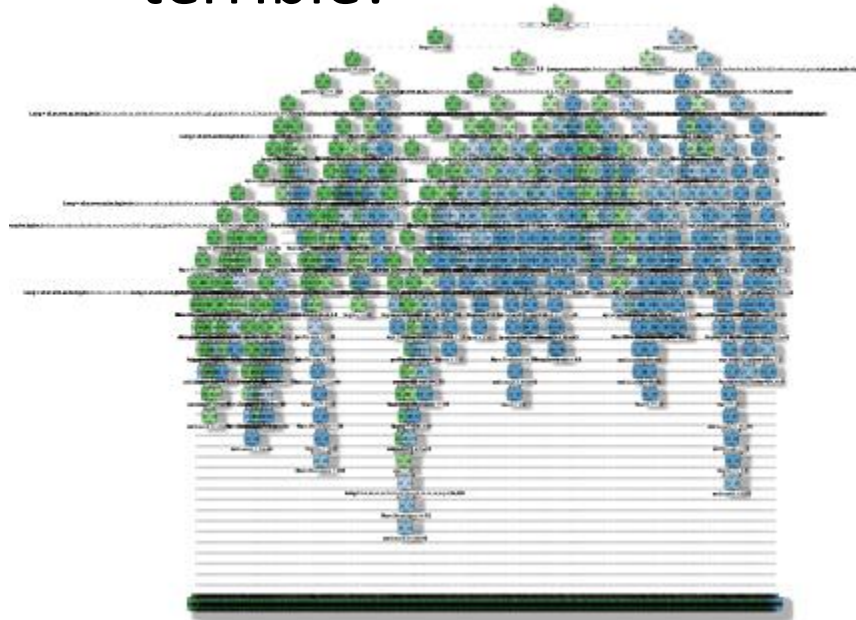


---

# Decision tree is all you need!?

# Interpretable v.s. Powerful

- A tree can still be terrible!



Rattle 2016-Aug-18 16:15:42 sklisarov

<https://stats.stackexchange.com/questions/230581/decision-tree-too-large-to-interpret>

- We use a forest!



# Goal of Explainable ML

---

- Completely know how an ML model works?
  - We do not completely know how brains work!
  - But we trust the decision of humans!

*The Copy Machine Study* (Ellen Langer, Harvard University)

“Excuse me, I have 5 pages. May I use the Xerox machine?”

60% accept

“Excuse me, I have 5 pages. May I use the Xerox machine,  
**because I’m in a rush?**”

94% accept

“Excuse me, I have 5 pages. May I use the Xerox machine,  
**because I have to make copies?**”

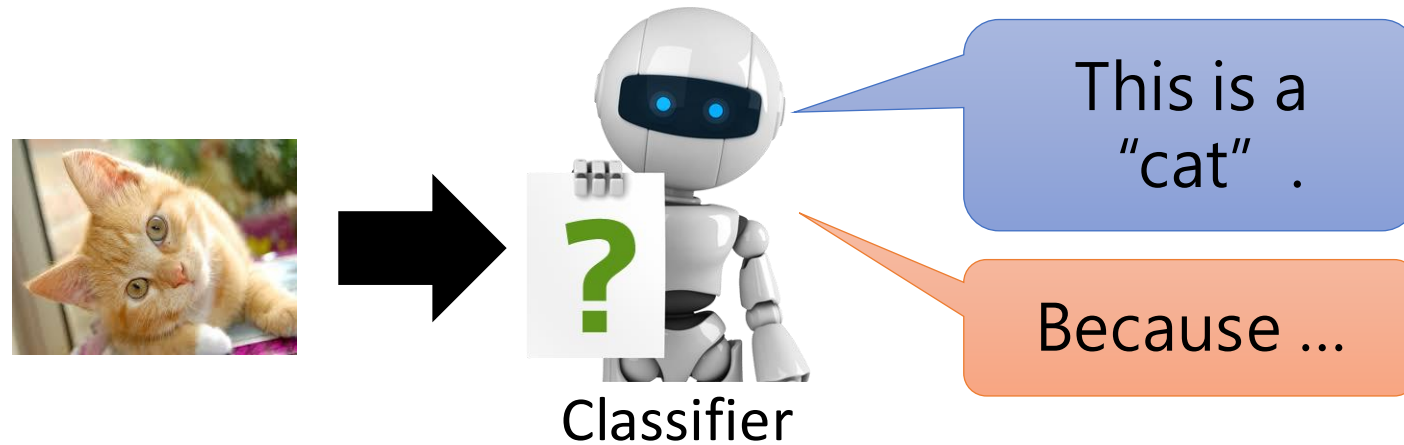
93% accept

<https://jamesclear.com/wp-content/uploads/2015/03/copy-machine-study-ellen-langer.pdf>

---

Make people (your customers,  
your boss, yourself) comfortable.

# Explainable ML



## Local Explanation

Why do you think this image is a cat?

## Global Explanation

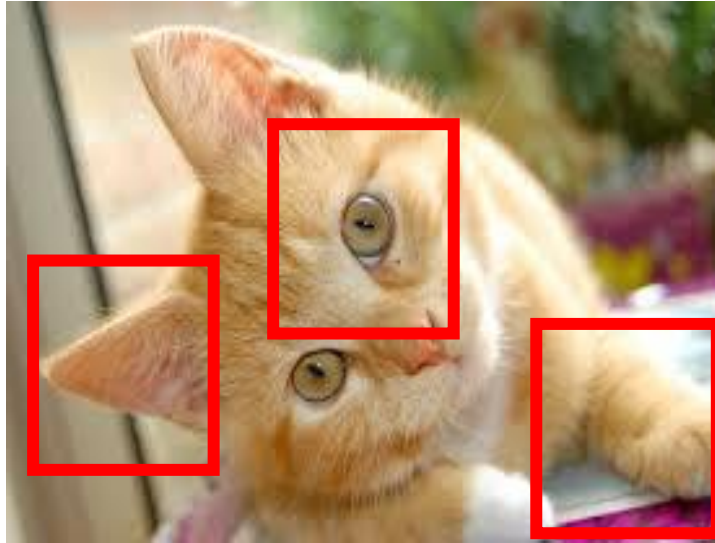
What does a “cat” look like?

(not referred to a specific image)

# Local Explanation: Explain the Decision

Questions: Why do you think this image  
is a cat?

# Which component is critical?



Which component is critical for making decision?

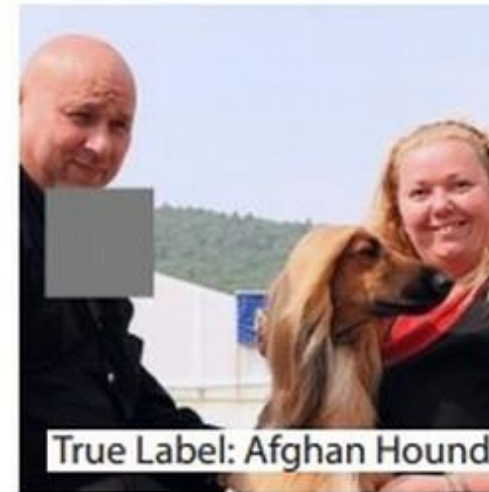
Object  $x$  →

Components:

$\{x_1, \dots, x_n, \dots, x_N\}$

Image: pixel, segment, etc.  
Text: a word

- Removing or modifying the components
  - Large decision change
- ➡ Important component



Reference: Zeiler, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional networks. In *Computer Vision—ECCV 2014* (pp. 818-833)

