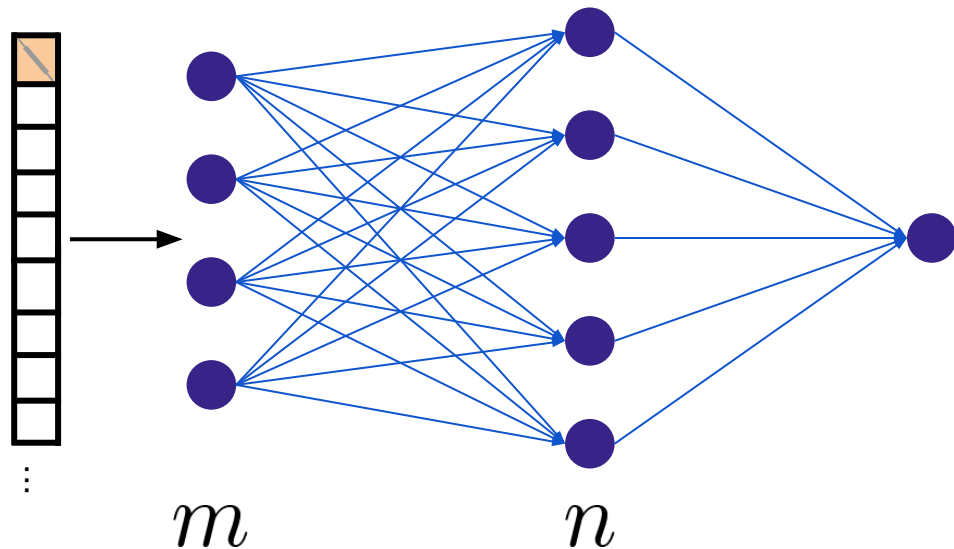# Transfer Learning

BOSTON UNIVERSITY
MACHINE INTELLIGENCE
COMMUNITY

Duy Nguyen, Justin Chen
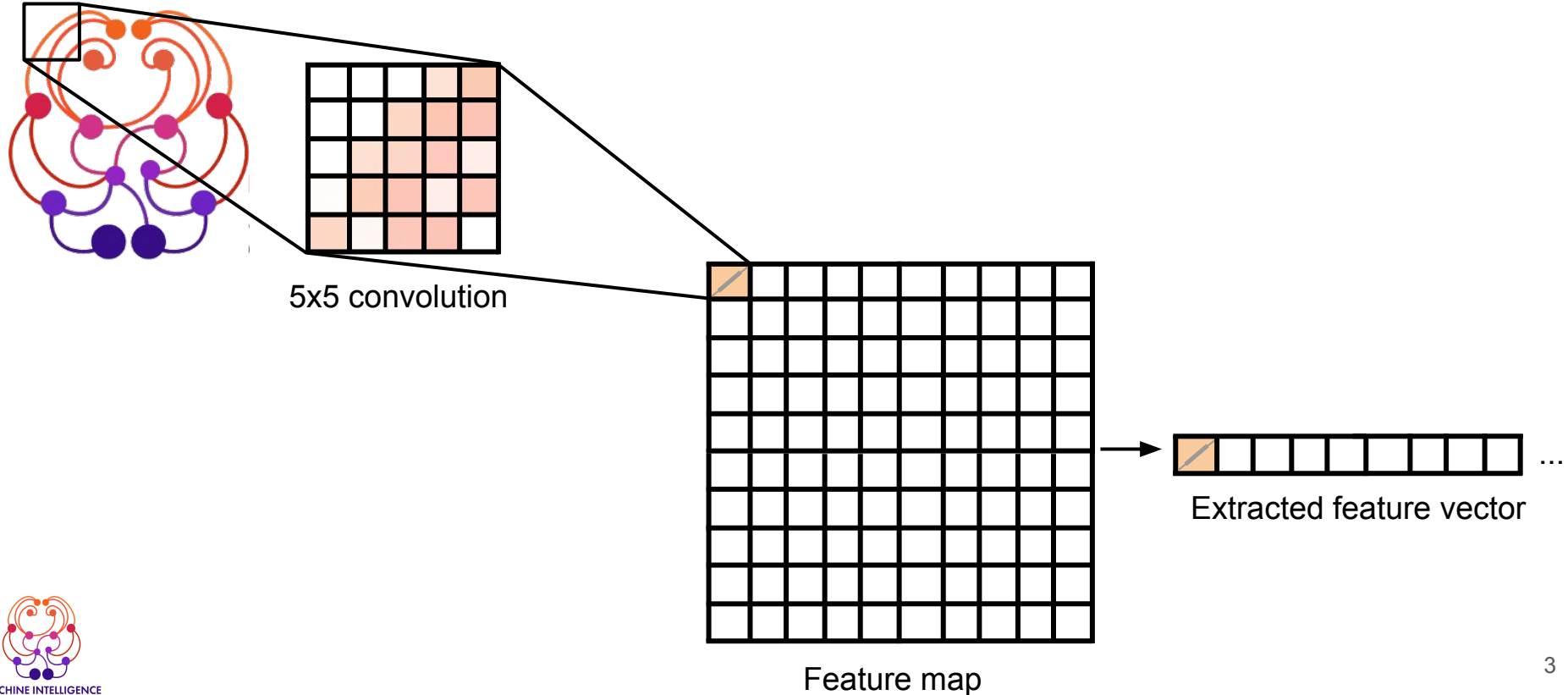Oct. 31, 2017

# CNN: Fully-connected Layer



$$\begin{bmatrix} \theta_{11}^1 & \theta_{12}^1 & \theta_{13}^1 & \dots & \theta_{1m}^1 \\ \theta_{21}^1 & \theta_{22}^1 & \theta_{23}^1 & \dots & \theta_{2m}^1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \theta_{n1}^1 & \theta_{n2}^1 & \theta_{n3}^1 & \dots & \theta_{nm}^1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$$

