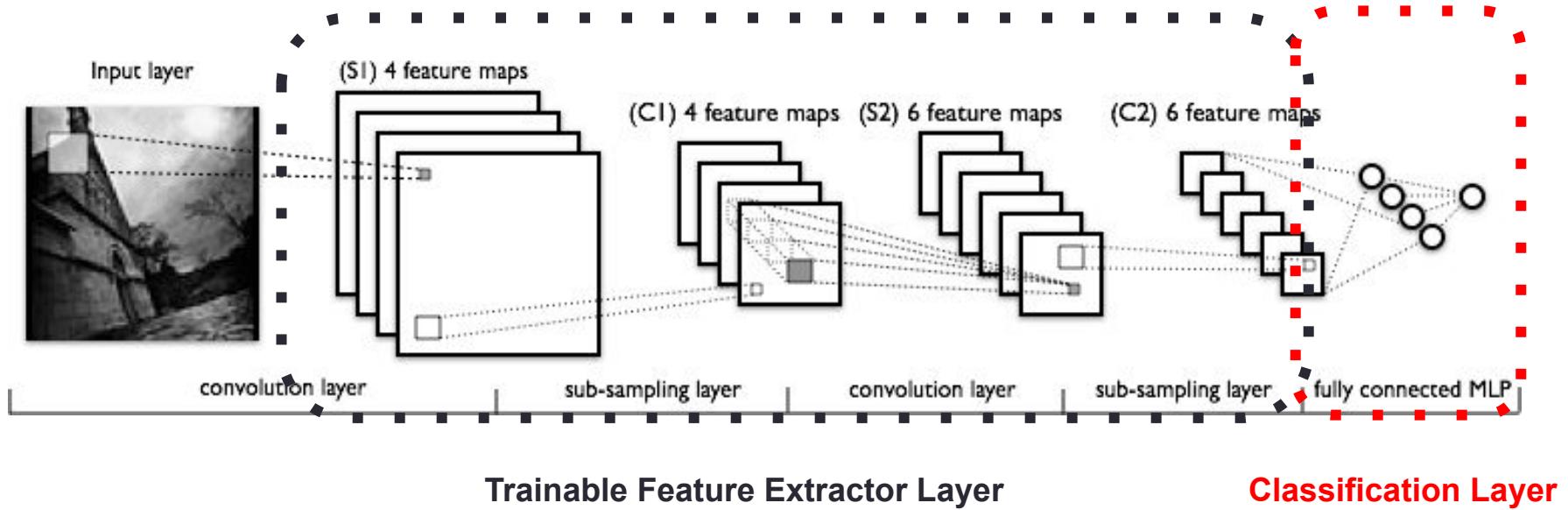


Transfer learning

Deep learning as a feature extractor

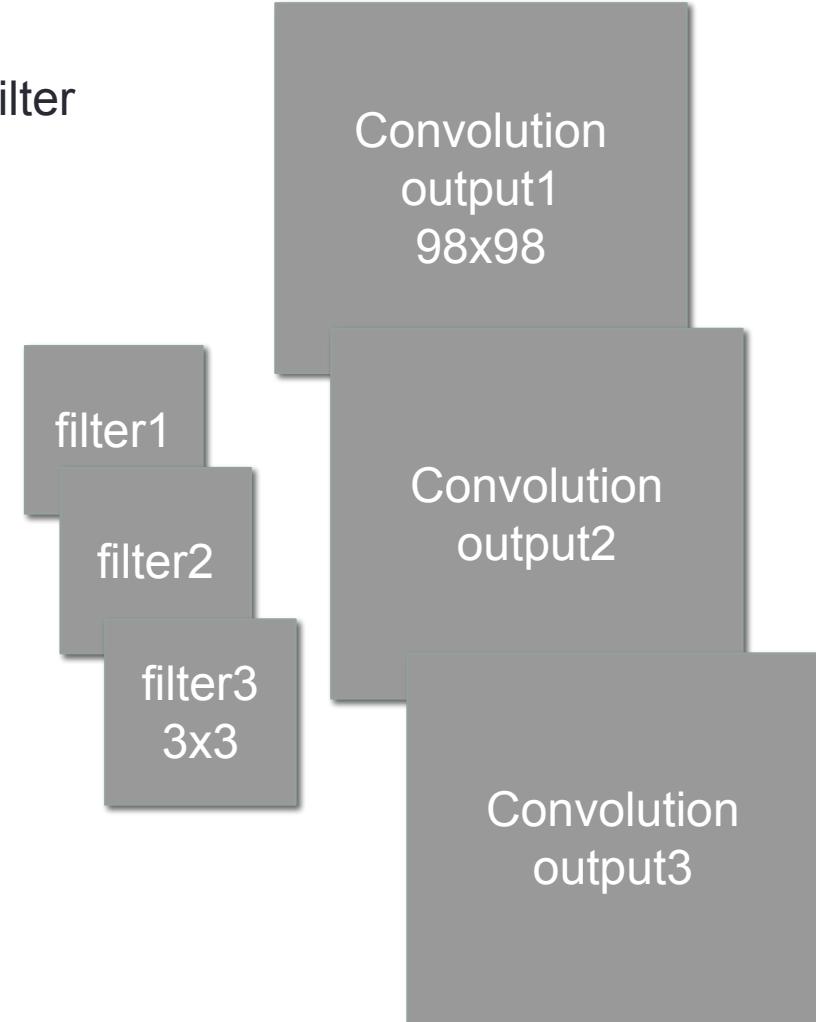
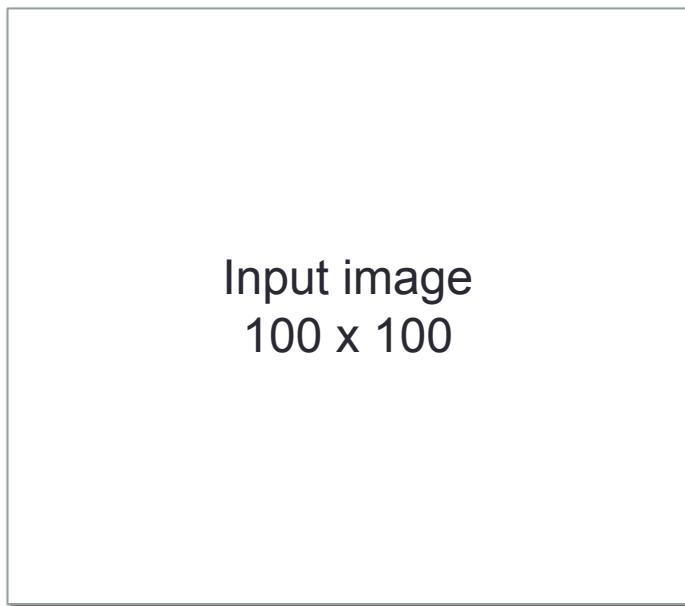
Convolution Neural Network (CNN)



Convolution layers followed by sub-sampling (pooling) layers
Output of each layer is called a feature map.

Convolution filters

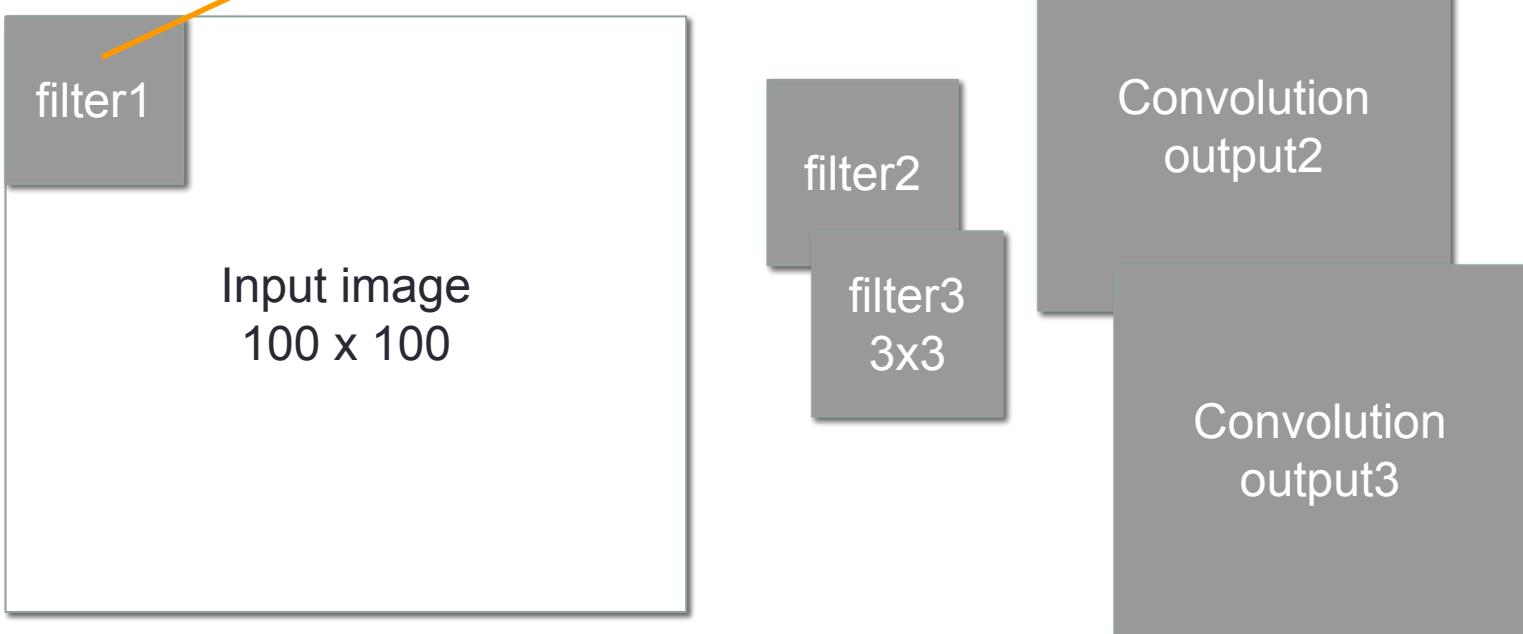
Multiply inputs with filter values
Output one feature map per filter



Convolution filters

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

$$\begin{aligned} & 1*2 + -1*3 + 1*4 \\ & + 1*6 + 1*7 + \\ & 2*8 = 32 \end{aligned}$$

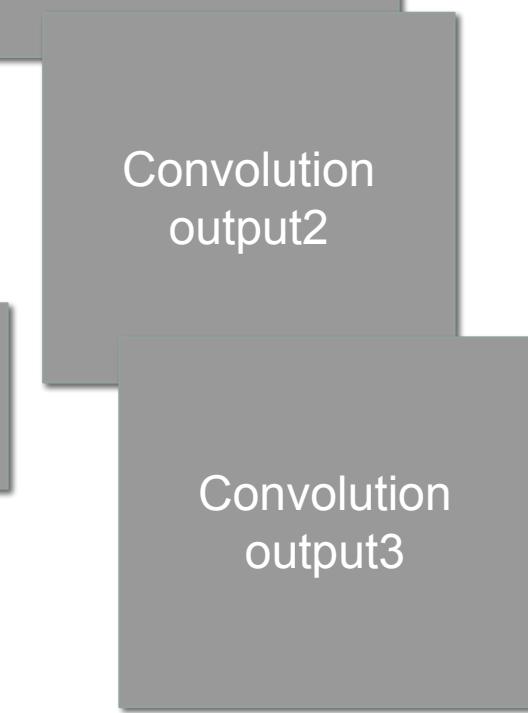
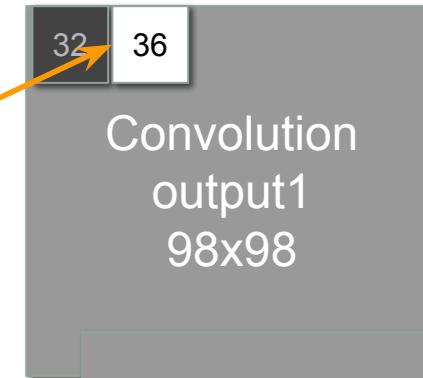
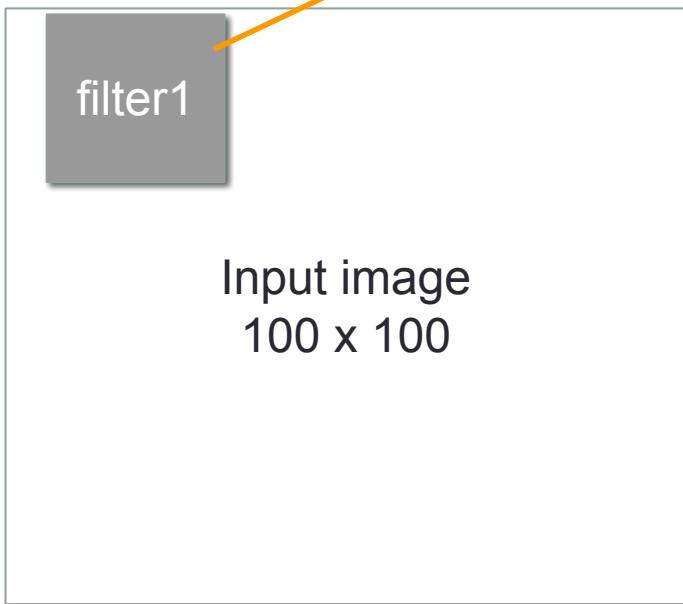