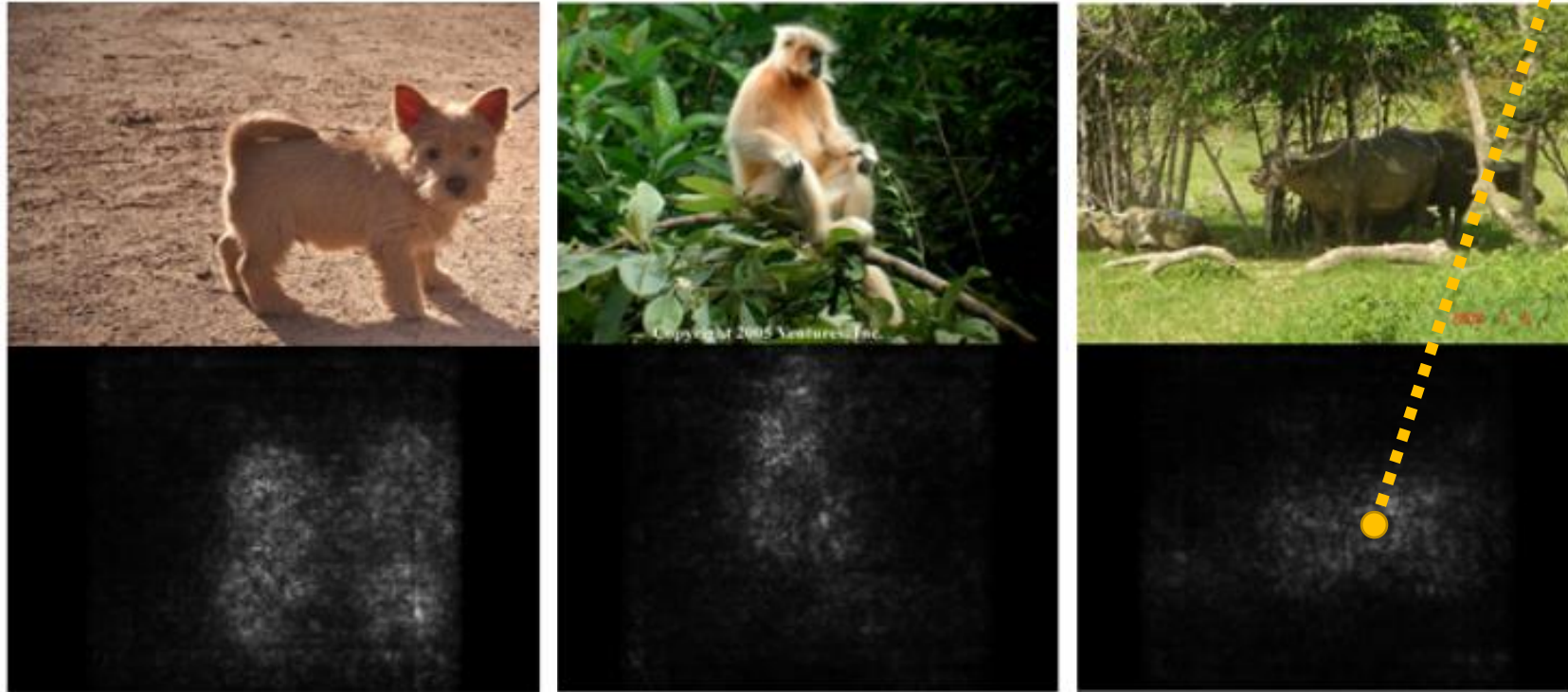
$$\{x_1, \dots, x_n, \dots, x_N\} \longrightarrow \{x_1, \dots, x_n + \Delta x, \dots, x_N\}$$

pixels

$$e \longrightarrow e + \Delta e$$

loss of an example (the difference between model output and ground truth)

$$\left| \frac{\Delta e}{\Delta x} \right| \longrightarrow \left| \frac{\partial e}{\partial x_n} \right|$$



[Saliency Map](#)

Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR, 2014

# Case Study: Pokémon v.s. Digimon



<https://medium.com/@tyreeostevenson/teaching-a-computer-to-classify-anime-8c77bc89b881>

# Task

Pokémon images: <https://www.Kaggle.com/kvpratama/pokemon-images-dataset/data>

Digimon images:  
<https://github.com/DeathReaper0965/Digimon-Generator-GAN>



Pokémon



Digimon

Testing  
Images:



# Experimental Results

```
model = Sequential()
model.add(Conv2D(32, (3, 3), padding='same', input_shape=(120,120,3)))
model.add(Activation('relu'))
model.add(Conv2D(32, (3, 3)))
model.add(Activation('relu'))
model.add(MaxPooling2D(pool_size=(2, 2)))

model.add(Conv2D(64, (3, 3), padding='same'))
model.add(Activation('relu'))
model.add(Conv2D(64, (3, 3)))
model.add(Activation('relu'))
model.add(MaxPooling2D(pool_size=(2, 2)))

model.add(Conv2D(256, (3, 3), padding='same'))
model.add(Activation('relu'))
model.add(Conv2D(256, (3, 3)))
model.add(Activation('relu'))
model.add(MaxPooling2D(pool_size=(2, 2)))

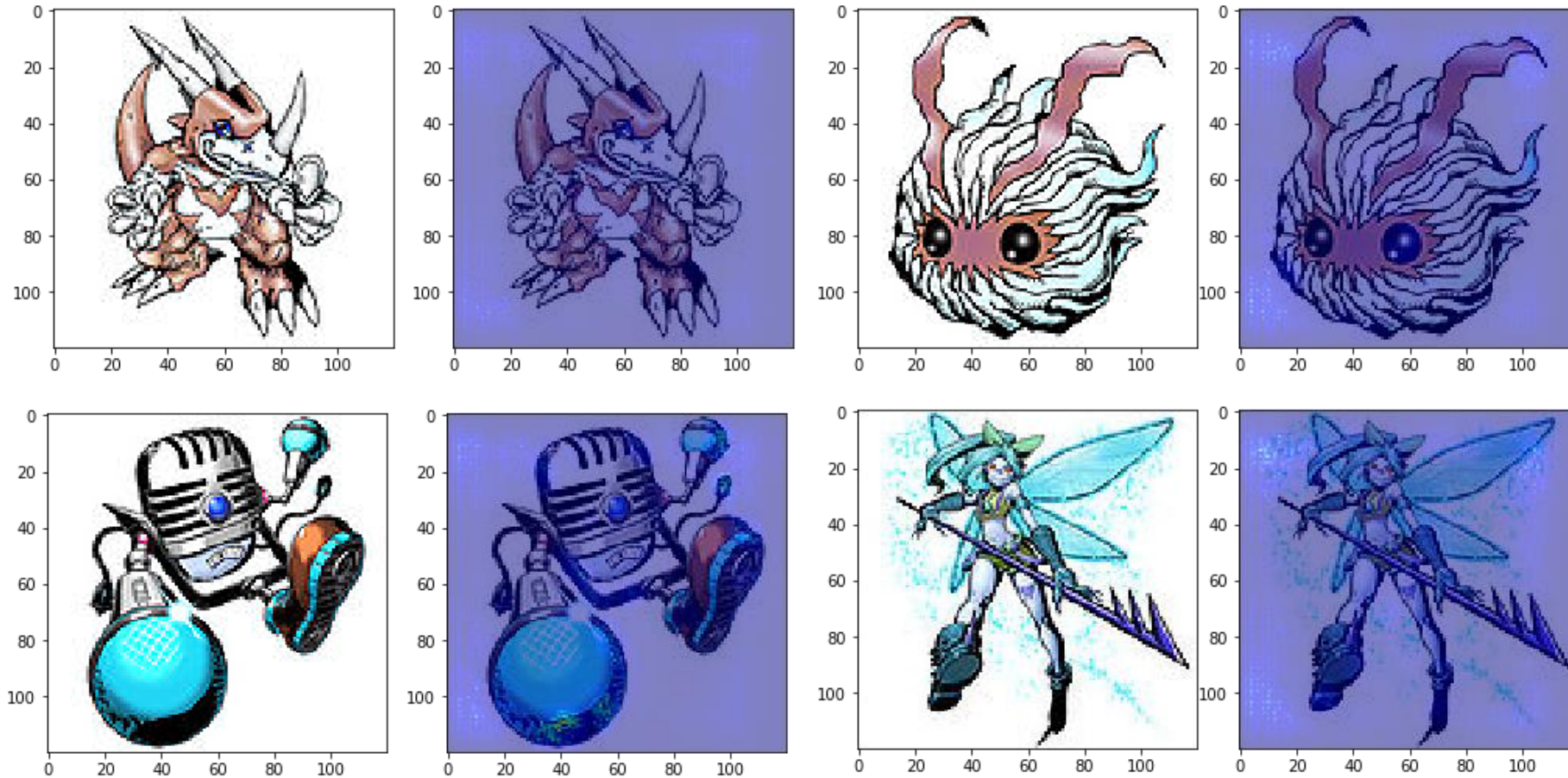
model.add(Flatten())
model.add(Dense(1024))
model.add(Activation('relu'))
model.add(Dense(2))
model.add(Activation('softmax'))
```