$$\begin{bmatrix} \theta_{11}^2 & \theta_{12}^2 & \theta_{13}^2 & \dots & \theta_{1n}^2 \end{bmatrix} \sigma(\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix})$$

$$= \begin{bmatrix} h_1 \end{bmatrix}$$
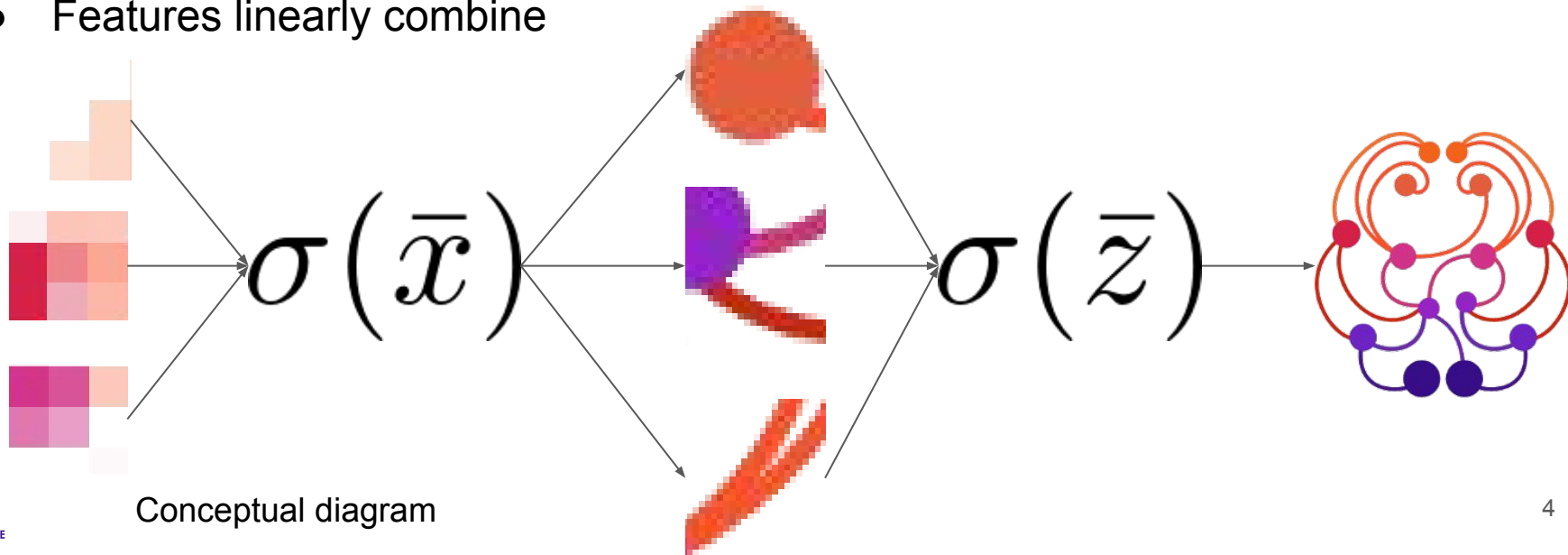
$m$