Convolution filters

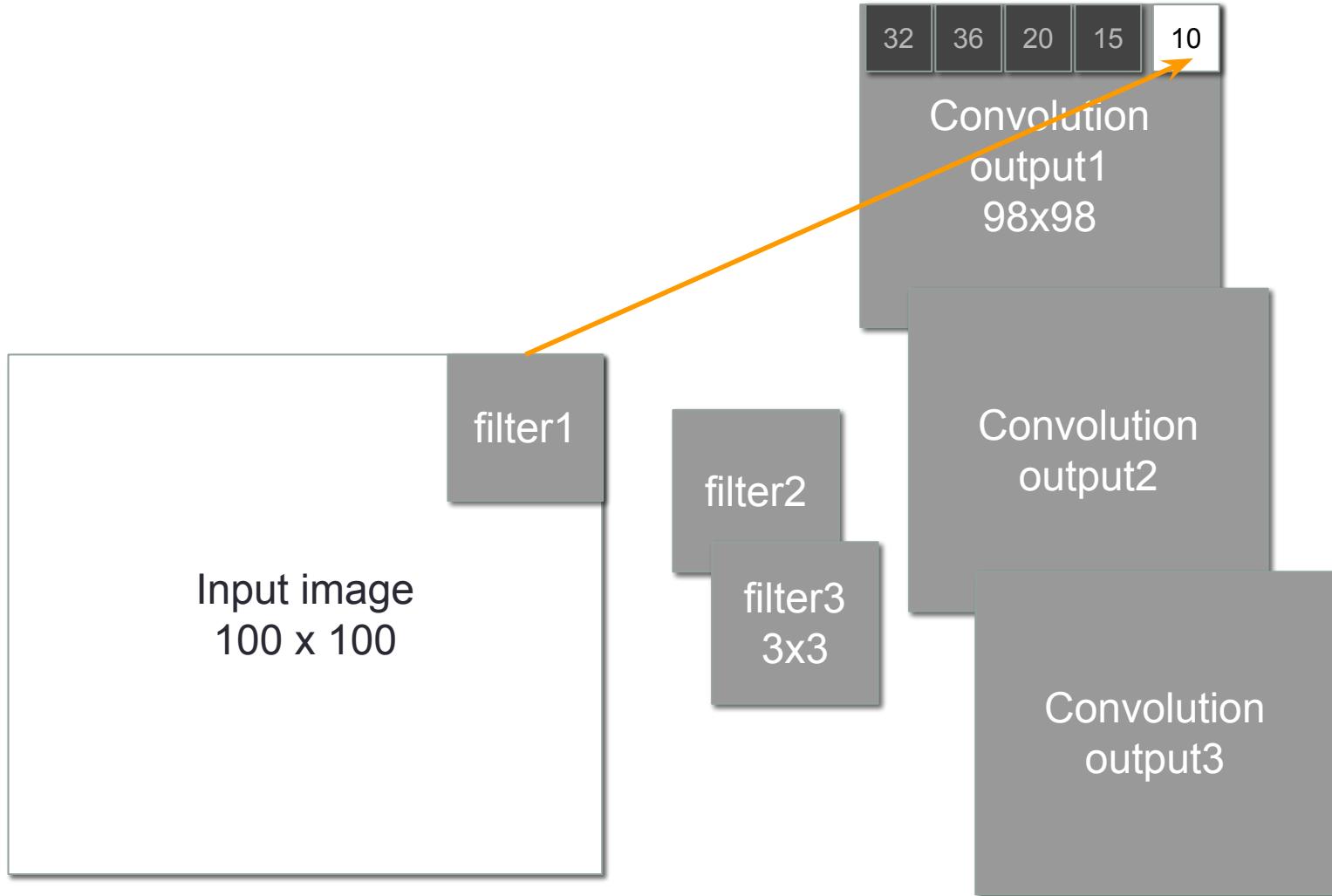
0	1	-1
1	0	1
1	2	0
2	3	1
5	6	3
8	9	8

$$\begin{aligned} & 1*3 + -1*1 + 1*5 \\ & + 1*3 + 1*8 + \\ & 2*9 = 36 \end{aligned}$$

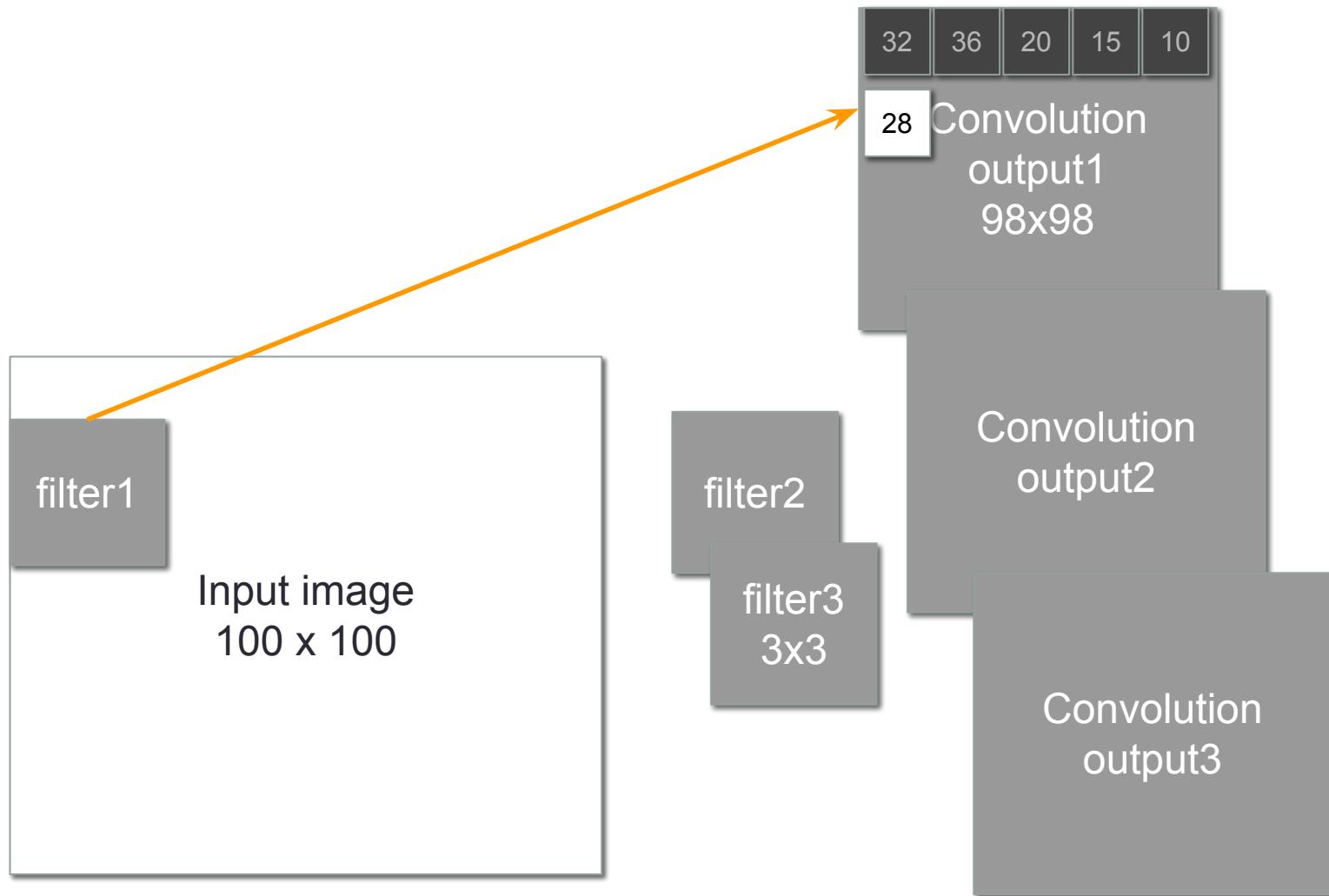
Stride of 1



Convolution filters

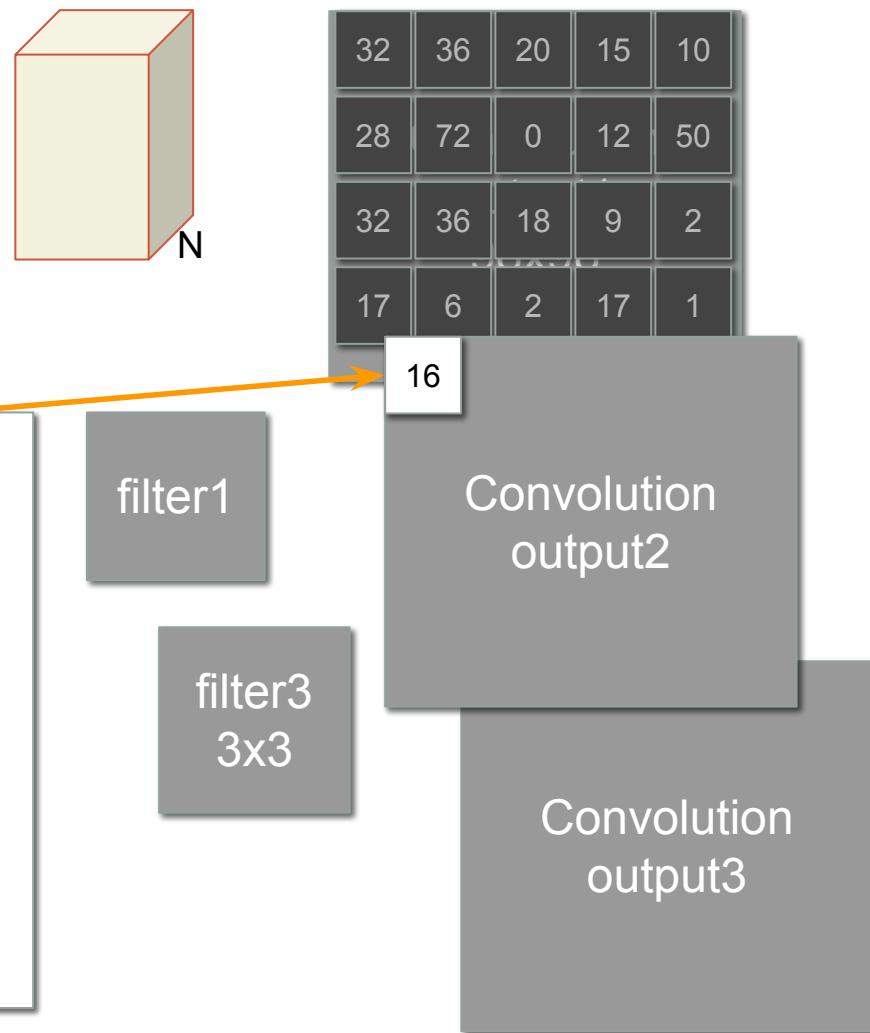


Convolution filters



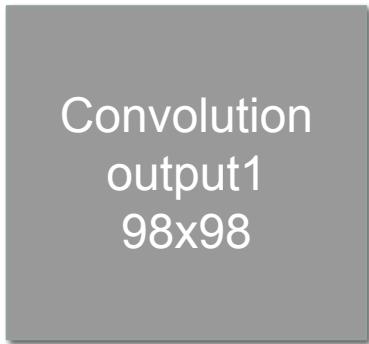
Convolution filters

N filters means N feature maps
You get a 3 dimensional output



Pooling/subsampling

Reduce dimension of the feature maps



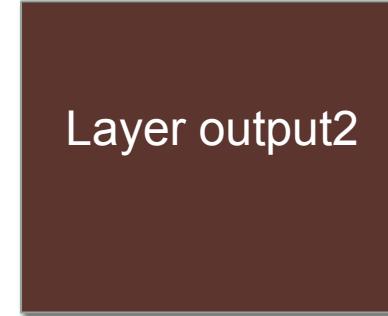
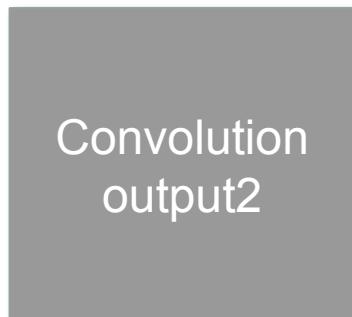
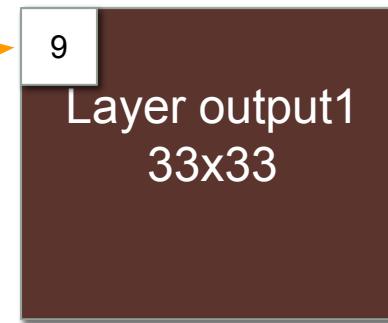
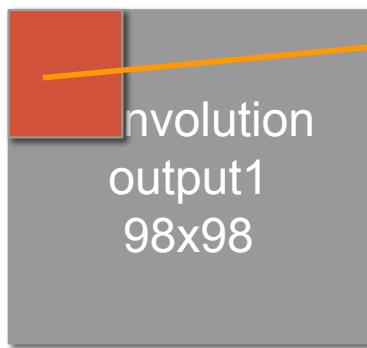
3x3 Max filter
with no overlap



Pooling/subsampling

1	2	3
4	5	6
7	8	9

Max = 9

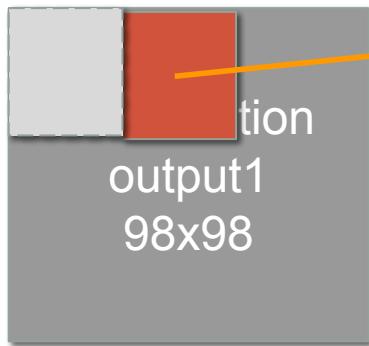


Pooling/subsampling

5	2	1
5	7	1
9	5	12

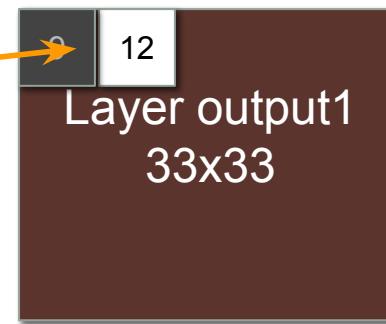
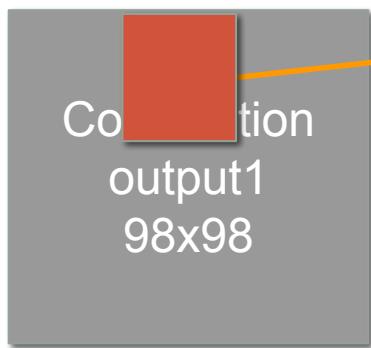
Max = 12