Training Accuracy: 98.9%

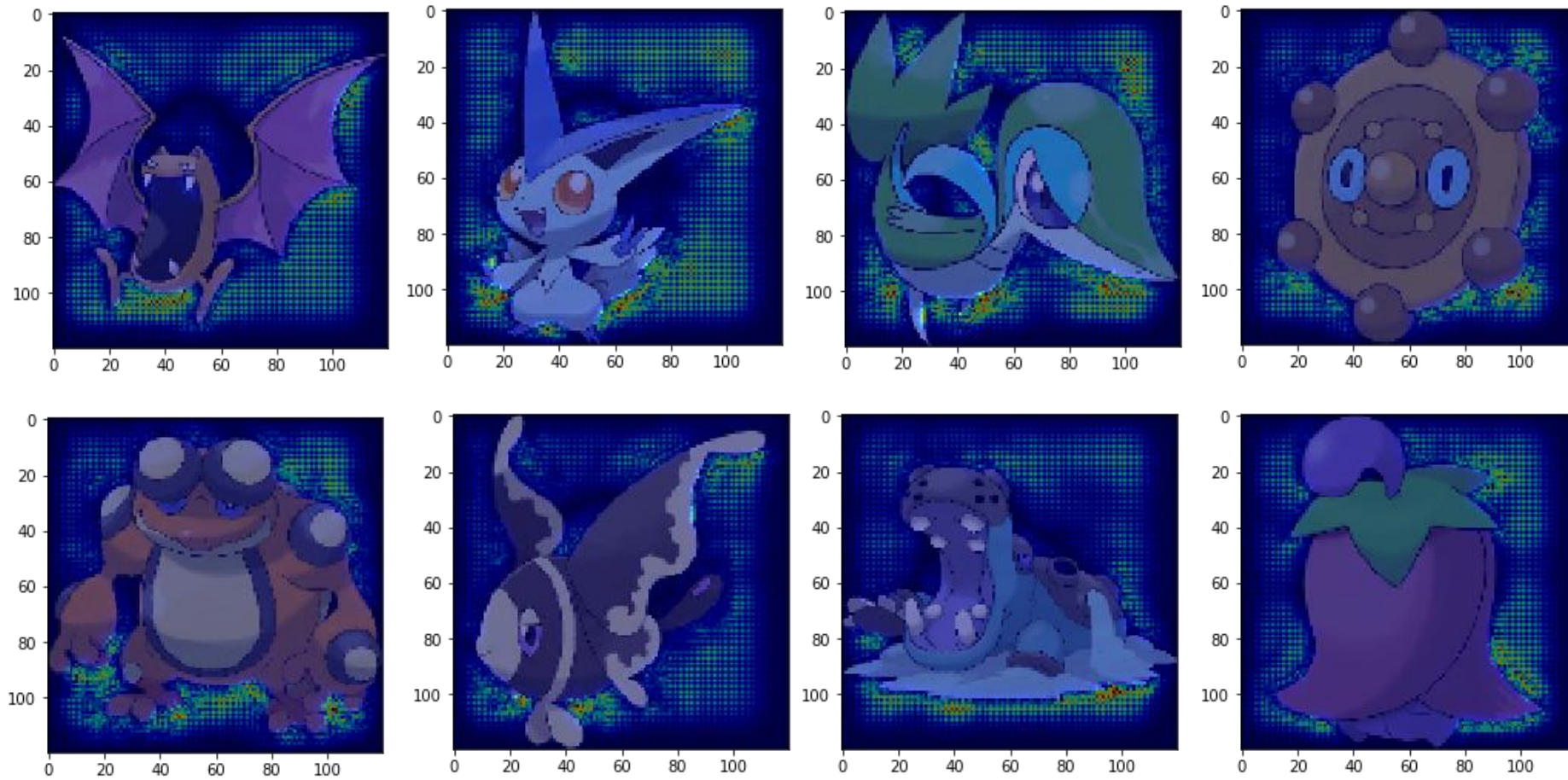
Testing Accuracy: 98.4%

Amazing!!!!!!

# Saliency Map



# Saliency Map



# What Happened?

- All the images of Pokémon are PNG, while most images of Digimon are JPEG.



png files have transparent background

loading the files

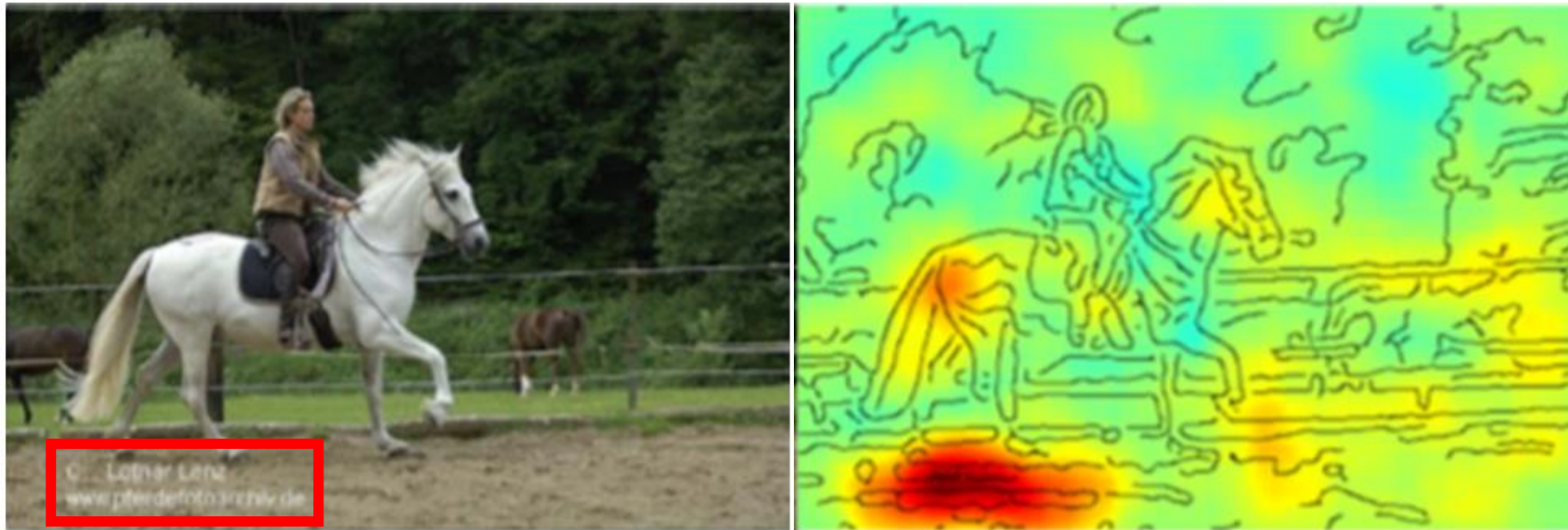


transparent background becomes black

Machine discriminates Pokémon and Digimon based on the background colors.