$n$

Ignoring bias values for simplification

In this example, we're doing regression

MACHINE INTELLIGENCE COMMUNITY

2

# CNN: Automatic Feature Extraction

5x5 convolution

Feature map

Extracted feature vector

MACHINE INTELLIGENCE
COMMUNITY

3

# Feature Hierarchy

- Learned features become progressively more complex throughout the network
- Earlier layers contain more primitive features
- Features linearly combine
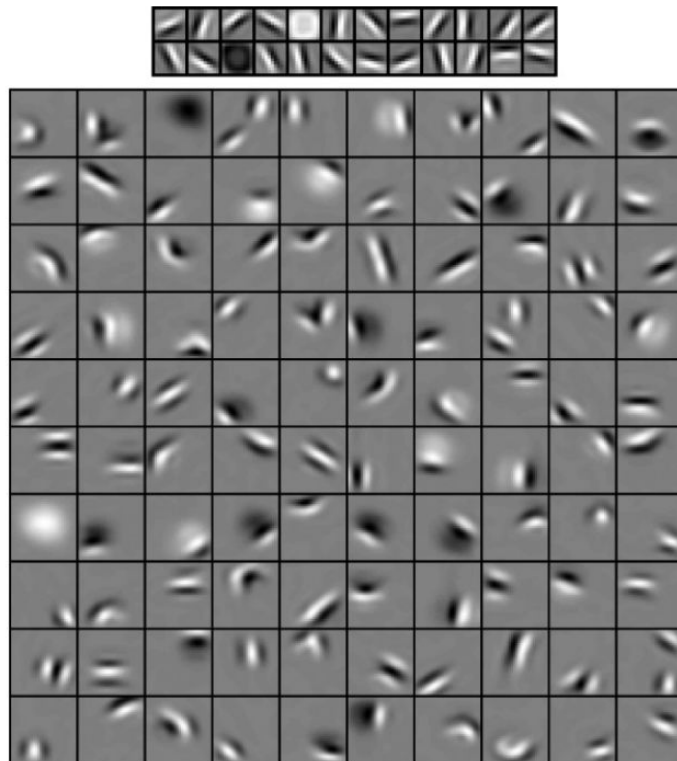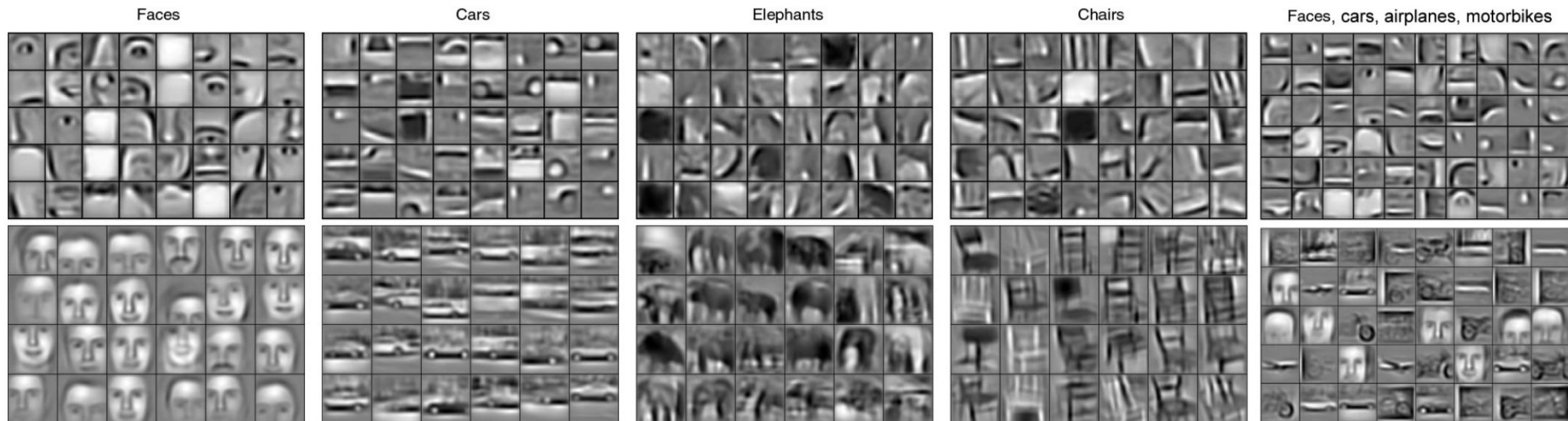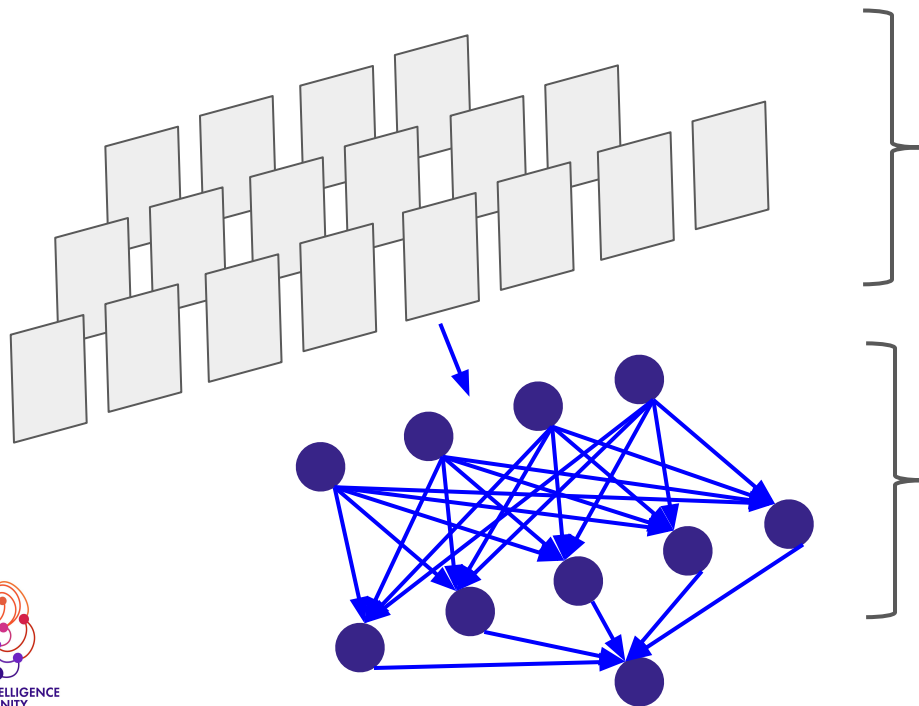


Conceptual diagram

# Primitive features

# Complex features



Image from Lee et al. 2011

# How to reuse learned models?

- Would be useful to be able to reuse learned parameters and learned features



- Trained feature extractor contains features common among domain
- Reuse this part