Stride = 3



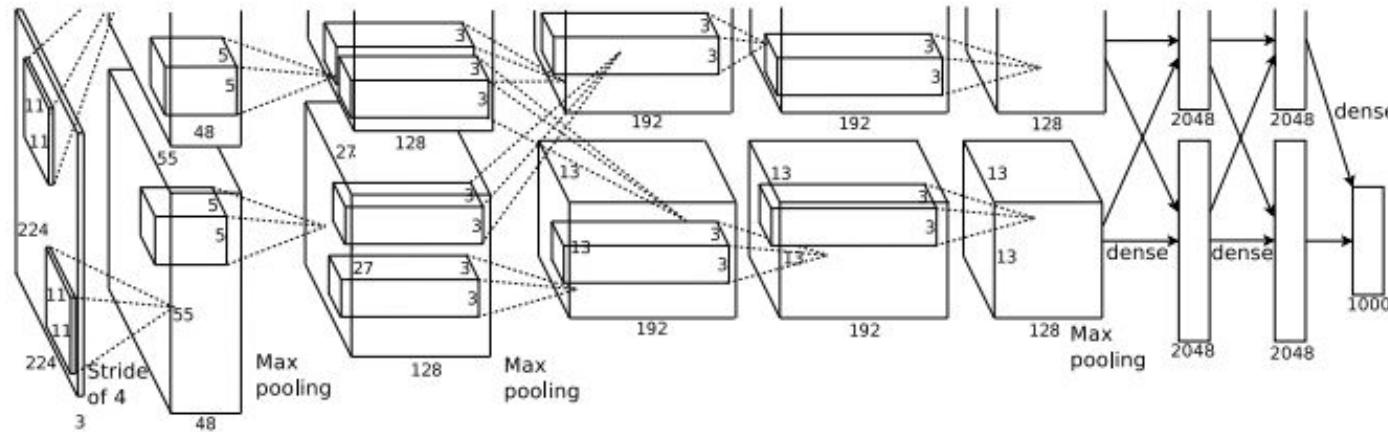
Pooling/subsampling

Can use other functions besides max
Example, average



AlexNet

Convolutions, max pooling, dropout, data augmentation, ReLU activations, SGD with momentum
Two pipelines to fit into two GPUs

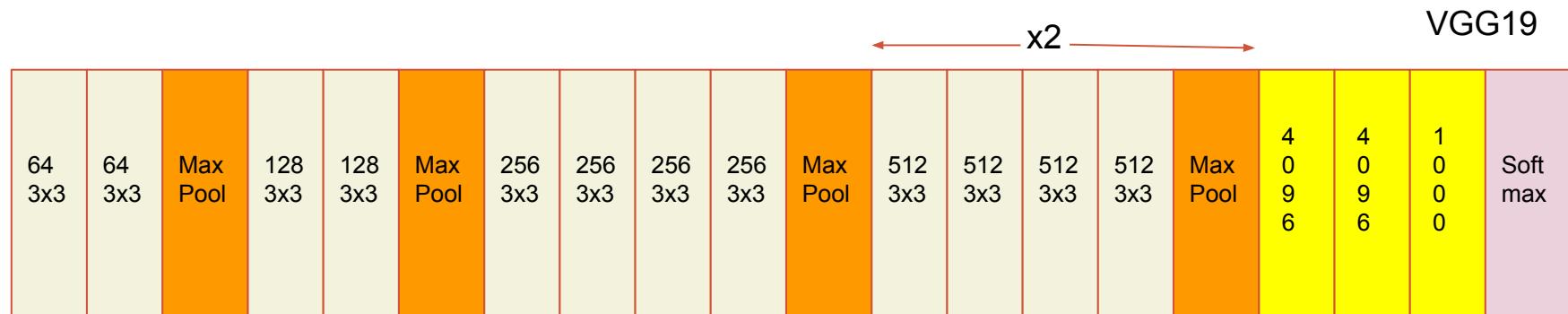


VGG

Uniform 3x3 convolutional filters

19 layers! Pushing the limits of conventional wisdom at that time.

Used by many since pre-train weights are publically available

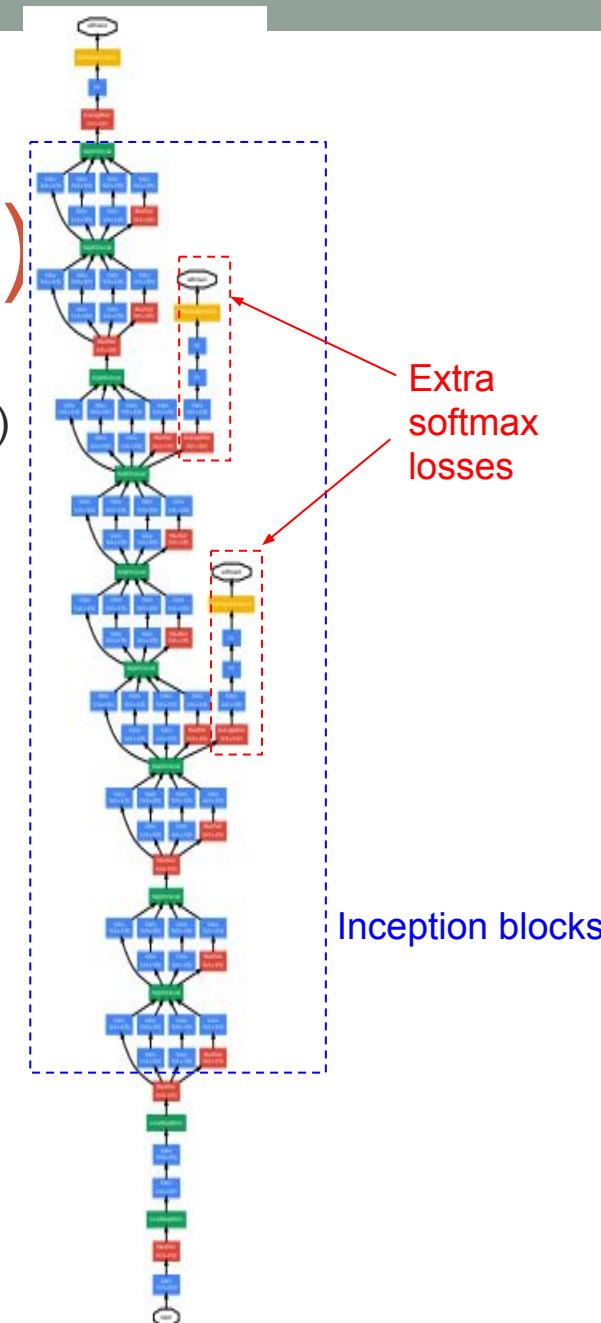
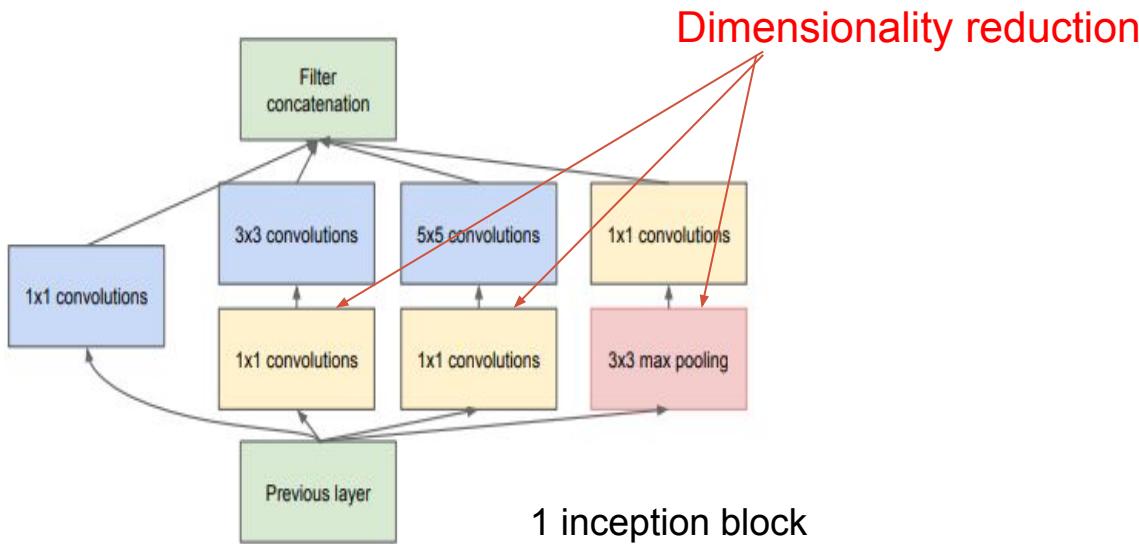


GoogLeNet (Inception v1)

Multiple filter sizes per layer (objects come in different scales)

Dimensionality reduction via 1x1 convolution

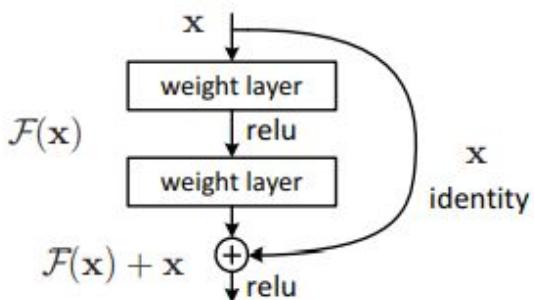
Multiple softmax losses to help the gradient problem



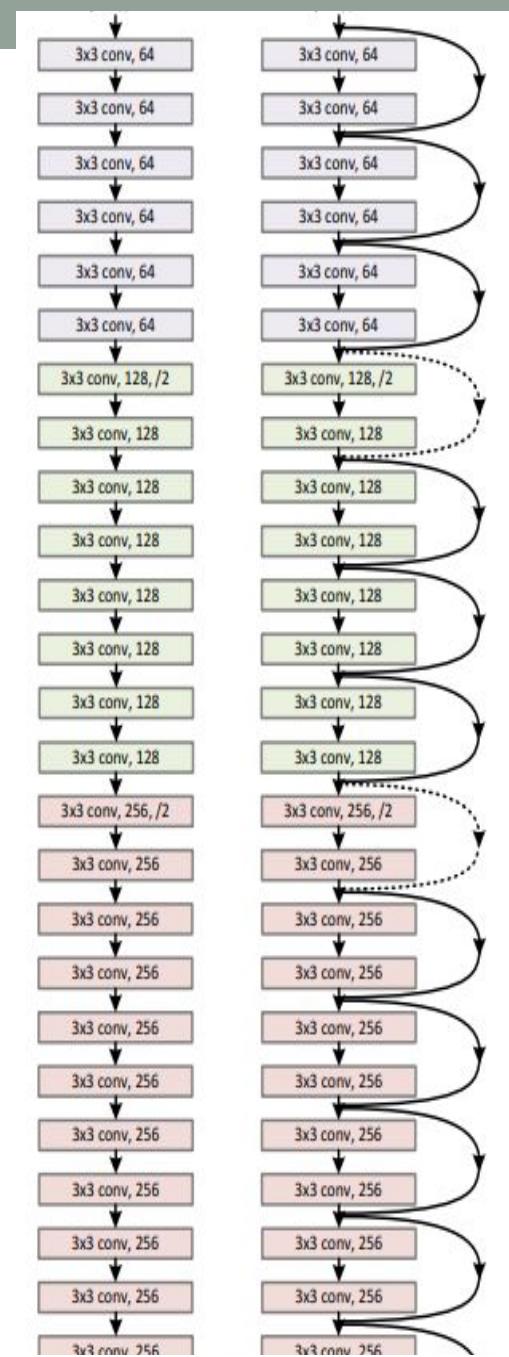
ResNet (Residual Network)

Batch norm

Extra “skip connections” to reduce
the vanishing gradient problem

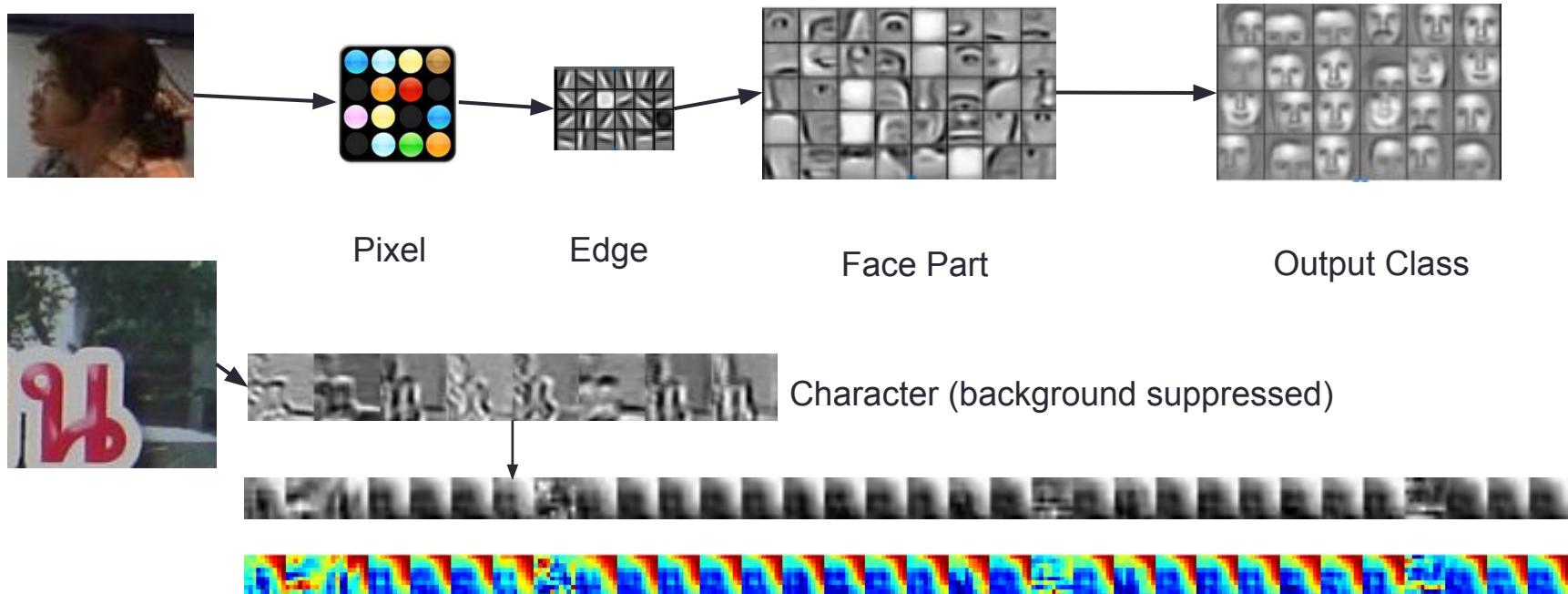


Kaiming He, et al. “Deep Residual Learning for Image Recognition” 2015



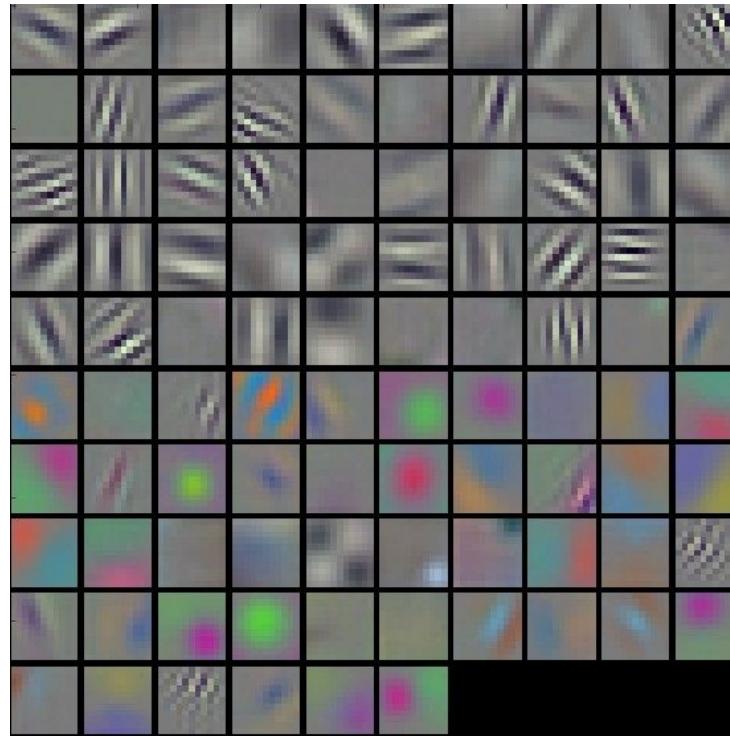
Convolution Neural Network (CNN)