# More Examples ...

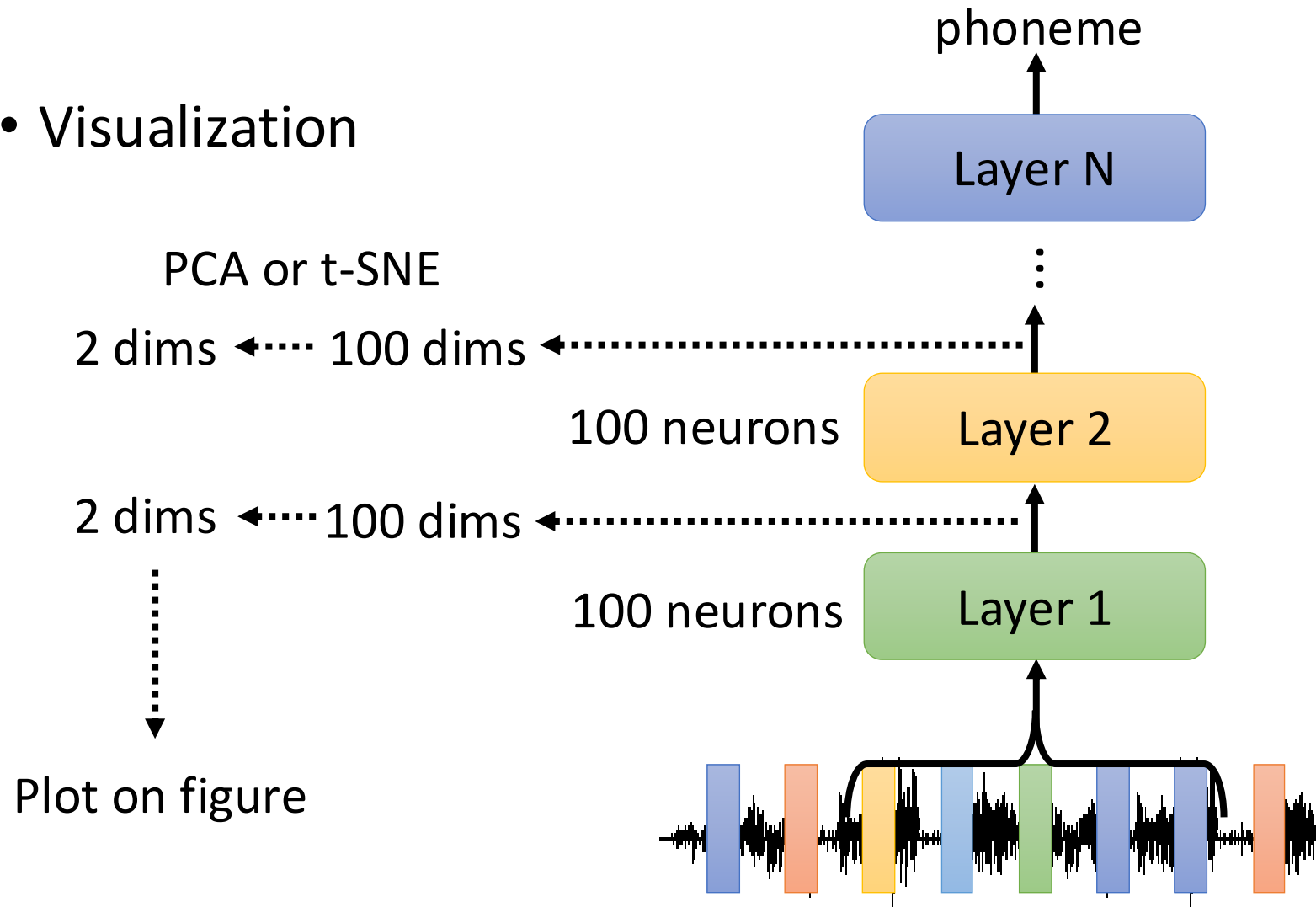
- PASCAL VOC 2007 data set



This slide is from: GCPR 2017 Tutorial — W. Samek & K.-R. Müller

# How a network processes the input data?

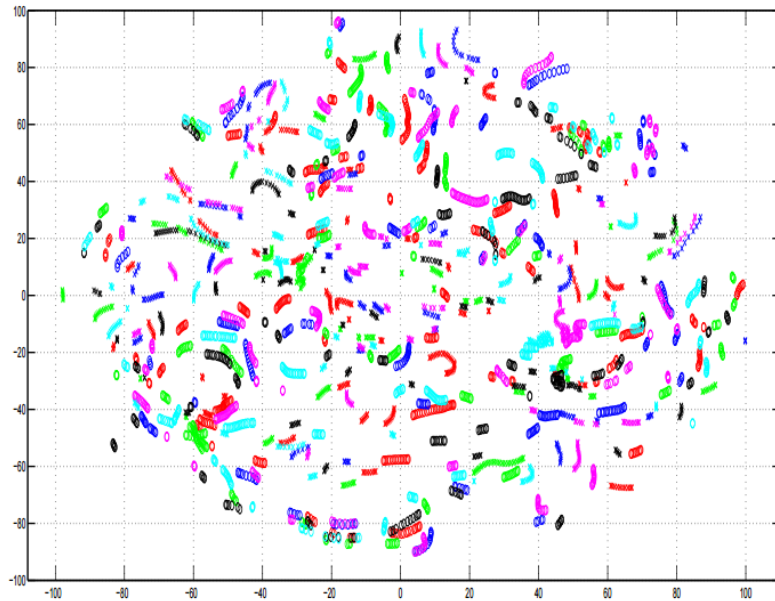
- Visualization



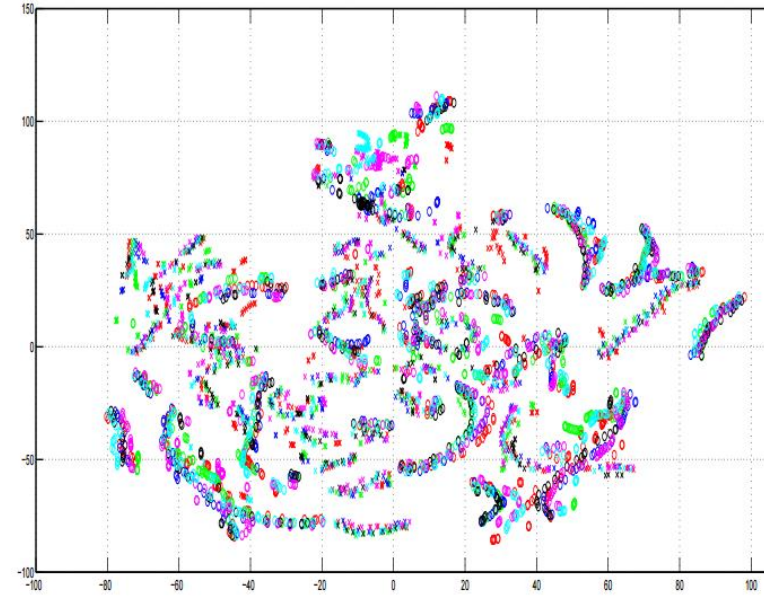
# How a network processes the input data?

A. Mohamed, G. Hinton, and G. Penn,  
“Understanding how Deep Belief Networks Perform  
Acoustic Modelling,” in ICASSP, 2012.

- Visualization  
Colors: speakers



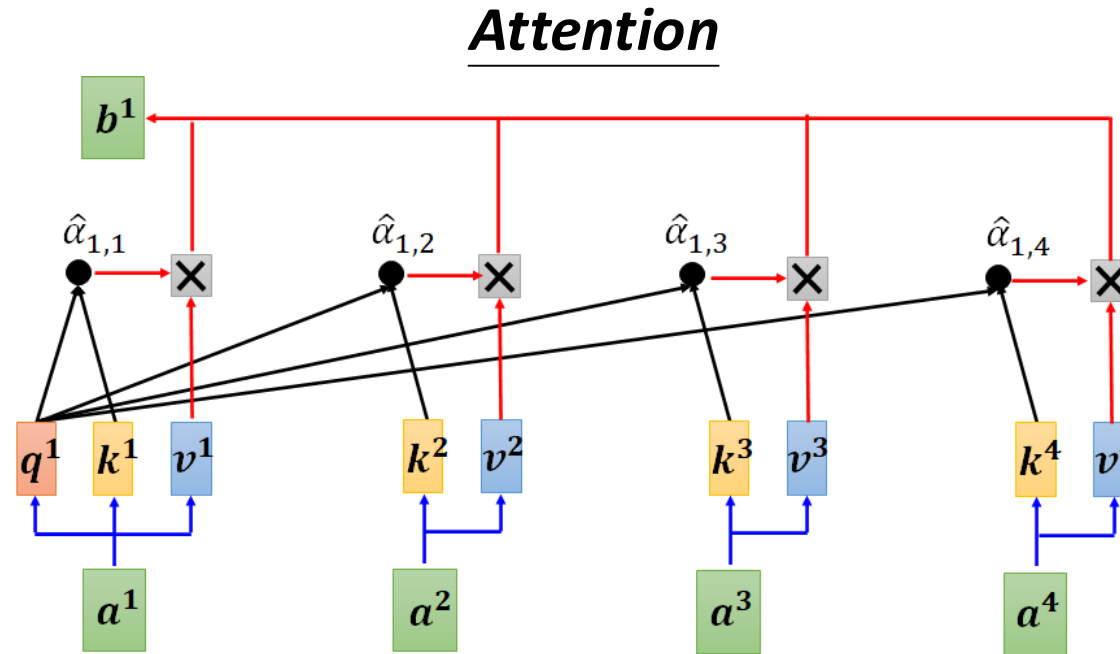
Input Acoustic Feature (MFCC)