- Trained fully-connected network contains features specific to dataset
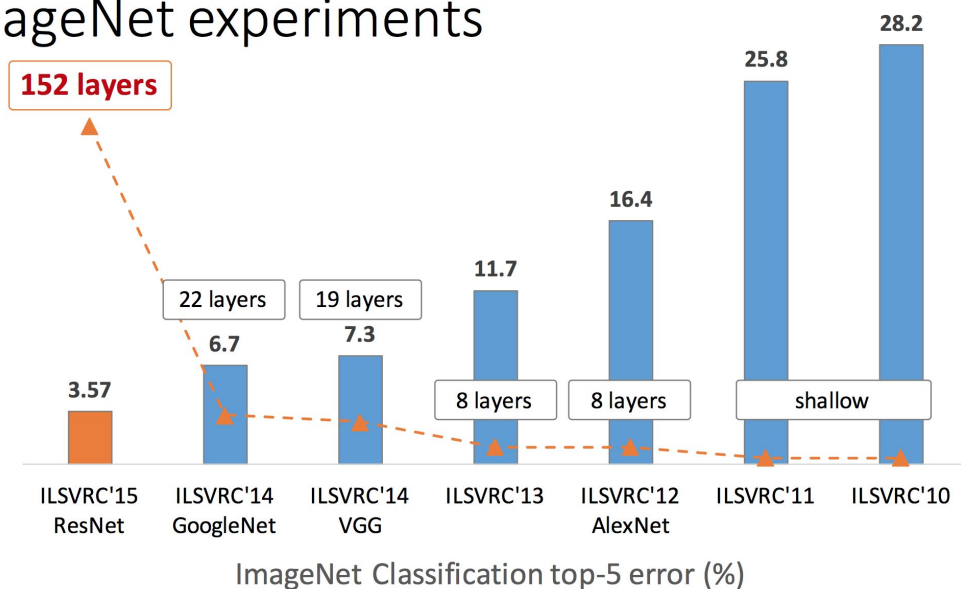- Can potentially reuse weights in latter layers as well

# Deep learning so far ...

- Convolutional neural networks allows us to capture compositional data and convert it into a hierarchical representation.
  - This approach have prove quite effective in classification task, especially in computer vision.
- While CNN is effective in computer vision task, it is also requires a lot of resources.
- CNN, Neural Networks, Backpropagation are not new. All of these techniques have been around since the 1980s. But why it's popular now?
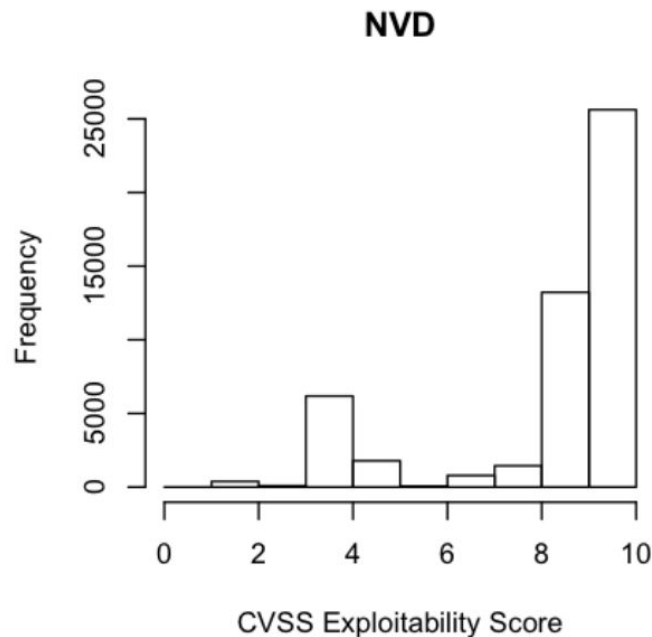
MACHINE INTELLIGENCE
COMMUNITY

# ImageNet

- Both a collection of 14+ millions images that are categorised **by human** into roughly 20 thousands categories. This dataset also fuels the ImageNet challenge.

## ImageNet experiments



ImageNet Classification top-5 error (%)

# The domain of data

- Finding large amount of data for every task we want our machine learning model to perform is hard.
  - e.g. doing machine translation on rare language pair.
  - e.g. data about cyber vulnerabilities exploits



**NVD**

CVSS Exploitability Score

Image from Allodi and Massacci 2013

# The domain of data

- Sometime our input data varies.
  - e.g. speech recognition is challenging because every person's speech is a little bit different. Furthermore, there are different accent of the same language, thus one model for one language would not work well.
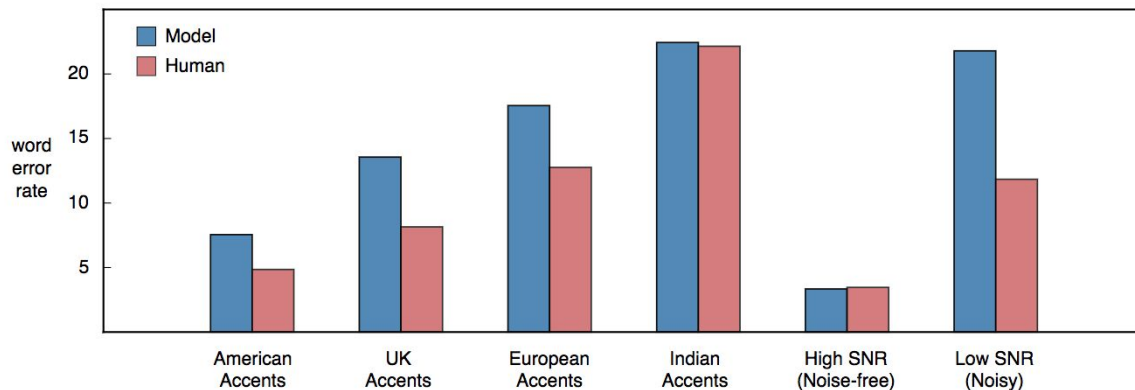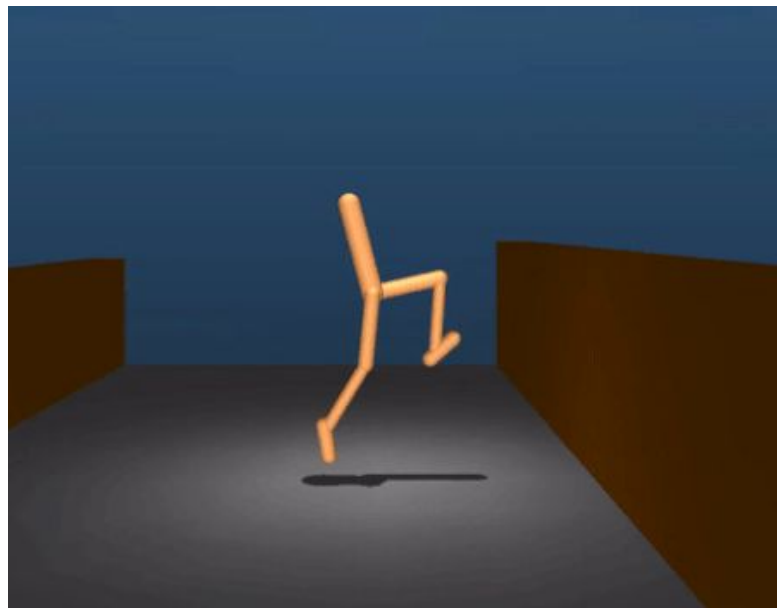