- Hierarchy of representations with increasing level of abstraction
- Each stage is a kind of trainable feature transform
- Image recognition: Pixel → edge → texture → part → object



What does the convolution learn

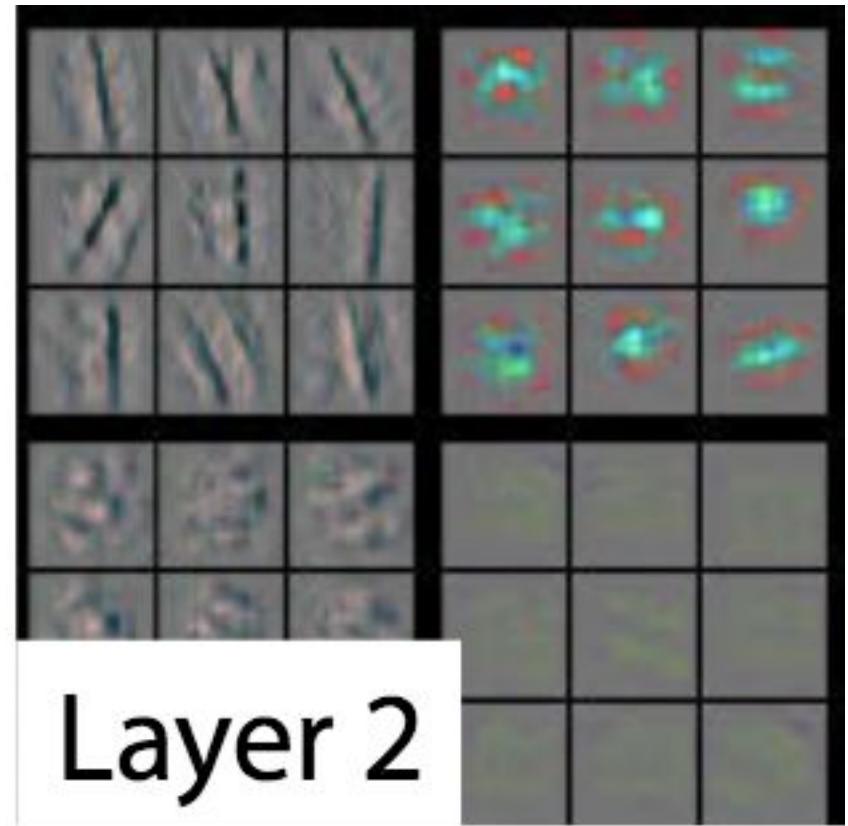
- Learning patterns



Higher layer captures higher-level concepts

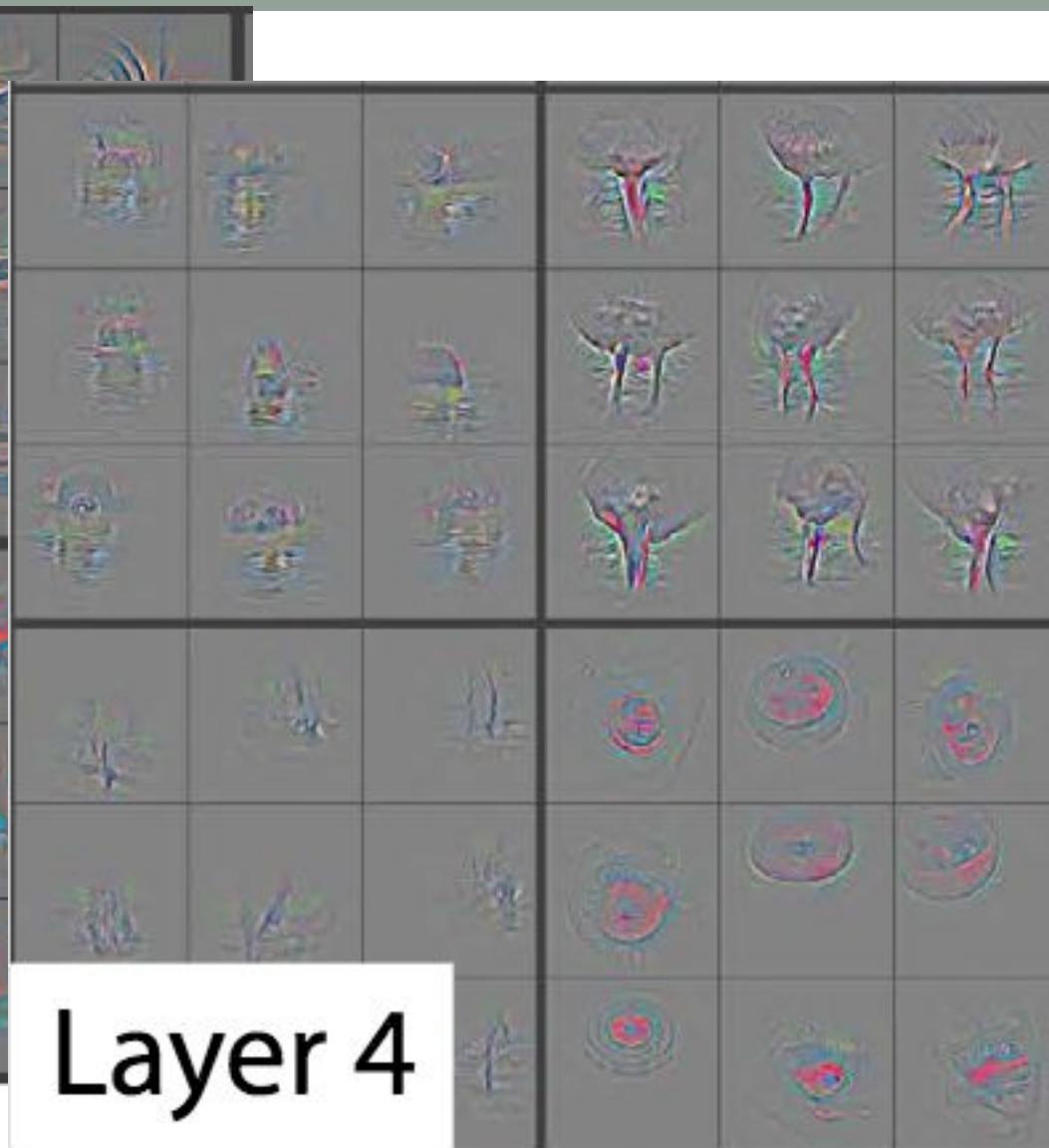


Layer 1





Layer 3



Layer 4

Have you ever seen this creature before?

Can you guess whether it is land or water animal?



You can transfer your knowledge in the past



Transfer learning

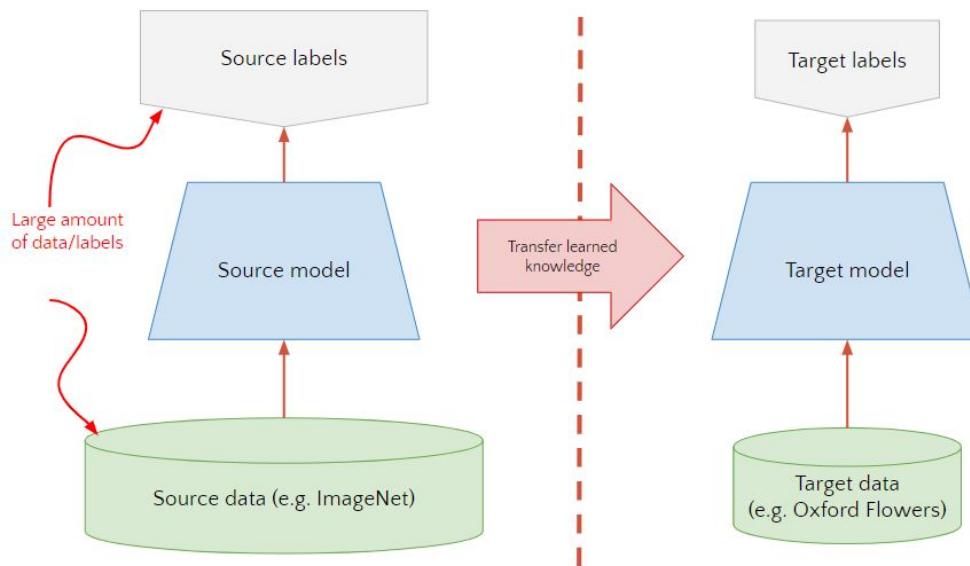
Myth: you can't do deep learning unless you have a million labelled examples for your problem.

Reality:

- You can **transfer** learned representations from a **related task**
- You can train on a nearby **surrogate objective** for which it is easy to generate labels

Transfer learning (basics)

- We know networks captures good representations
- Can we use it for other tasks?
- Use trained networks to initialize a new network for a different task.
- Re-train the network using SGD on new data.



For CV tasks, we call the pre-trained network **backbones**

Transfer learning idea

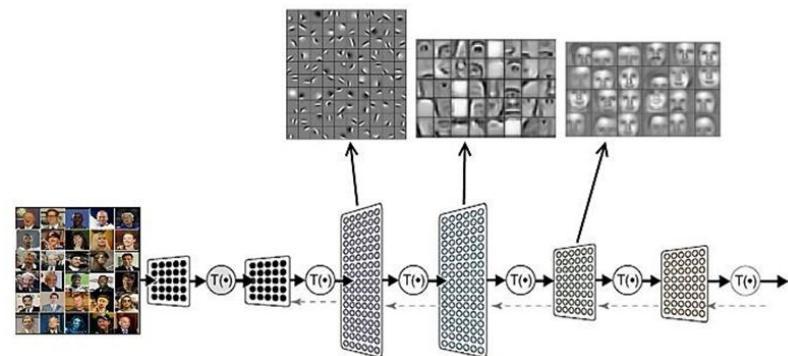
Instead of training a deep network from scratch for your task, you can

- Take a network trained on **a different domain** for a **different source task**
- **Adapt (fine-tune)** it for your domain and your **target task**

This lecture will talk about how to do this.

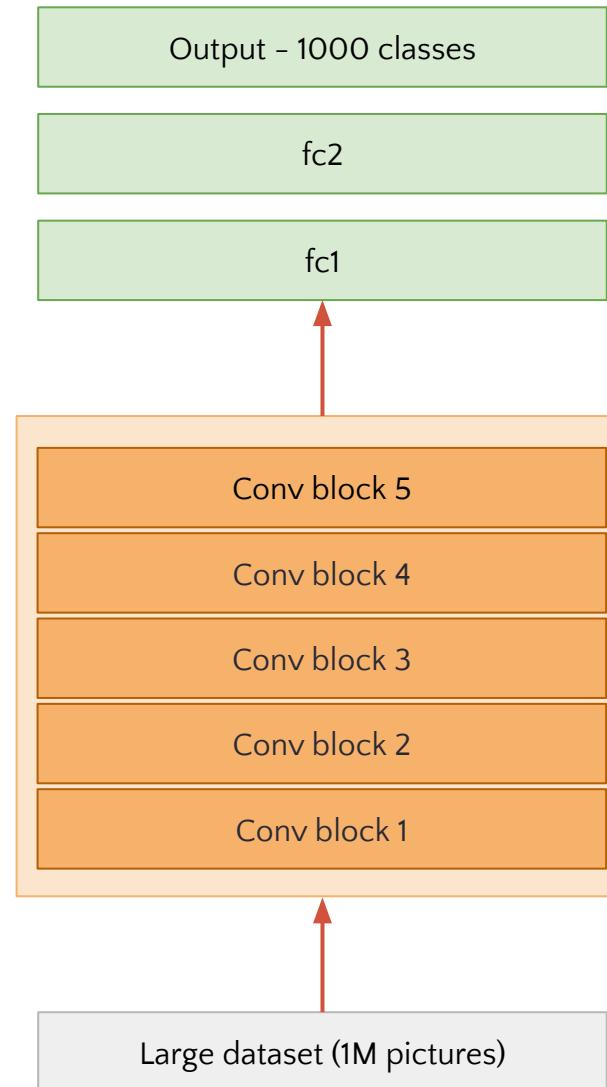
Variations:

- Different domain, same task
- Different domain, different task

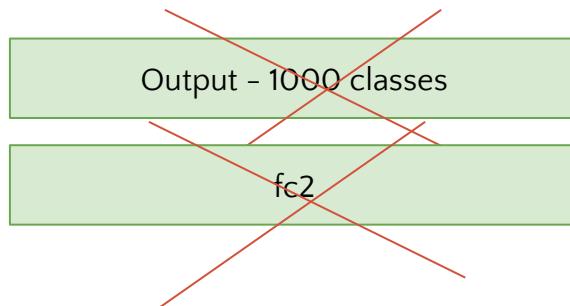


Transfer learning

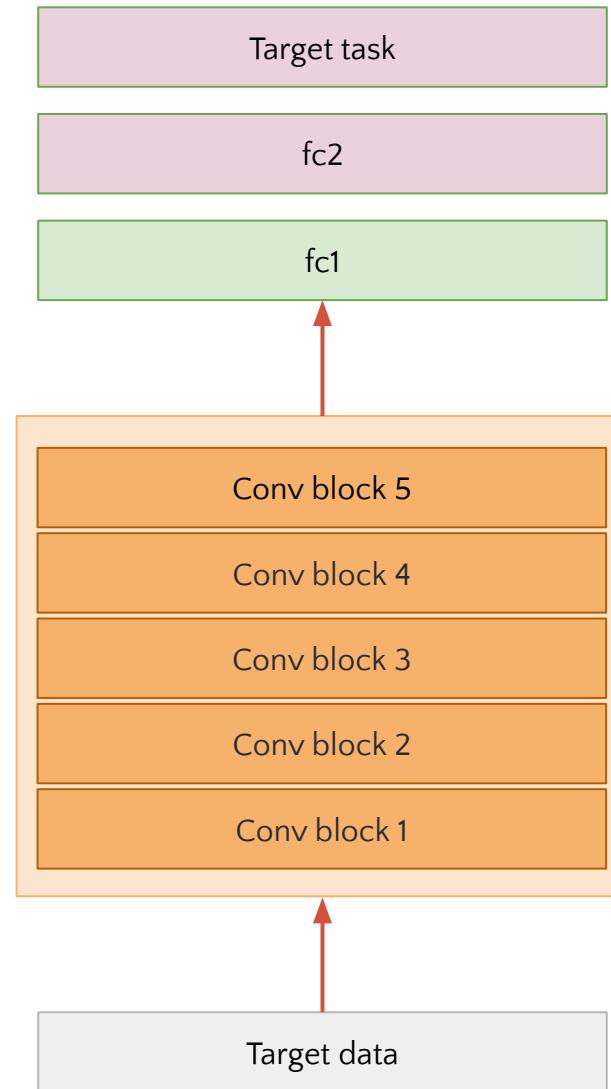
1. A model is trained on large dataset
Ex ImageNet



Transfer learning



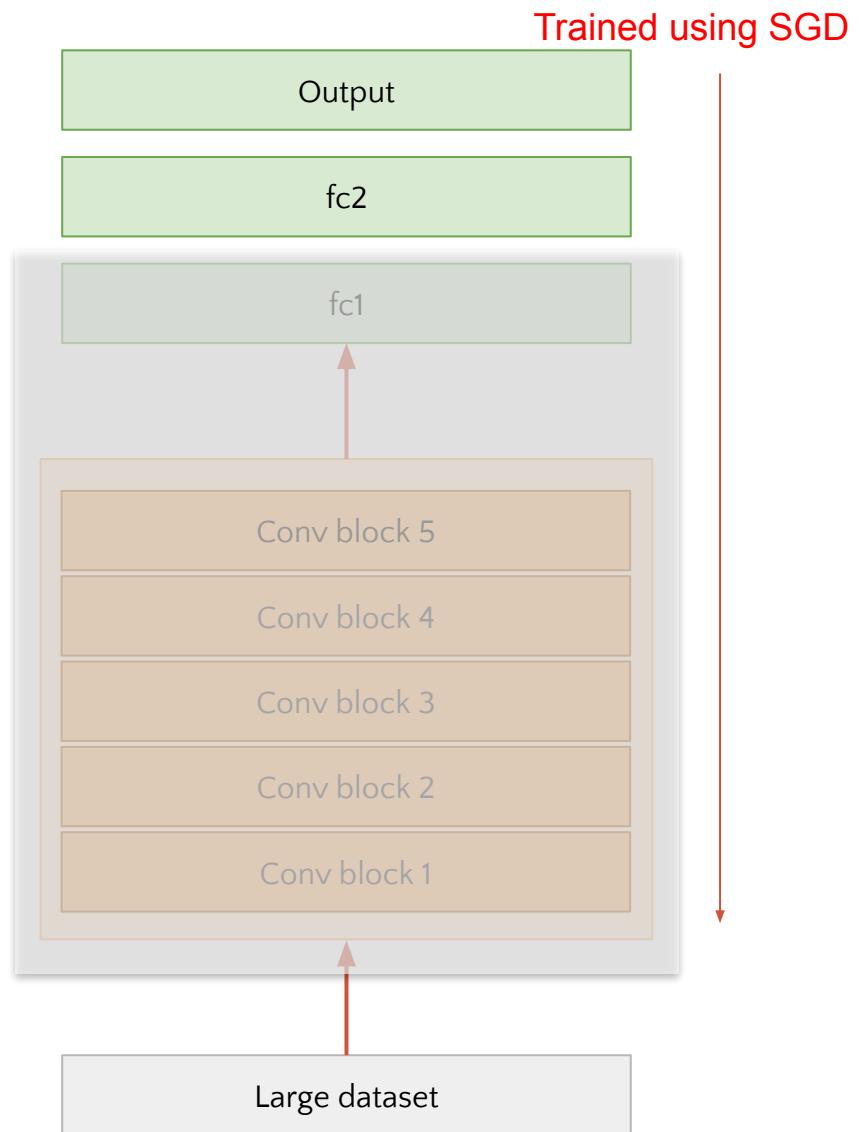
2. Replace the top layers with target task



Transfer learning

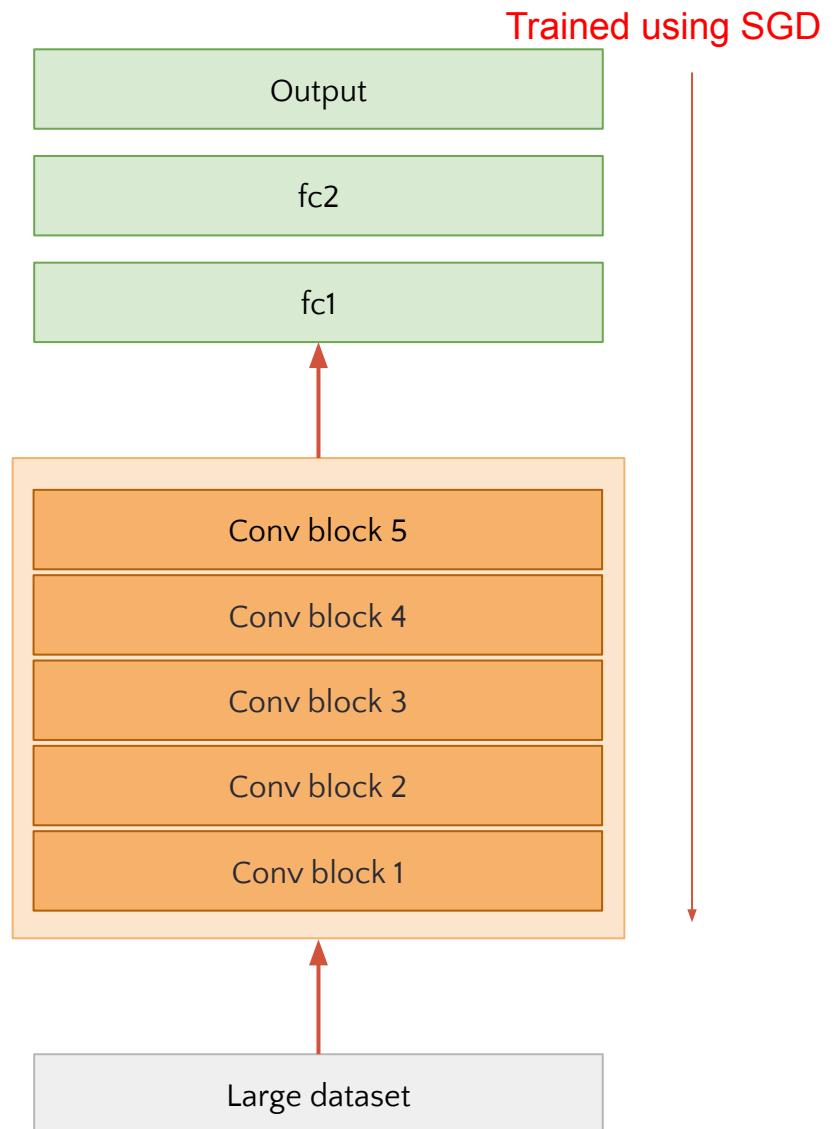
Freeze weights

3A. Train only the layer that is replaced (freezed weights)



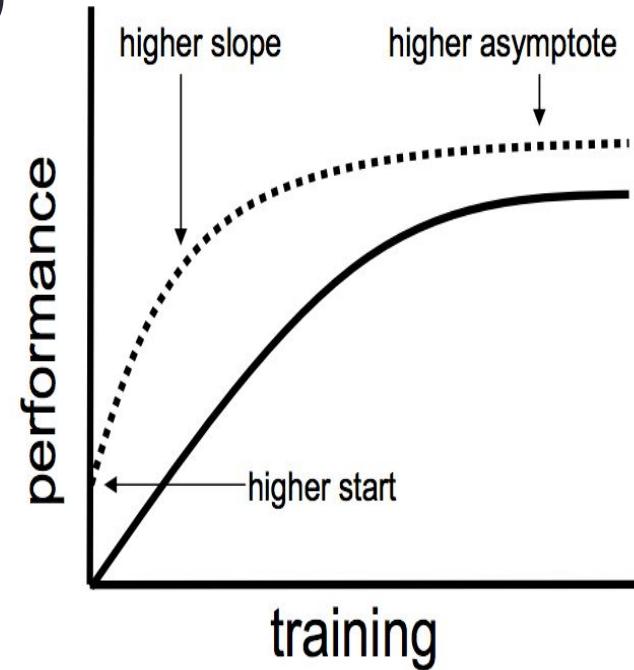
Transfer learning

3B. Train all layers
(unfreezed weights)



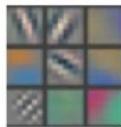
Benefits to transfer learning

1. Higher start
2. Higher slope (converge faster)
3. Higher asymptote (if small data)

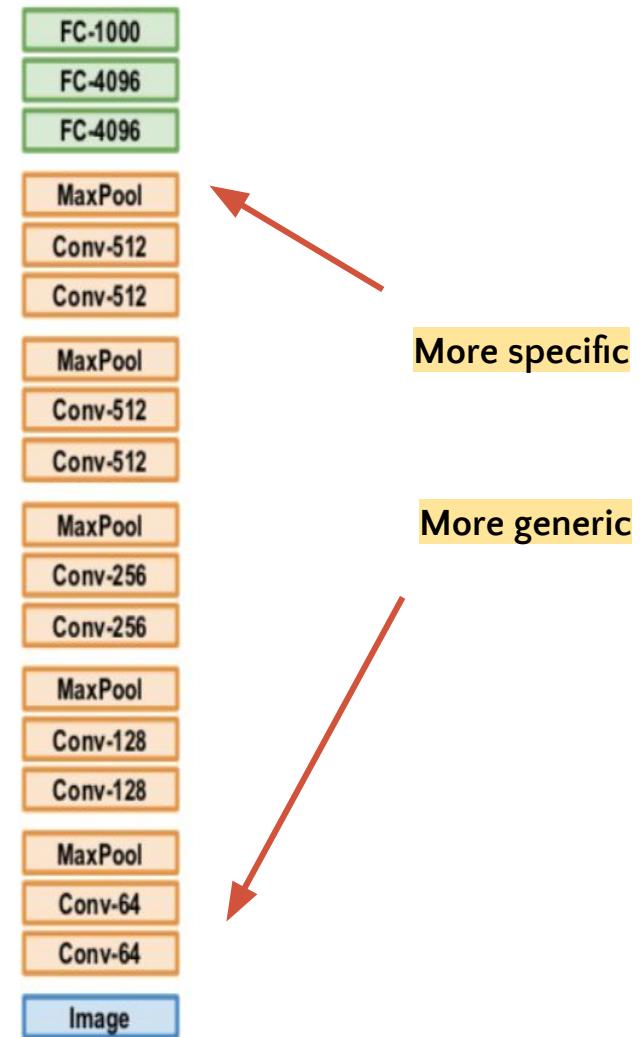
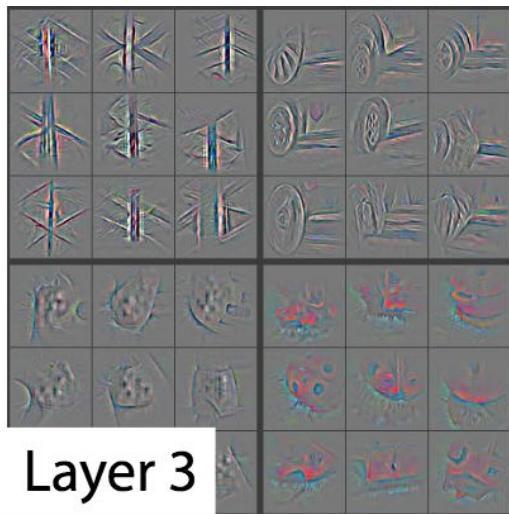
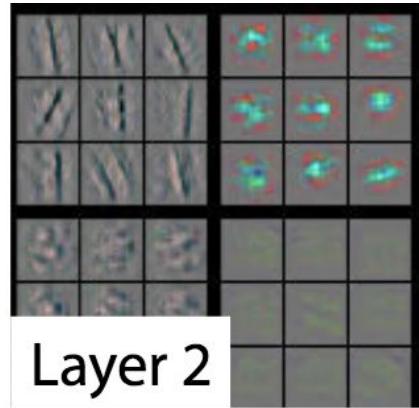


Layers

Lower layers: more general
Higher layers: more specific



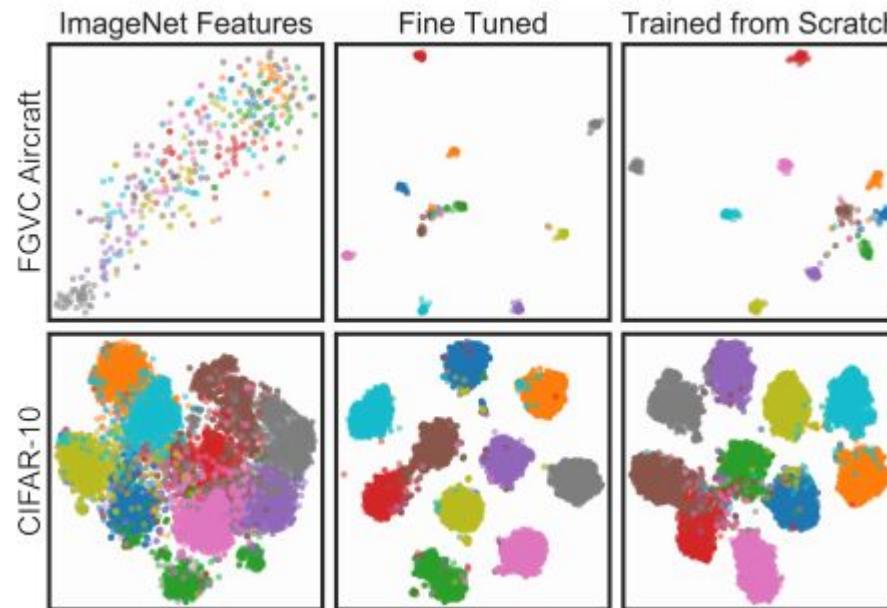
Layer 1



Domains

Similar domain can transfer easier
Fine-tuning (adaptation) is crucial

Aircraft images
Different from ImageNet



Natural images
Similar to ImageNet

Transfer learning in speech recognition

Assamese and Bengali language



Bengali



Assamese

Bengali	WER
No pre-training (10 hr Bengali)	71.8%
No pre-training (60 hr Bengali)	64.5%
Transfer learning (10 hr Assamese -> 10 hr Bengali)	66.0%
Transfer learning (60 hr Assamese -> 10 hr Bengali)	64.6%
+ Adaptation	63.7%
Transfer learning (60 hr Bengali -> 10 hr Bengali)	61.6%