8-th Hidden Layer

# How a network processes the input data?

- Visualization



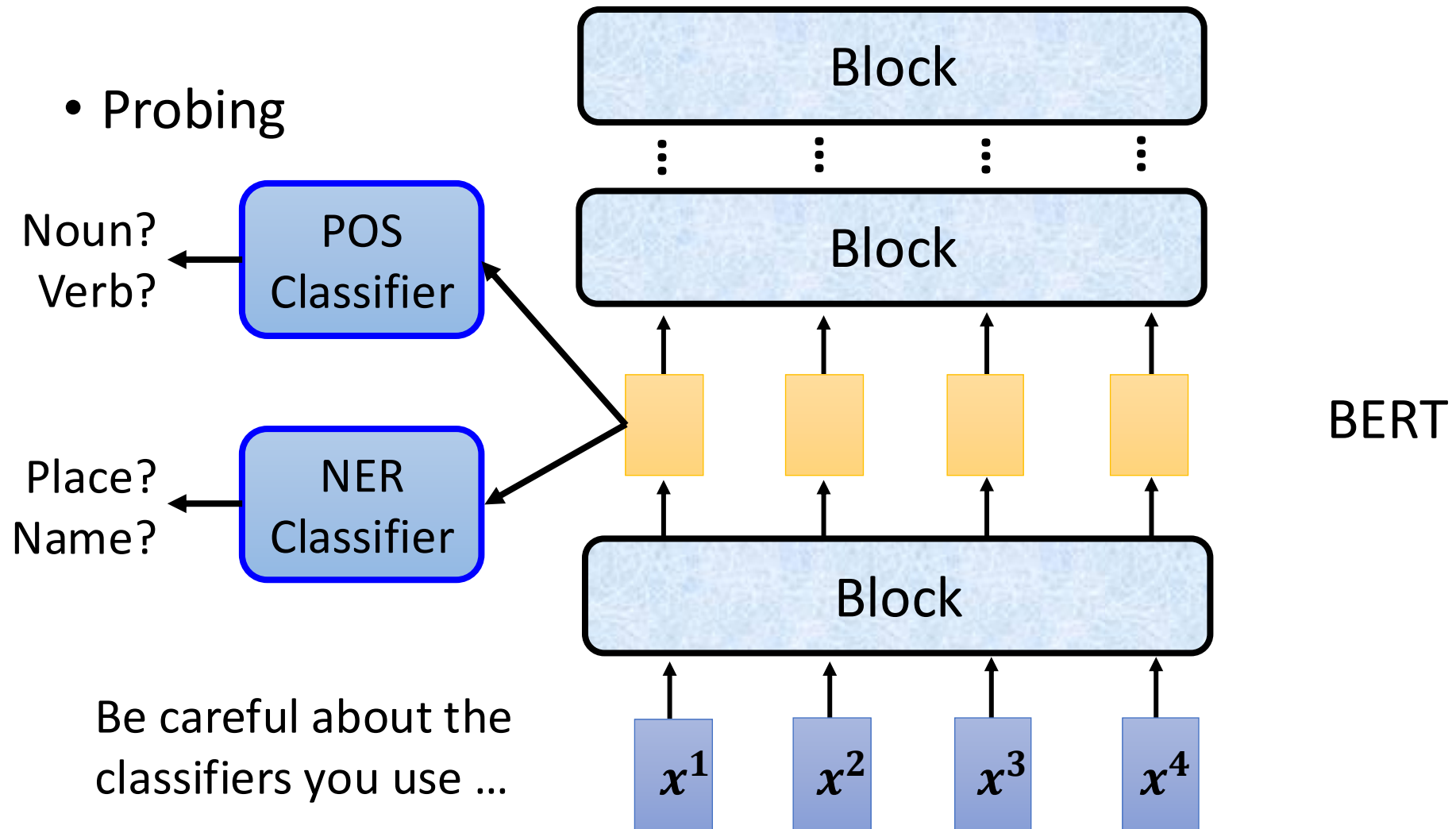
Attention is not Explanation

<https://arxiv.org/abs/1902.10186> (02. 2019)

Attention is not not Explanation

<https://arxiv.org/abs/1908.04626> (05. 2019)

# How a network processes the input data?



Be careful about the classifiers you use ...

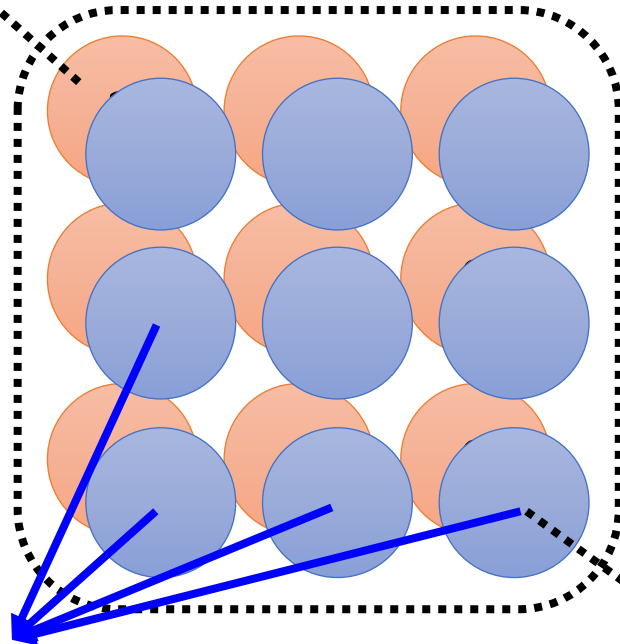
---

# Global Explanation: Explain the whole Model

Question: What does a “cat” look like?

# What does a filter detect?

output of filter 2



Large values

output of filter 1

➔ Image  $X$  contains the patterns filter 1 can detect.

unknown

image  $X$

input

filters

Convolution

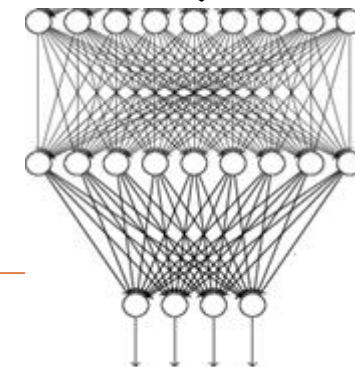
filters

Convolution

Max Pooling

Max Pooling

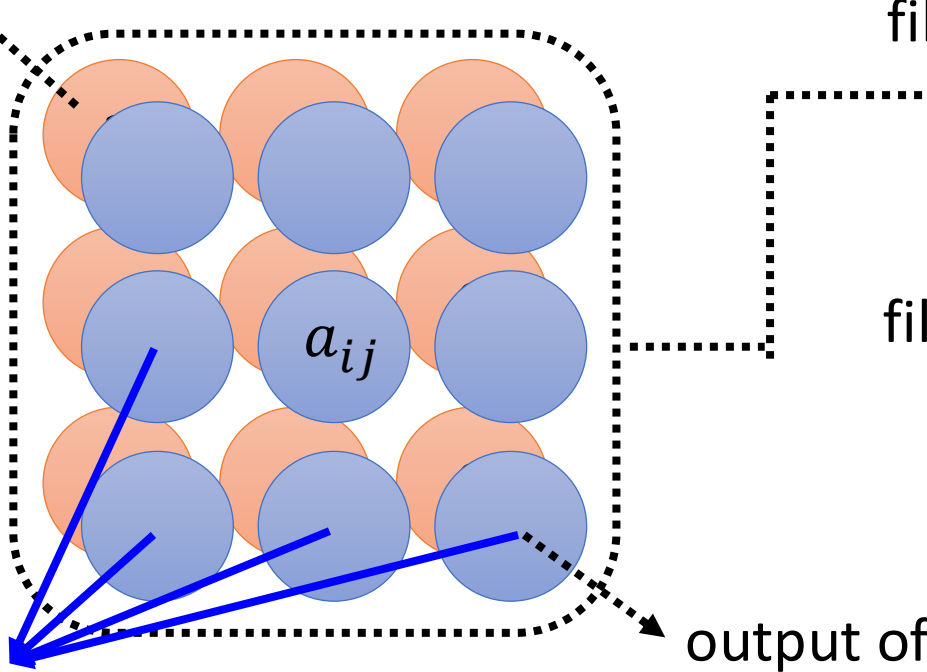
flatten



Let's **create** an image including the patterns.