Chart from Awni Hannun

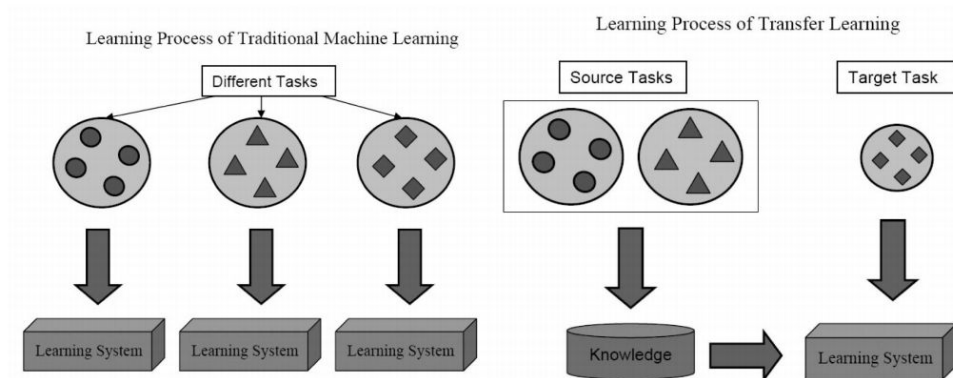# Learning from simulation

- Use a simulation of the real world to train a model.
  - Faster training time
  - Less costly
  - Able to repeat multiple time
- But the data from the simulation is still not the same data from the real world.
  - We need the ability to adapt the model from using the data in the simulated world to the data



*DeepMind's simulated environment*

Image from [DeepMind](DeepMind)

# So, transfer learning

- Utilise the feature extraction property of Convolutional Neural Networks
  - Reduce training time and computational cost.
  - Reduce the amount of data used.
- Improve performance with task where the data domain varies.
  - In general, instead of mapping data to output, we distill some "knowledge" about the data and use that to make a better machine learning model.



Image from Pan and Yang 2010

13

# More formally

$$\mathcal{D} = \{\mathcal{X}, P(X)\}$$

Domain      Feature space      Marginal probability distribution

When $D_s \neq D_t$ or $T_s \neq T_t$, **transfer learning** aims to improve the target predictive function (the conditional probability distribution) in $D_t$ using knowledge from $D_s$ and $T_s$

Task      Output/Label space

$$\mathcal{T} = \{\mathcal{Y}, P(Y|X)\}$$

Conditional probability distribution
(this is what's usually learnt)

14

# Different scenarios of transfer learning

The domains are different

$$\mathcal{X}_S \neq \mathcal{X}_T$$

The feature spaces of the data are different from each other.
E.g. different languages.

$$P(X_S) \neq P(X_T)$$

The marginal probability distributions of the data are different.
E.g. different topic of document.