Closer domain is better

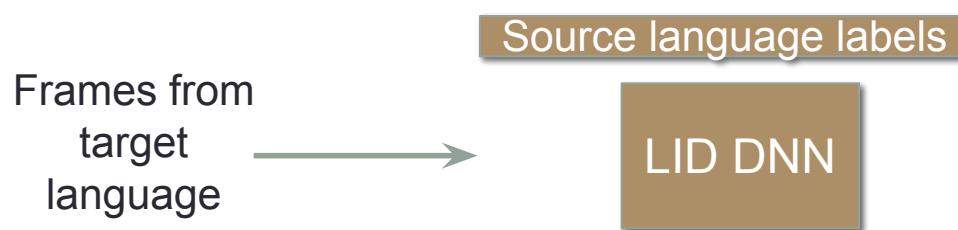
- Pre-train on 60 hours of data, and adapt on languages with 10 hours of data
- Numbers in **blue** is for 10->10 hours (baseline)

		Pretrain on			
		Bengali	Assamese	Lao	Turkish
Adapt to	Bengali	66.0	63.8	65.1	64.2
	Assamese	61.2	65.2	62.9	62.1
	Lao	59.8	60.1	62.3	60.0
	Turkish	61.8	63.1	63.3	63.9

- Transfer learning always improves performance.
- Similar languages perform better on transfer learning

Language Identification for data selection

- Train a classifier (LID) on source languages to predict the language given input frames
- Compute posteriors of the target language data using that classifier
- The best language for the target language should have the highest average posterior score



Predicting the best language

		Pretrain on			
		Bengali	Assamese	Lao	Turkish
Adapt to	Bengali	66.0	63.8	65.1	64.2
	Assamese	61.2	65.2	62.9	62.1
	Lao	59.8	60.1	62.3	60.0
	Turkish	61.8	63.1	63.3	63.9
		Language prediction score			
		Bengali	Assamese	Lao	Turkish
Input frames	Bengali	0.57	0.21	0.09	0.13
	Assamese	0.21	0.57	0.11	0.11
	Lao	0.08	0.11	0.71	0.10
	Turkish	0.13	0.12	0.10	0.65

The LID scores correspond to the best language to use most of the time.

Which task to transfer?

TASKONOMY

Home API Live Demo Pretrained Task Bank Paper Dataset Team

Sample Images (click to use).



Image to upload:

Choose File No file chosen

Run inference

Input Image



Surface Normals

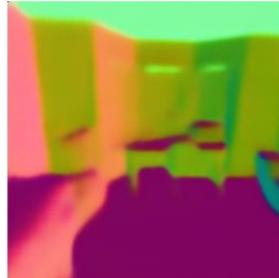
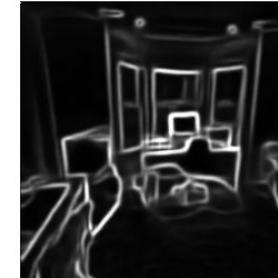


Image Reshading



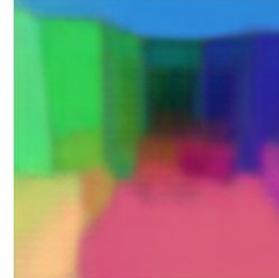
2D Texture Edges



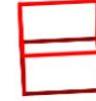
Vanishing Points



Unsupervised 2.5D Segm.



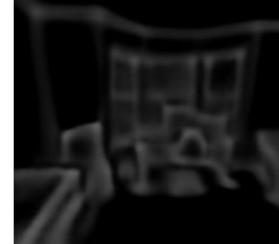
Room Layout



Scene Classification

Top 5 prediction:
home_office
office
television_room
computer_room
office_cubicles

3D Keypoints



3D Occlusion Edges



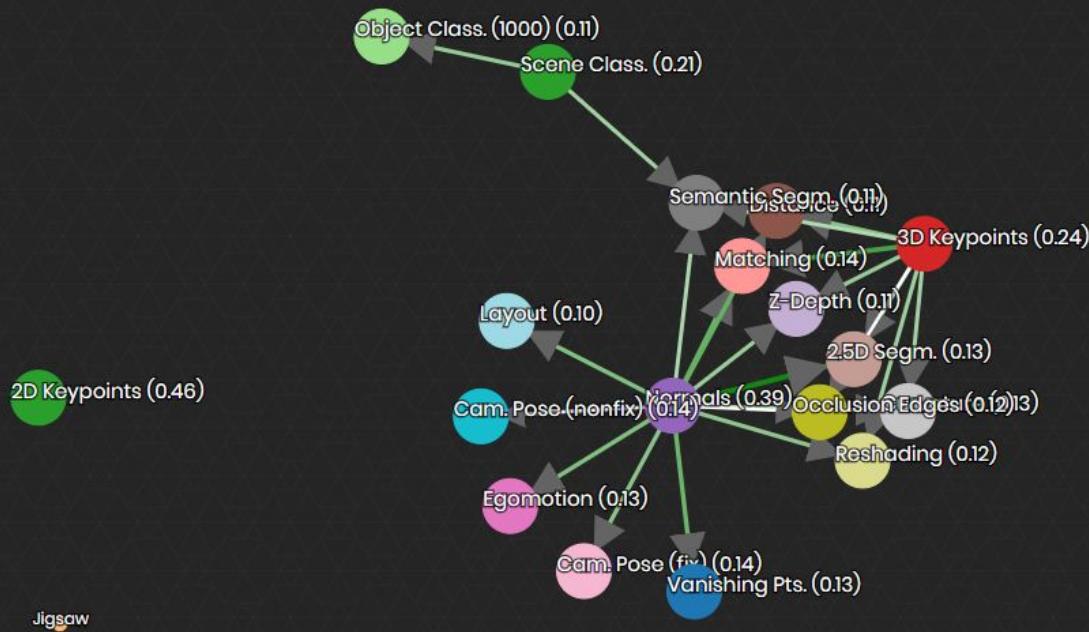
Choose task dictionary ▾



Solve! ↗

Instructions/Definitions ⓘ

Color Nodes:



0% 100%

Adaptation tricks

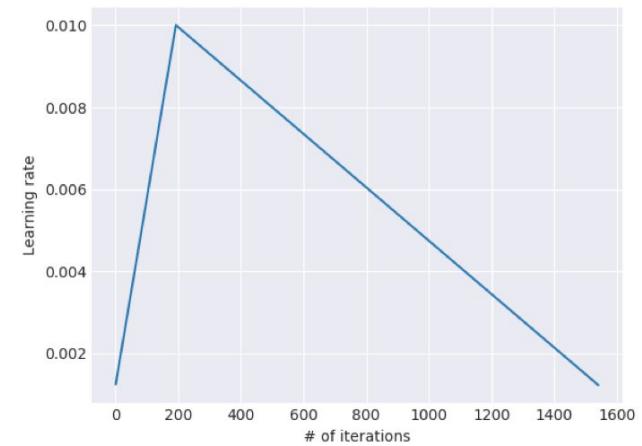
Triangle learning rate

At the start some of the model is randomly initialized. Cannot trust the gradient

Discriminative fine-tuning

Instead of using the same learning rate for all layers of the model, discriminative fine-tuning allows us to tune each layer with different learning rates.

- Layers that are closer to inputs → large learning rate
- Layers that are closer to outputs → small learning rate



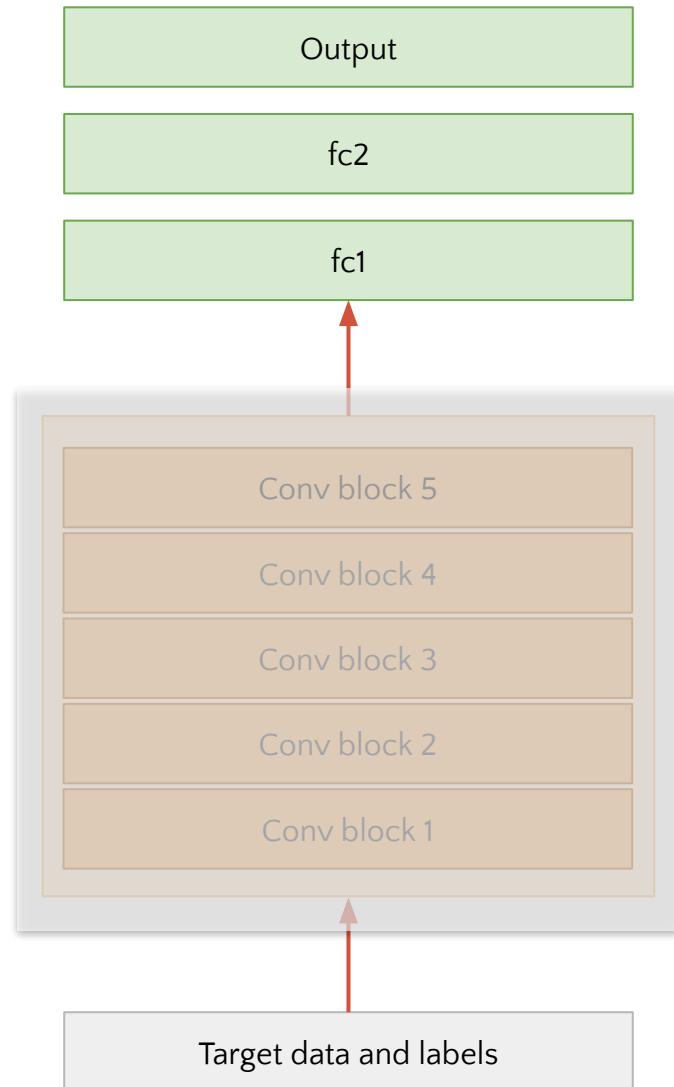
Advance adaptation ideas

- Chain-taw
 - Freeze the old layers, train the new weights for a bit
1. fine-tunes any new layers.
 2. fine-tunes each layer individually of base model starting from the first to the last.
 3. the entire model is trained with all layers.

Chain-taw

1. fine-tunes any new layers.

Freeze weights

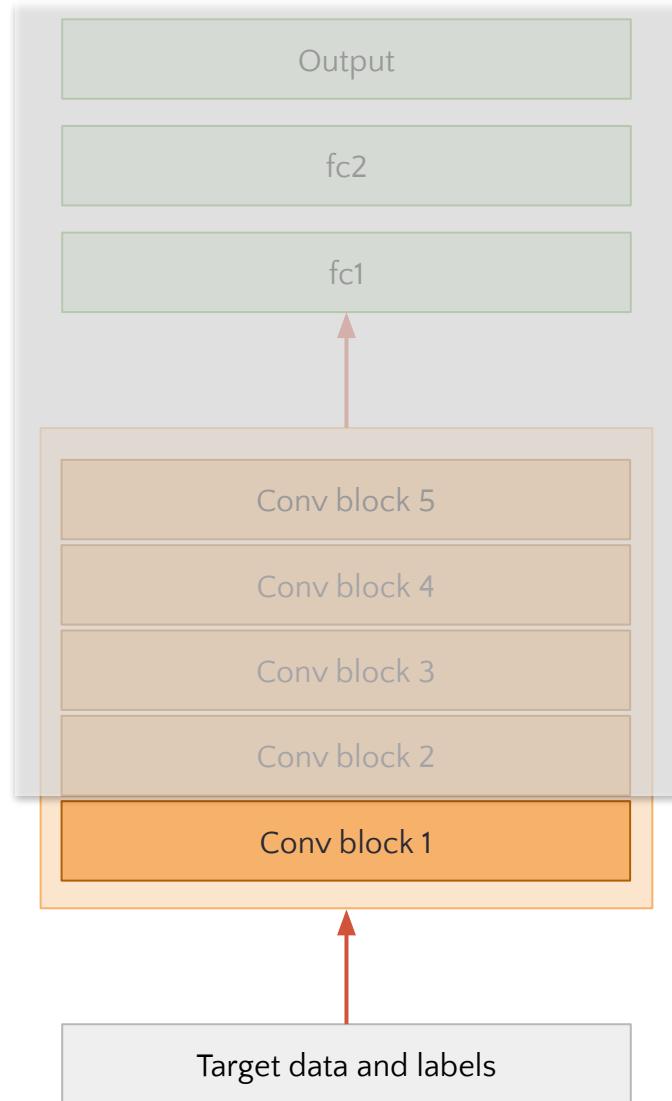


Chain-taw

Freeze weights



2. fine-tunes each layer individually of base model starting from the first to the last.