# What does a filter detect?

output of filter 2



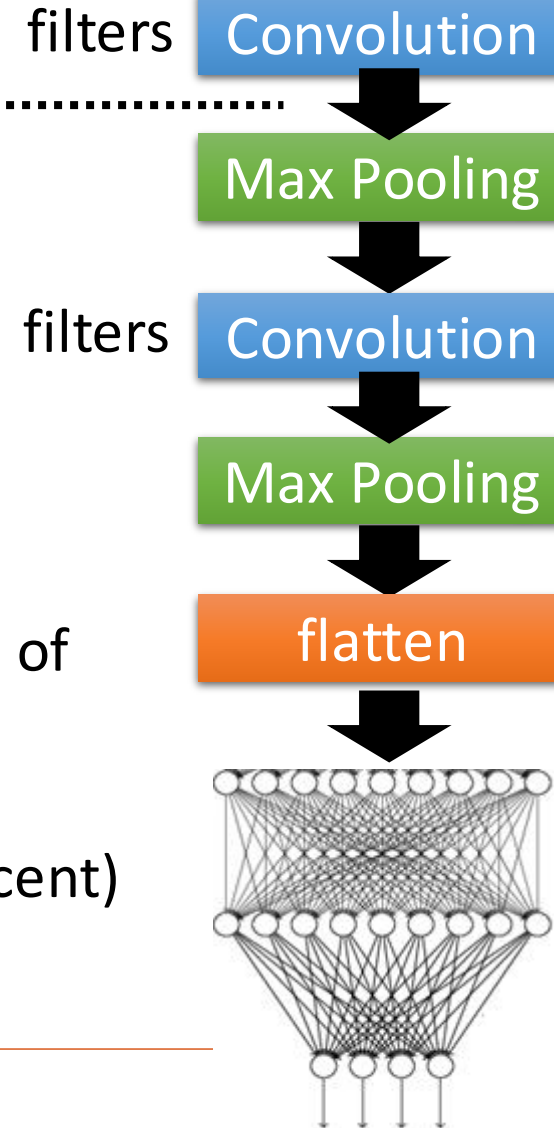
Large values

output of filter 1

$$X^* = \mathit{arg} \max_X \sum_i \sum_j a_{ij} \quad (\text{gradient ascent})$$

The image contains the patterns filter 1 can detect.

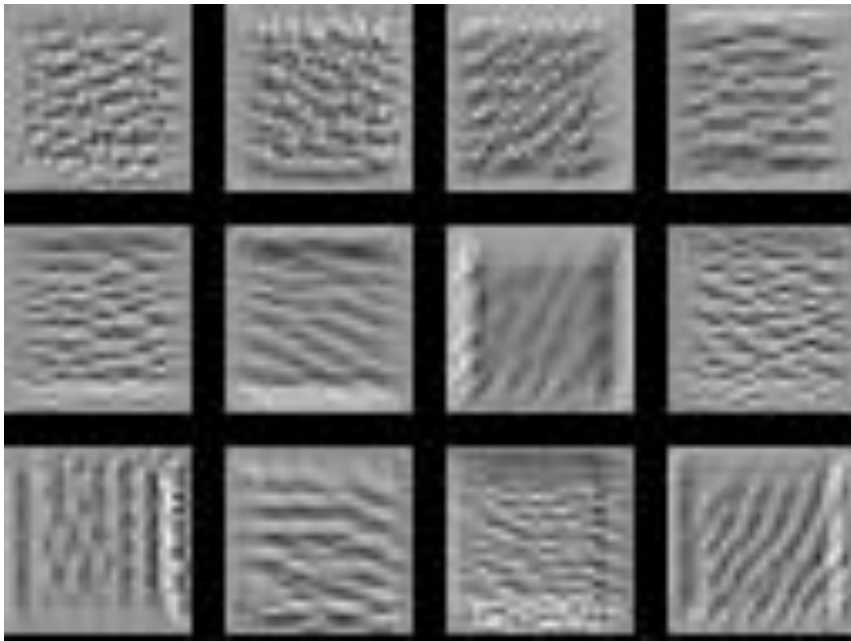
unknown image  $X$  input



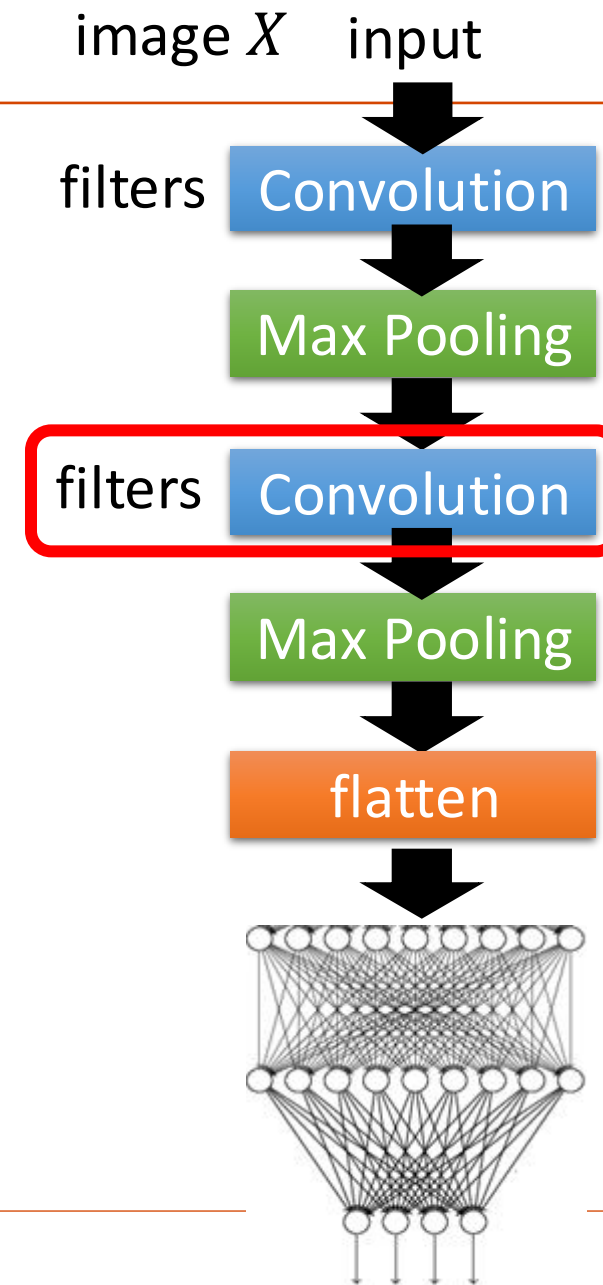
# What does a filter detect?