The tasks are different

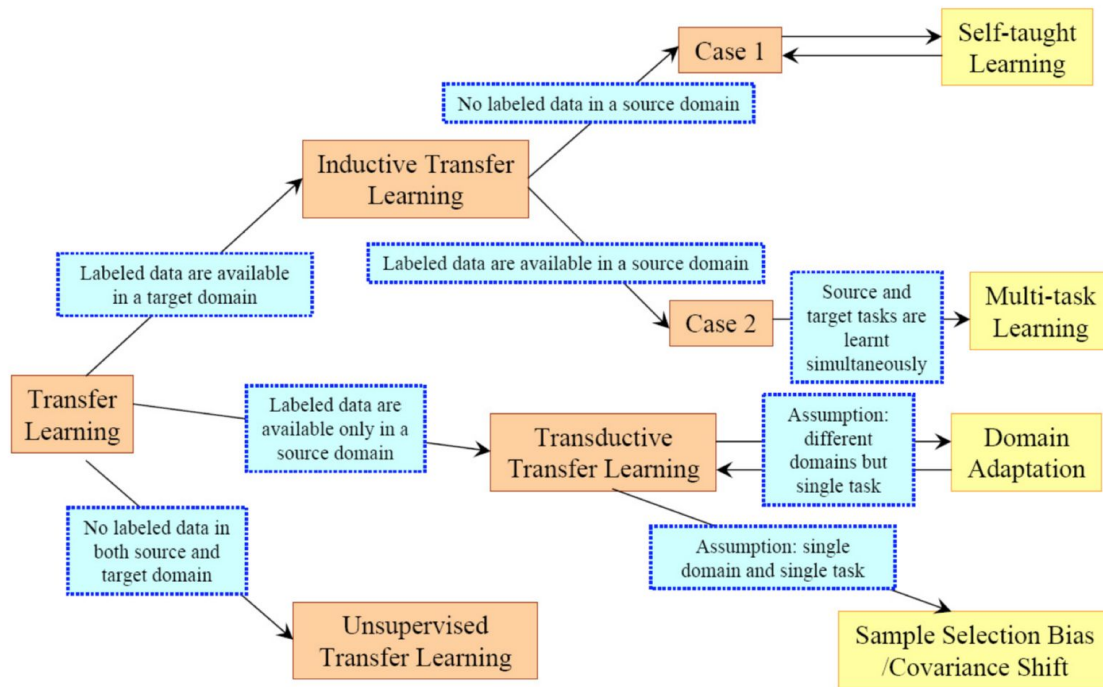$$\mathcal{Y}_S \neq \mathcal{Y}_T$$

The label spaces are different.
E.g. classify the same document into different label depend on different task.

$$P(Y_S|X_S) \neq P(Y_T|X_T)$$

The learned functions are different.
E.g. the source and target distribution of the data is unbalanced compared to their class.
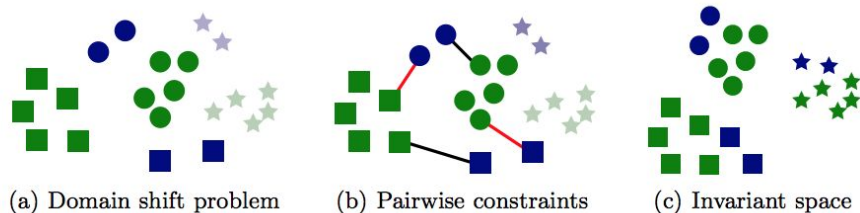
MACHINE INTELLIGENCE
COMMUNITY

# Different setting of transfer learning

- **Both domain and task are different ->** **Unsupervised Transfer Learning**.
- **Domains are different ->** **Transductive Transfer Learning**
- **Tasks are different ->** **Inductive Transfer Learning**
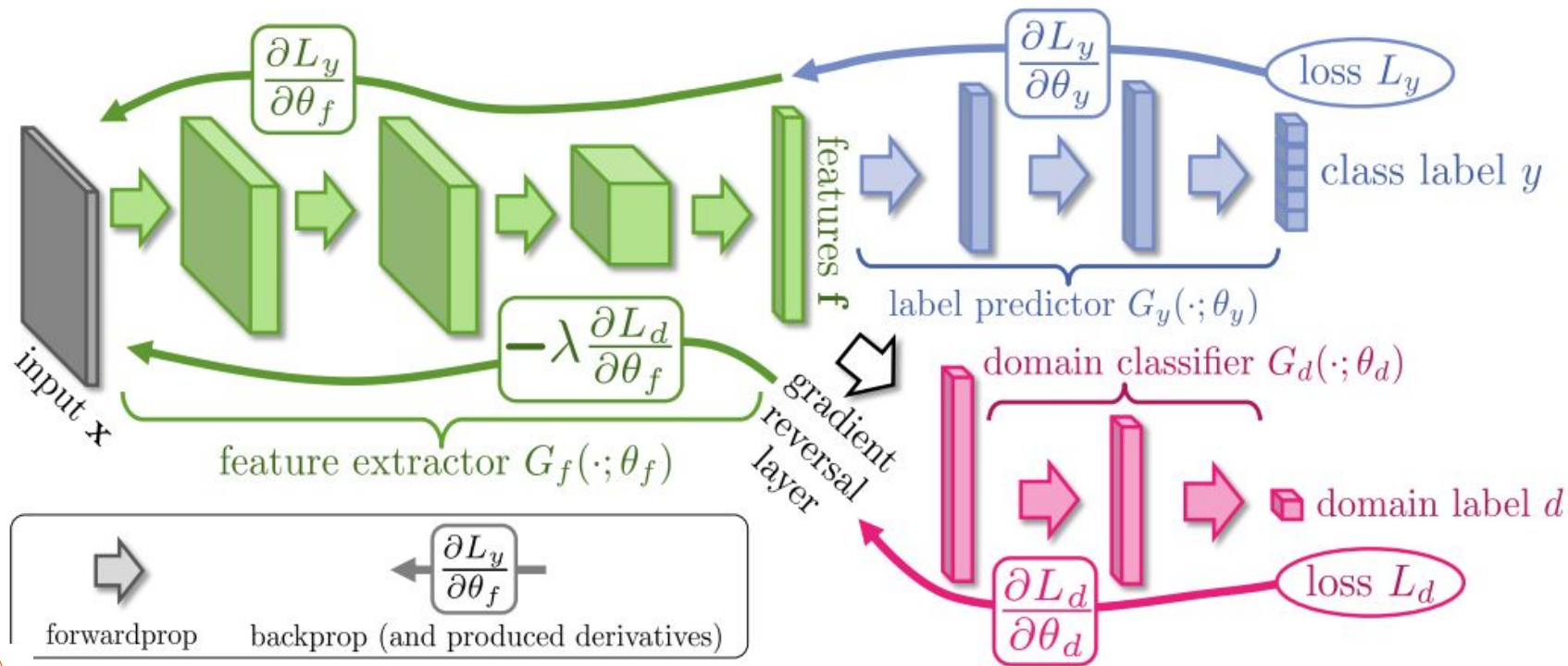
Image from Pan and Yang 2010