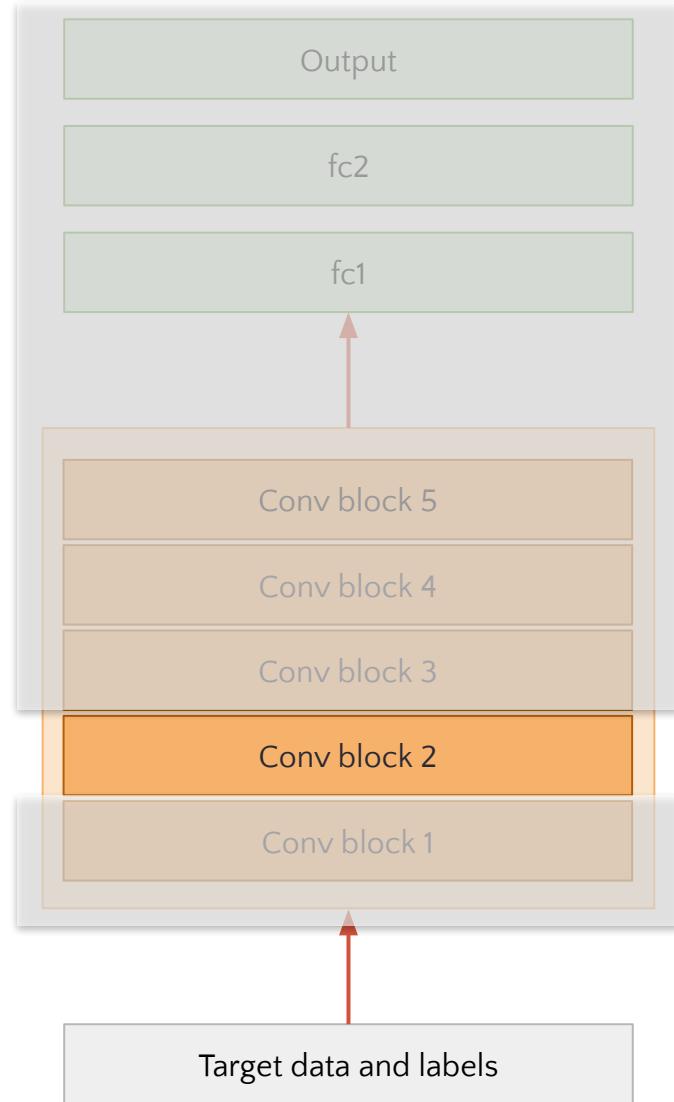


Chain-taw

Freeze weights



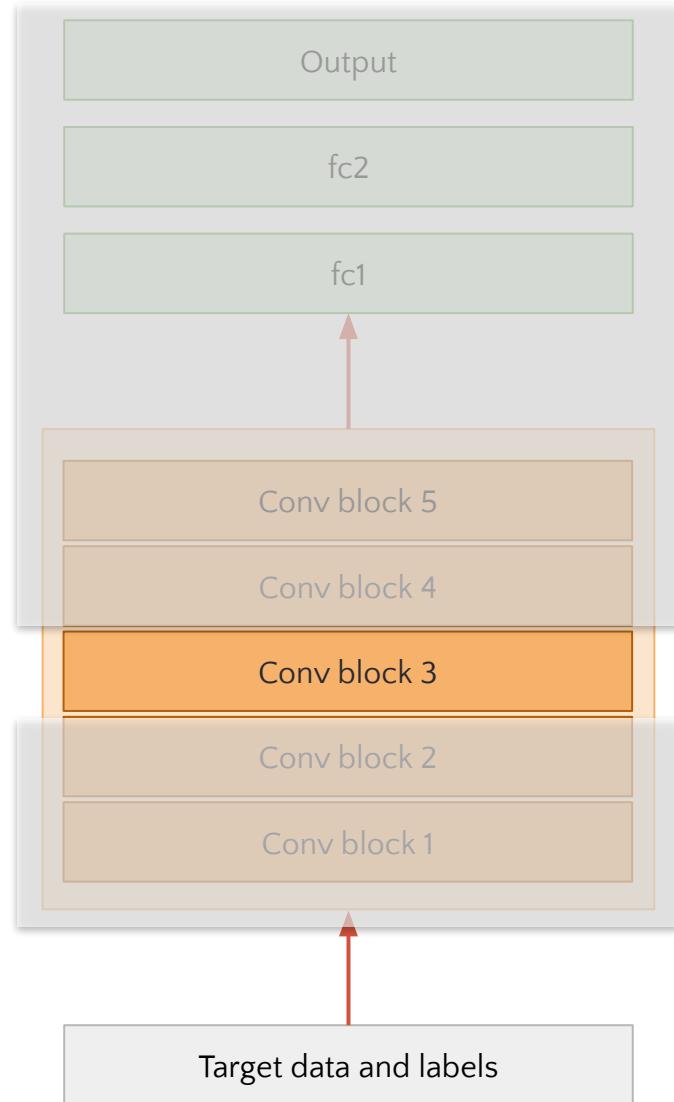
2. fine-tunes each layer
individually of base model
starting from the first to the last.



Chain-taw

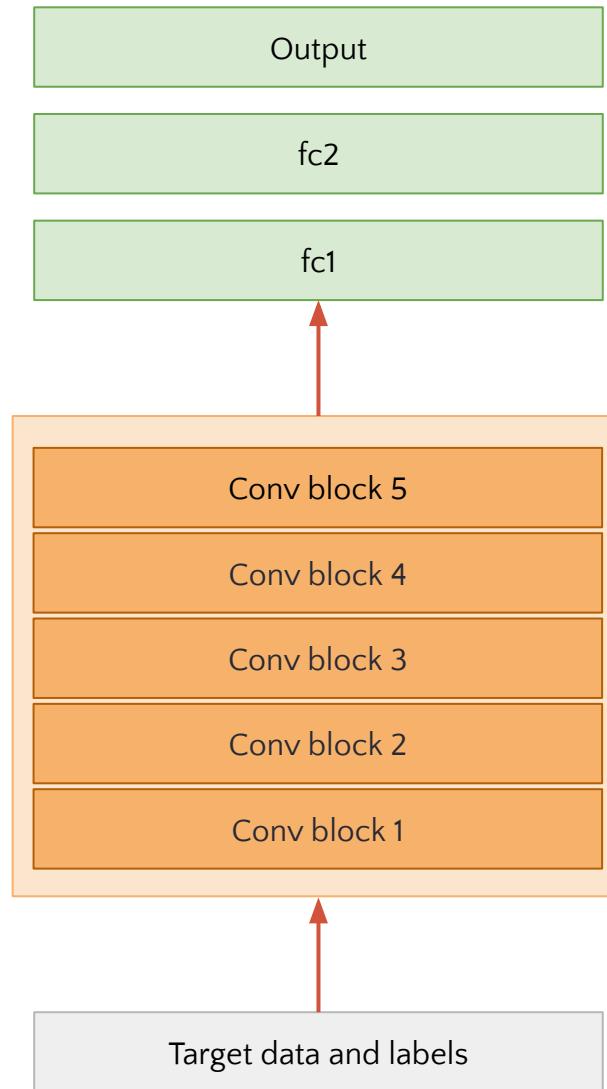
Freeze weights

2. fine-tunes each layer individually of base model starting from the first to the last.



Chain-taw

3. the entire model is trained with all layers.



Lab

Housing price prediction

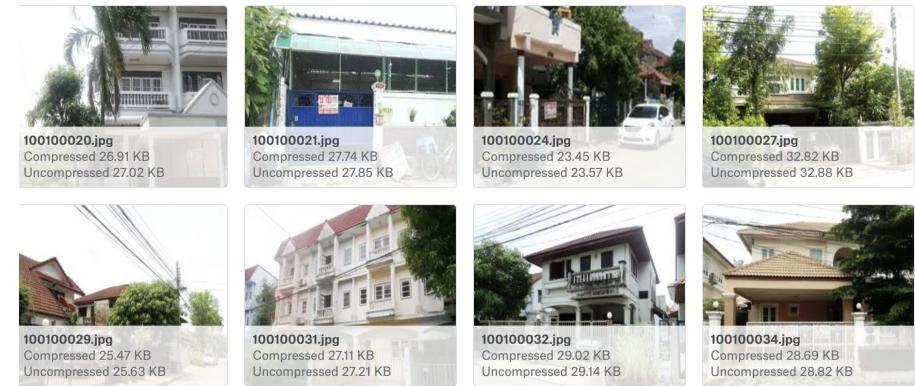


Structured vs unstructured data

Trainset_Homedottech_Hackathon.csv (1.18 MB)

20 of 36 columns

ListingID	ListingInfoID	ListingTypeID	BuildingNameTH	MaxRentPrice
100m	110m	1.00b	1.10b	1 5
1	110053340	1100053070	1	5000000.0
2	100500412	1005004121	1	
3	100100069	1001000691	1	
4	110025306	1100025175	1	
5	110039244	1100039089	1	
6	101400073	1014000731	1	
7	110049848	1100049613	1 မ.န်ဘား ဘုရားလ	
8	110048432	1100048229	1	
9	110002794	110002715	1	
10	110013078	1100012954	1	
11	110053026	1100052756	2	800000.0
12	110048599	1100048396	1	5500000.0



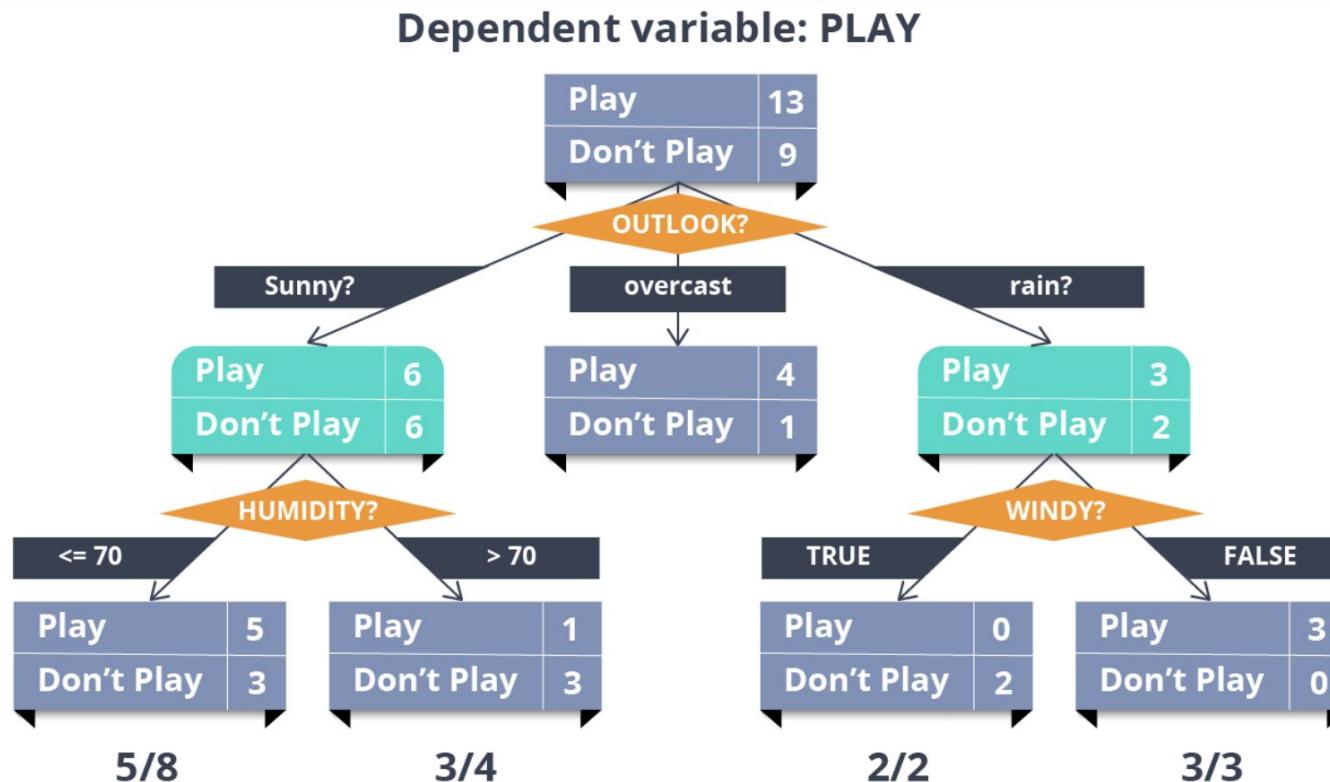
A combination of non deep learning method (XGBoost) + Deep learning
XGBoost - good for tabular data
Deep learning - good for unstructured data

Decision Trees

Introduction to Extreme Gradient Boosting

Decision Trees

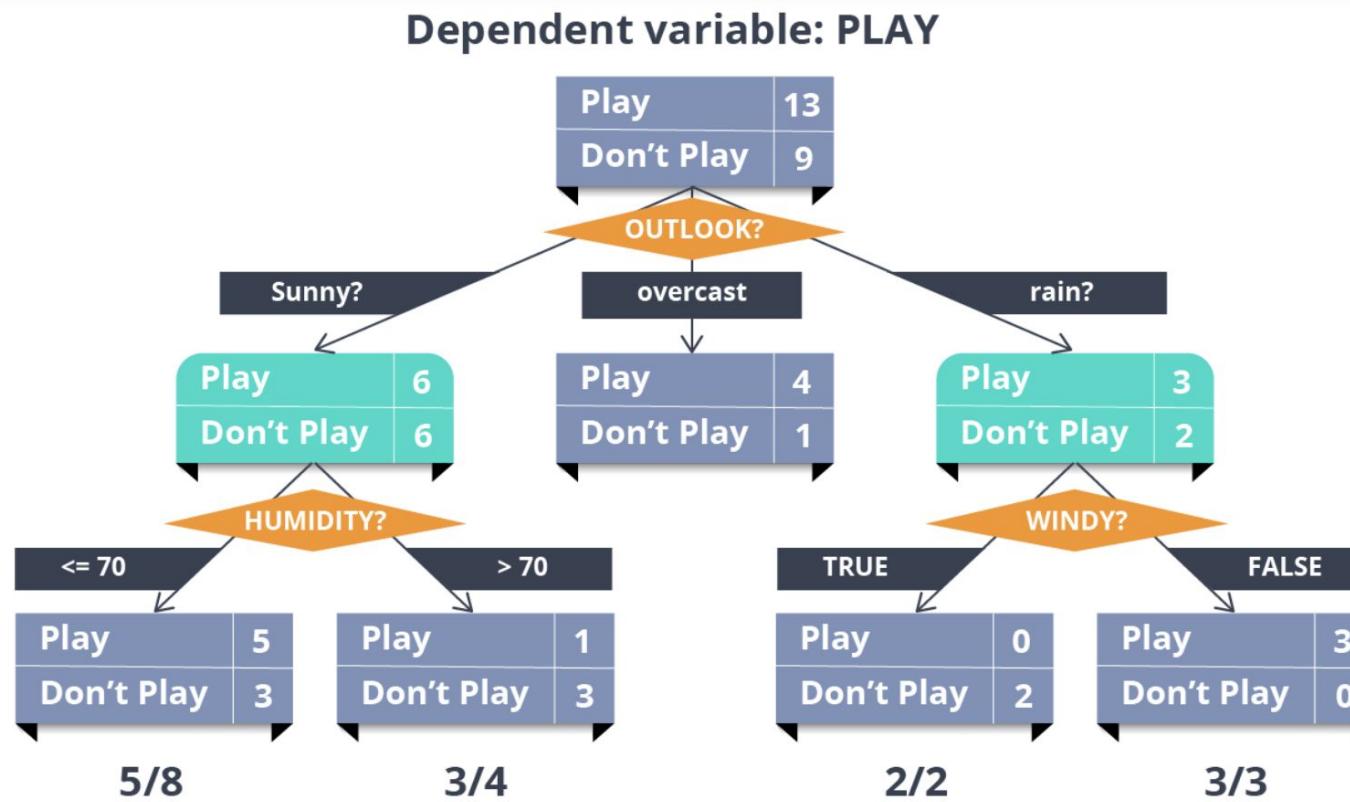
A tree structure that separates data into groups by the feature attributes
Can be used for classification and regression



What's a good decision tree?