E.g., Digit classifier

$X^*$  for each filter



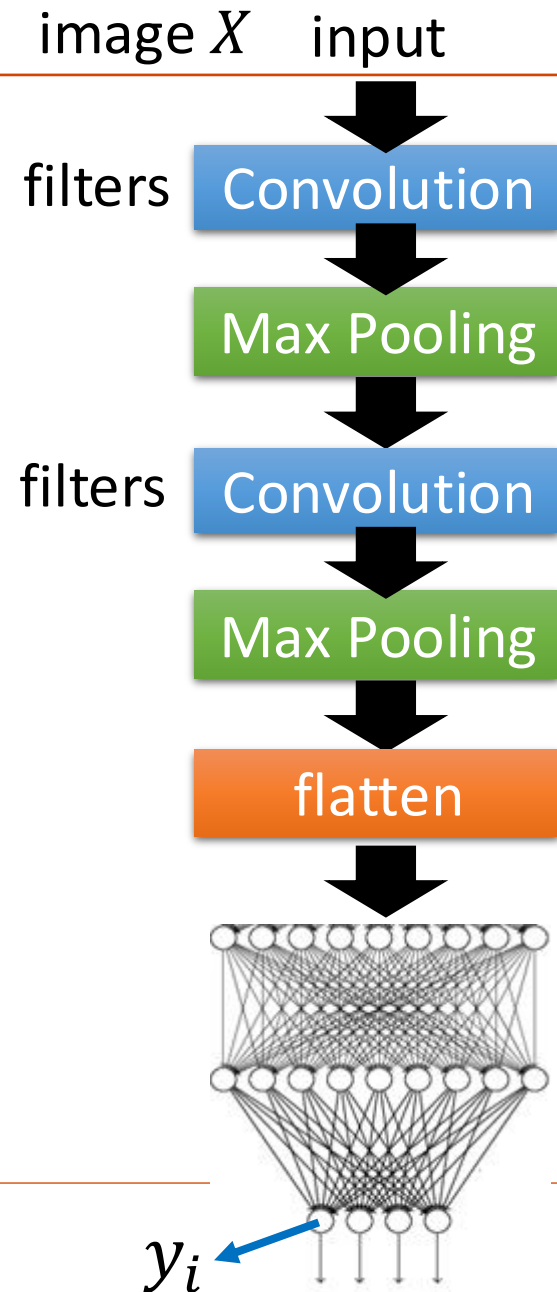
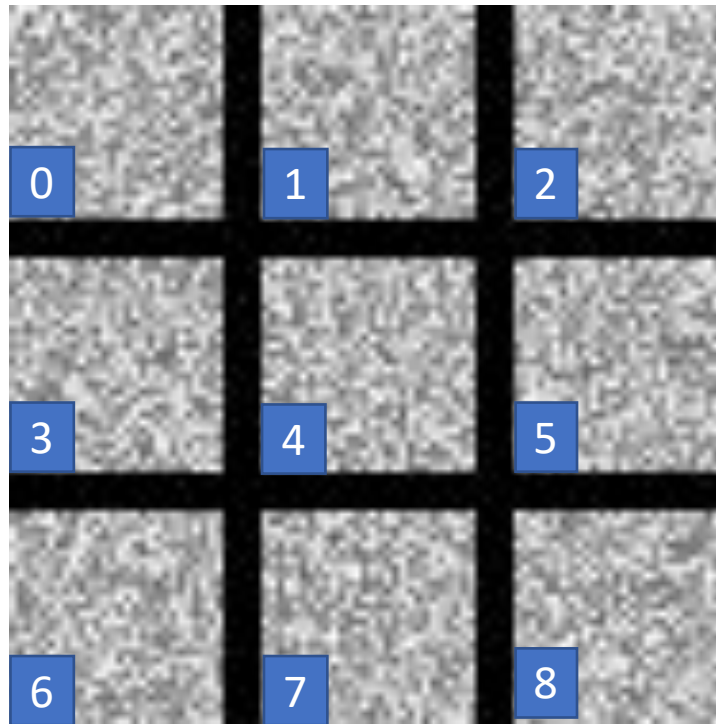
Use NMIST handwritten database



# What does a digit look like for CNN?

E.g., Digit classifier

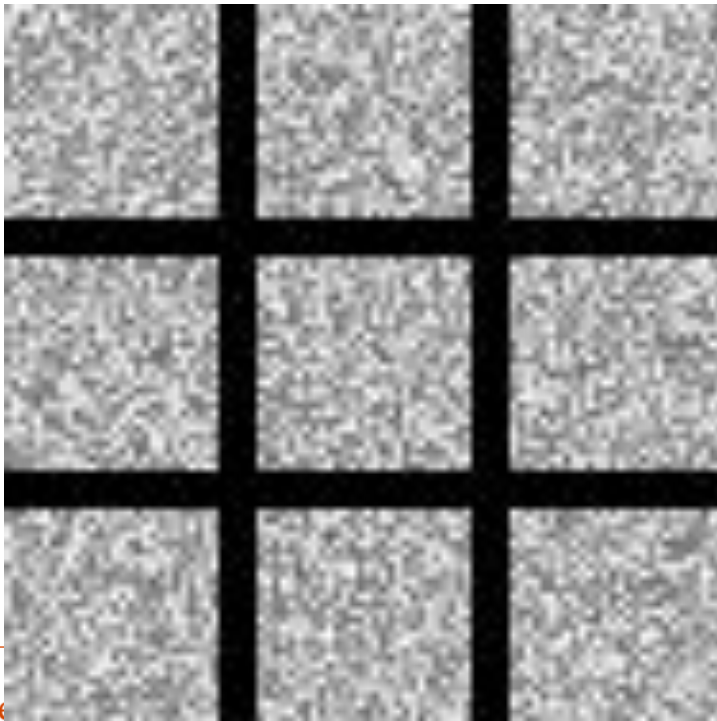
$$X^* = \underset{X}{\operatorname{arg\,max}} y_i \quad \text{Can we see digits?}$$



# What does a digit look like for CNN?

Find the image that maximizes class probability

$$X^* = \arg \max_X y_i$$

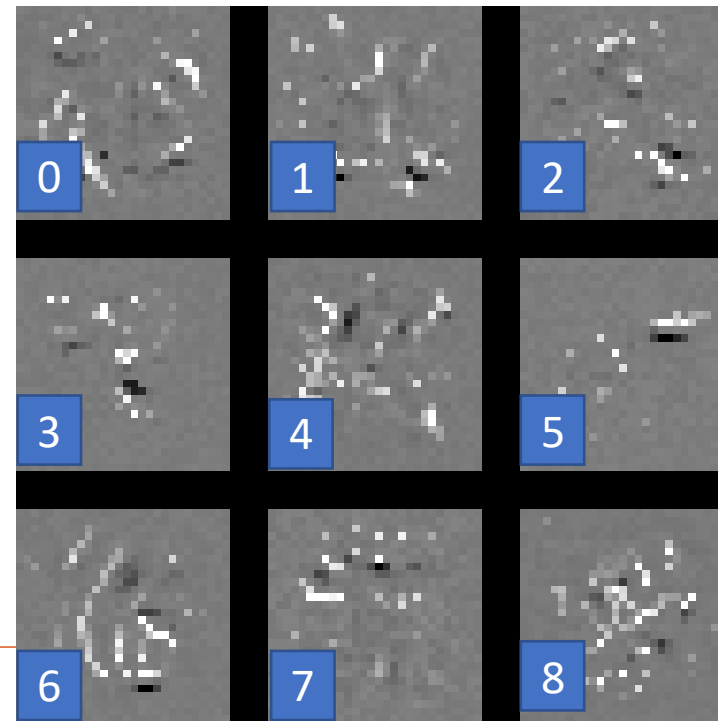


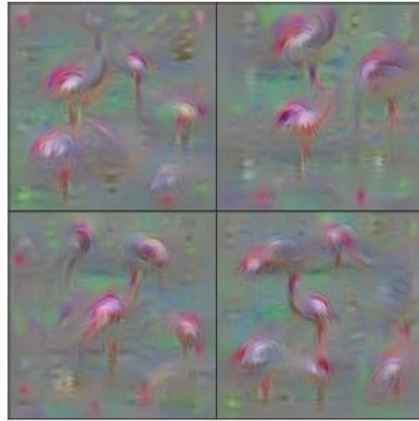
The image should look like a digit.