# Domain Adaptation

- We are learning the same task but on different domain of data.
- To make our model adapt to the change in domain, we could make the two domain "closer" to each other.
  - Transform data from the domains into an invariant space (K. Saenko et al. 2010)



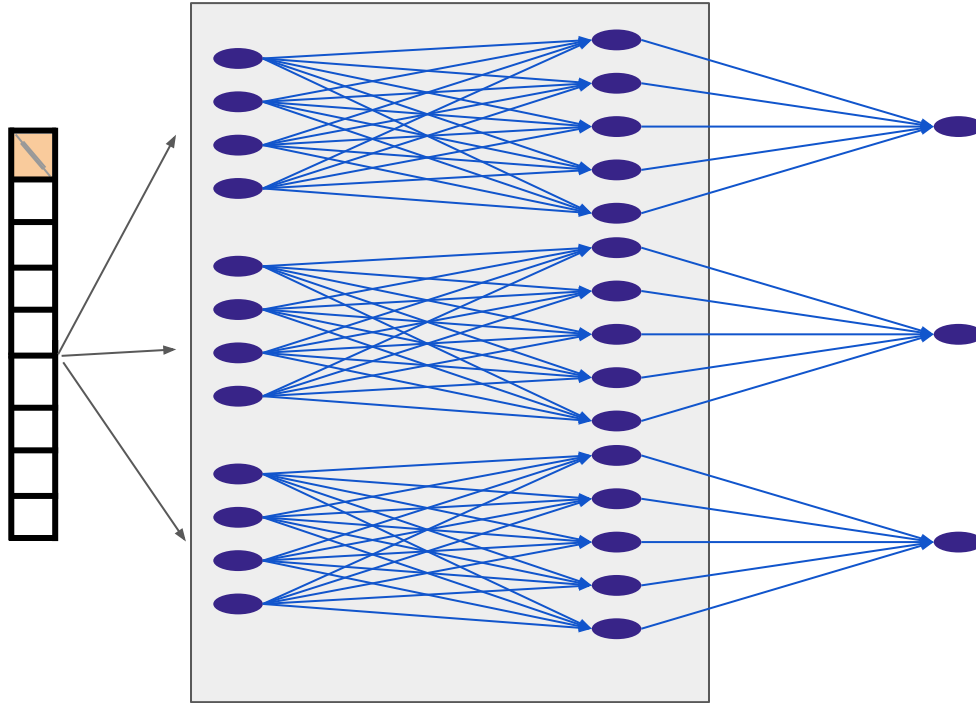(a) Domain shift problem    (b) Pairwise constraints    (c) Invariant space

  - Make it an objective of the model to "confuse" the domains of the data (Ganin 2015)

Image from Saenko et al. 2010

# Domain confusion using backpropagation

Image from Ganin and Lempitsky 2015

# Multi-task learning



- Instead of training a model for each task, let's train one model for all the task and share the parameters for the models.
- This take advantage of the hierarchical representation of the data through the network.

# Few-shot, One-shot, Zero-shot learning

- Data is expensive. Can we use transfer learning to learn with less data?

- **One-shot** learning means that the model have one example of the class before it has to make the prediction and **zero-shot** learning is making prediction on brand new data.
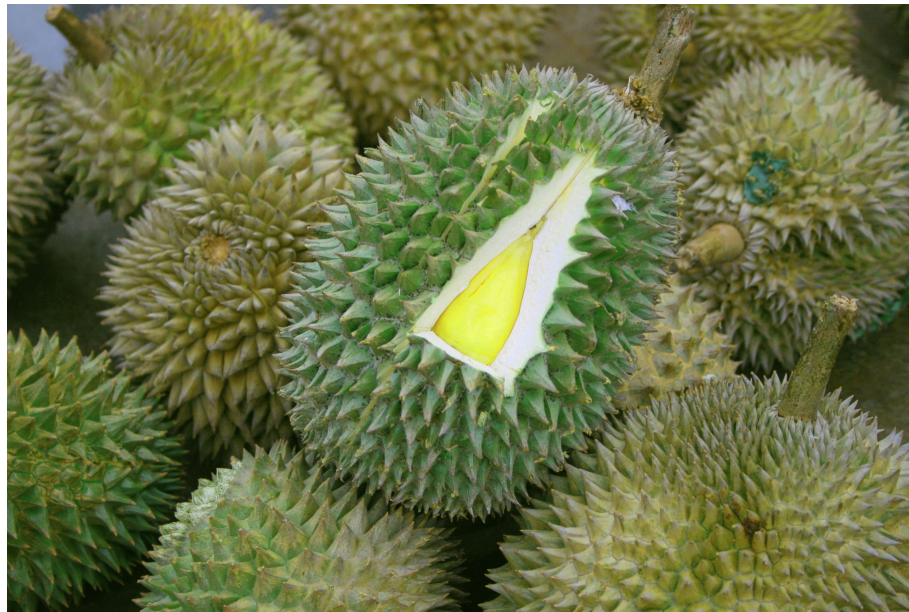
# Example of few shot learning in human



- If you are not familiar with this spiky fruit on the left here, it's a durian.
- Human has the ability to learn/understand concepts, ideas, or objects from very small "training" example.