$$X^* = \arg \max_X y_i + \underline{R(X)}$$

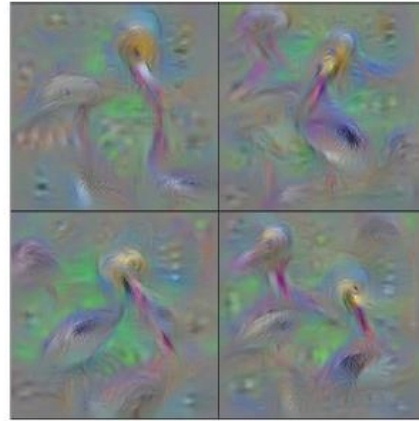
$$R(X) = - \sum_{i,j} |X_{ij}|$$

How likely  
X is a digit

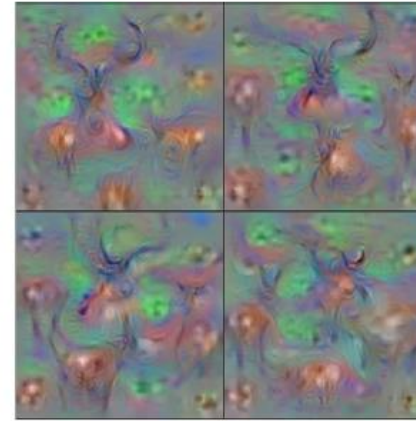




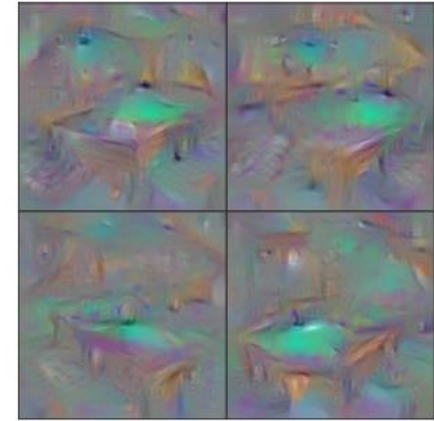
Flamingo



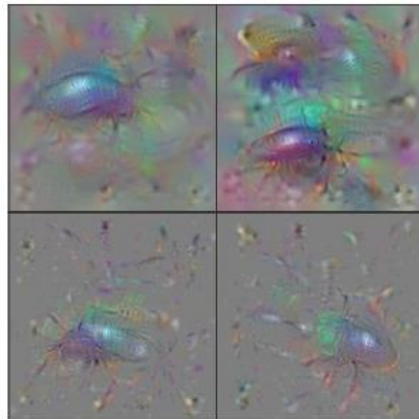
Pelican



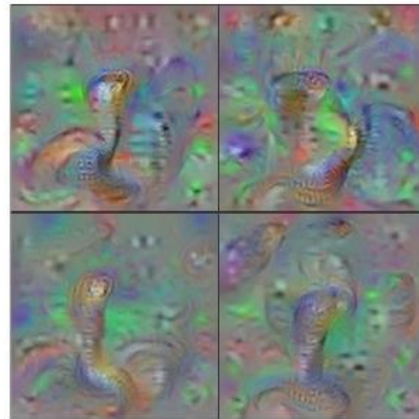
Hartebeest



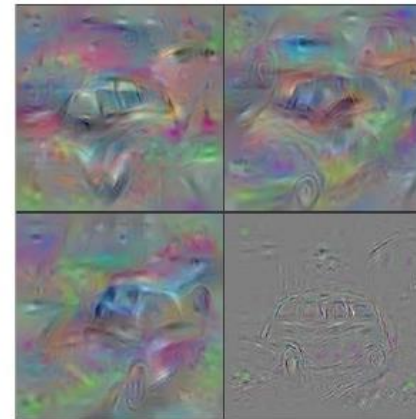
Billiard Table



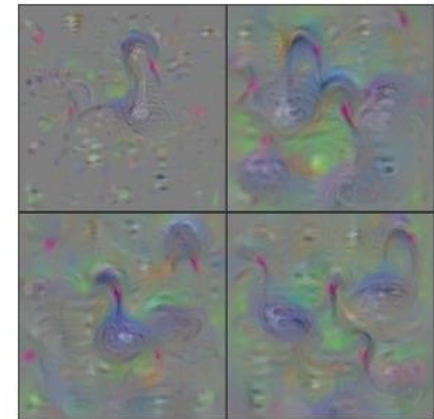
Ground Beetle



Indian Cobra



Station Wagon



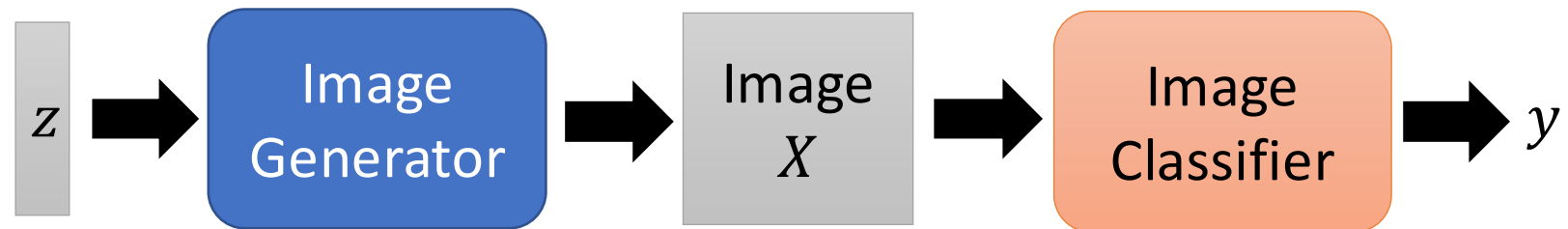
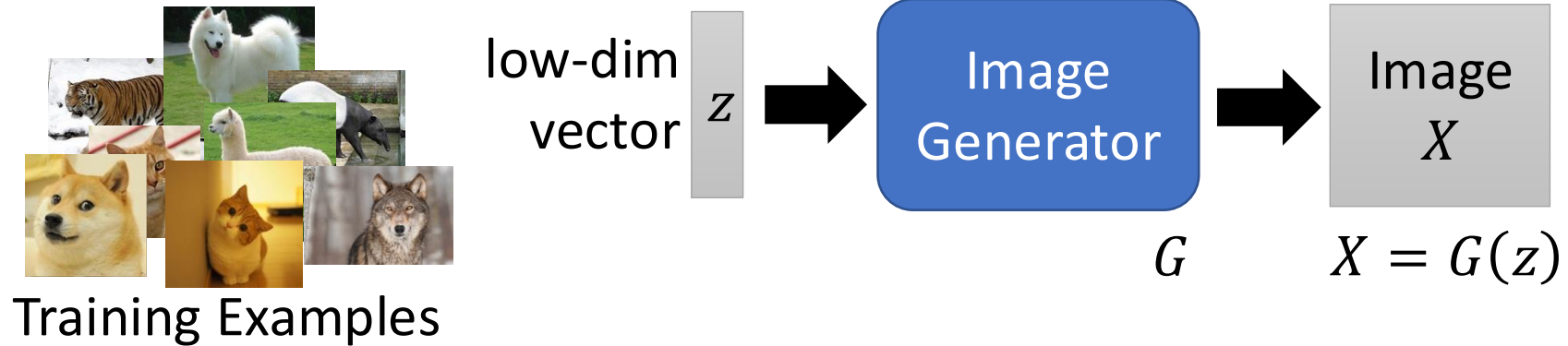
Black Swan

With several regularization terms, and hyperparameter tuning .....

<https://arxiv.org/abs/1506.06579> (2015)

# Constraint from Generator

- Training a generator



$$X^* = \arg \max_X y_i \longrightarrow z^* = \arg \max_z y_i$$

Show image:

$$X^* = G(z^*)$$

# Plug & Play Generative Networks: Conditional Iterative Generation of Images in Latent Space

(2017)



redshank

ant

monastery



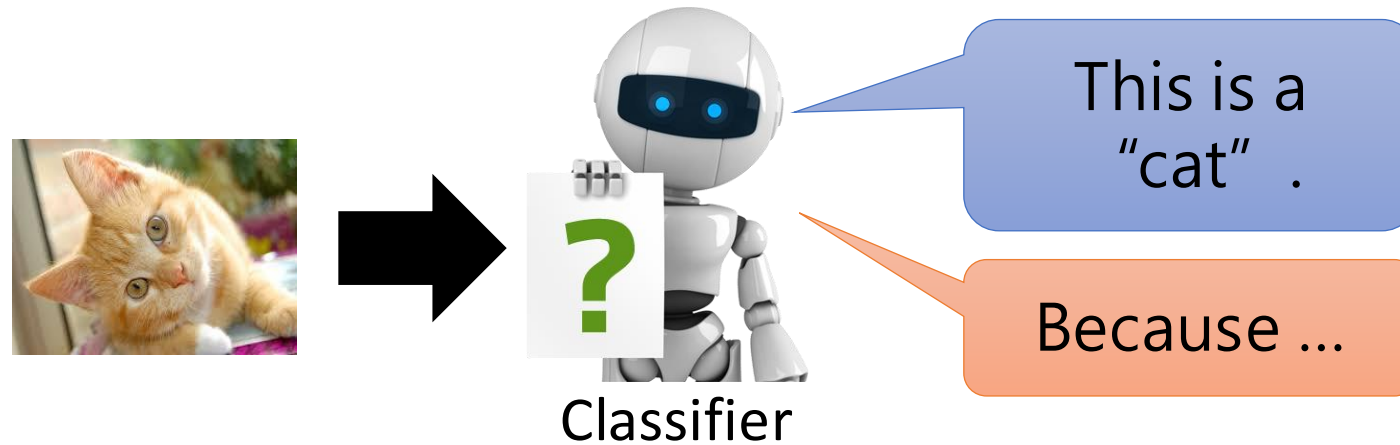
volcano

[https://arxiv.org/abs/  
1612.00005](https://arxiv.org/abs/1612.00005)

Syracuse University

# Conclusion Remarks

---



## **Local Explanation**

Why do you think this image is a cat?

## **Global Explanation**

What does a "cat" look like?

(not referred to a specific image)

---

Any Questions?

bidong@syr.edu