Image from Wikimedia Commons

# Example of few shot learning in human

- So if you see the image of this spiky fruit again, you would easily be able to recognise that's it's a durian.
  - In essence, you only require one or few training example to "learn" new thing.
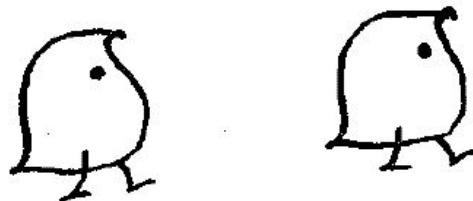
# Example of zero shot learning in human

- Furthermore, human can generalise from other knowledge.
  - The wug test developed by Berko Gleason (1958) show that very young child can apply morphological rule to unknown words.
  - Here, we don't know the plural form of the word "wug", but we are able to formulate it anyway.

THIS IS A WUG.

NOW THERE IS ANOTHER ONE.
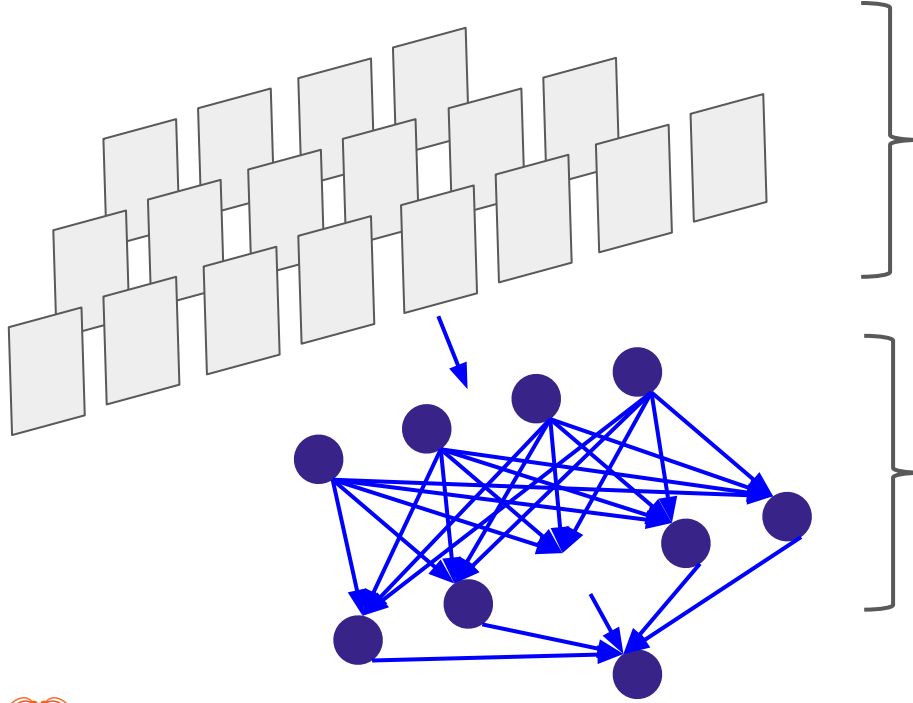
THERE ARE TWO OF THEM.

THERE ARE TWO _____.

Image from Berko 1958

# Word-embedding

- Turns words to vector representation.
- Models like *word2vec* exploit the fact that semantically similar words tend to be around similar words.
- Vector representation have semantics value.



Image from Olah 2014

# Pre-training and fine-tuning



- **Fine-tuning** is when you have a different but related dataset and you want to train the whole model to better fit what you have.
- **Pre-train** models have all the weights already trained. Just make a new Fully Connected layer at the end and retrain the layer.

# Pretrained models with PYTÓRCH

```python
1  import torchvision.models as models
2  resnet18 = models.resnet18(pretrained=True)
3  alexnet = models.alexnet(pretrained=True)
4  squeezenet = models.squeezenet1_0(pretrained=True)
5  vgg16 = models.vgg16(pretrained=True)
6  densenet = models.densenet161(pretrained=True)
7  inception = models.inception_v3(pretrained=True)
```

```python
2  for param in model_conv.parameters():
3      param.requires_grad = False
```

```python
1  num_ftrs = model_conv.fc.in_features
2  model_conv.fc = nn.Linear(num_ftrs, 3)
```

- Import and instantiate pretrained model
- Remove the output layer of model
- Append a randomly initialized classifier and fine-tune

Try it out!

# References & Further Reading

1.  http://machinelearning.wustl.edu/mlpapers/paper_files/ICML2011Glorot_342.pdf
2.  http://www.deeplearningbook.org/contents/representation.html
3.  http://anthropology.uwo.ca/faculty/creider/027/wugs.pdf
4.  https://www.cs.princeton.edu/~rajeshr/papers/cacm2011-researchHighlights-convDBN.pdf
5.  http://ruder.io/transfer-learning
6.  http://ieeexplore.ieee.org/document/5288526/

MACHINE INTELLIGENCE
COMMUNITY

# Shoutout to our Sponsors

Boston University Computer Science

Boston University SPARK!

Boston University Software Application and Innovation Lab

# Upcoming Events

**MIC Paper signup:**          https://goo.gl/iAm6TL
**BUMIC Projects signup:** https://goo.gl/GmP9oK

BUMIC paper discussion:

Paper: **Decoupled Neural Interfaces using Synthetic Gradients**
Location: Fishbowl Conference Room
Date: 11.06.17     Time: 7 PM

Next workshop:

Topic: Sequential Data
Location: BU Hariri Seminar Room
Date: 11.07.17     Time: 7 PM

MACHINE INTELLIGENCE
COMMUNITY