Separates the data nicely

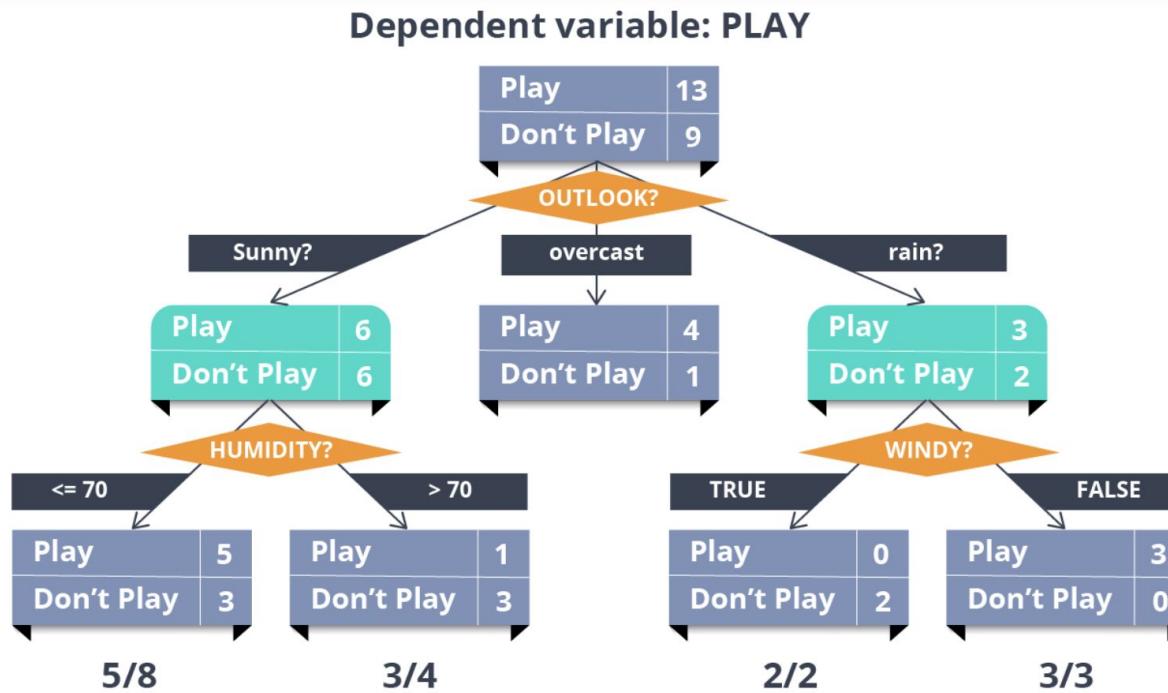
Within a certain budget (smaller trees) - less overfitting



How to create a good decision tree?

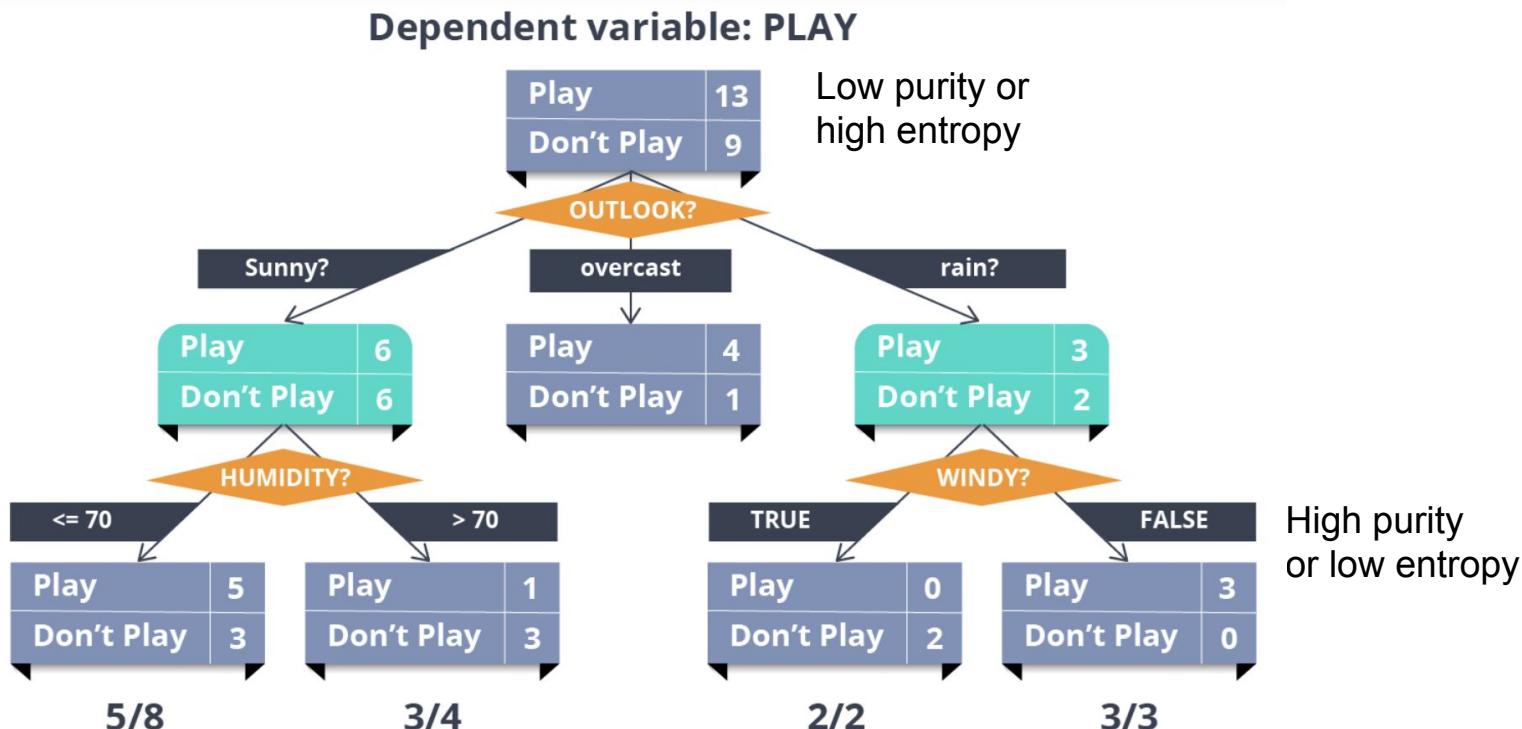
- Pick the attribute that best separates the classes
- Keep doing it until a leave contains entirely one class or you decide it's not worth it to add more nodes

How to determine the best attribute automatically?



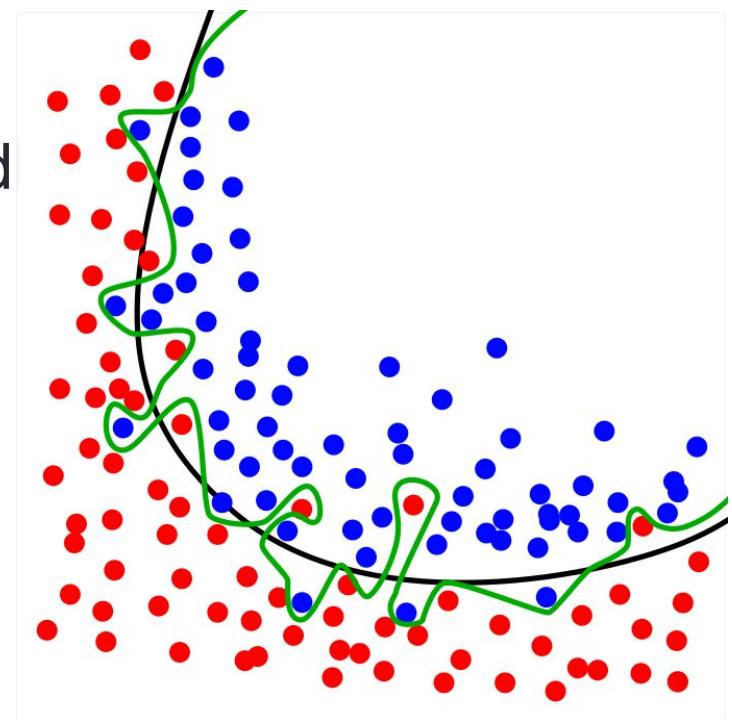
Purity (Entropy)

Want trees that give high purity with the minimum number of nodes



Problems with decision trees

- Can overfitting easily
- Susceptible to noise or bad labelled data



Tree ensemble model

Tree ensemble model

Ensemble types are models that **combine multiple** models together

A group of experts voting on a subject
Can lead to less overfitting

Tree ensemble = Multiple trees = Random Forest!

Bagging

Create multiple subsets of data

Each subset is used to train a different tree

The final answer is the average or mode

Less overfitting and can handle mislabeled data

Random Forest

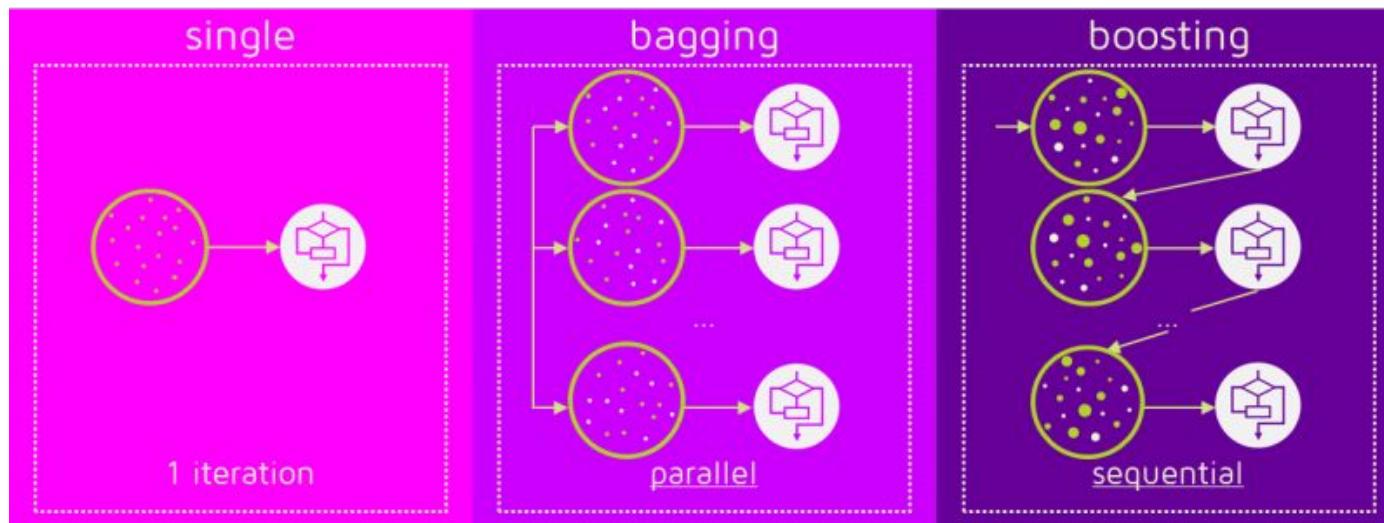
We can also use bagging on features

Each tree has **different training samples AND set of features**

Boosting vs Bagging

Boosting is another way to create multiple trees

But boosting is iterative, the next tree is based on the errors from the previous trees



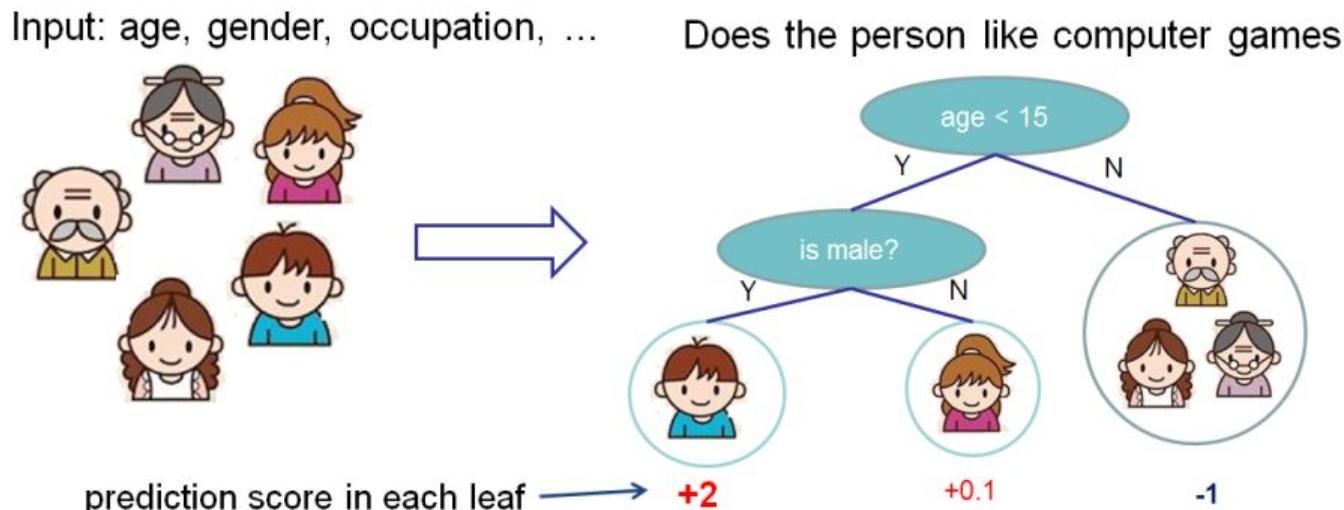
Gradient Boosting

A method of boosting that use gradient-based methods

Tree Gradient Boosting

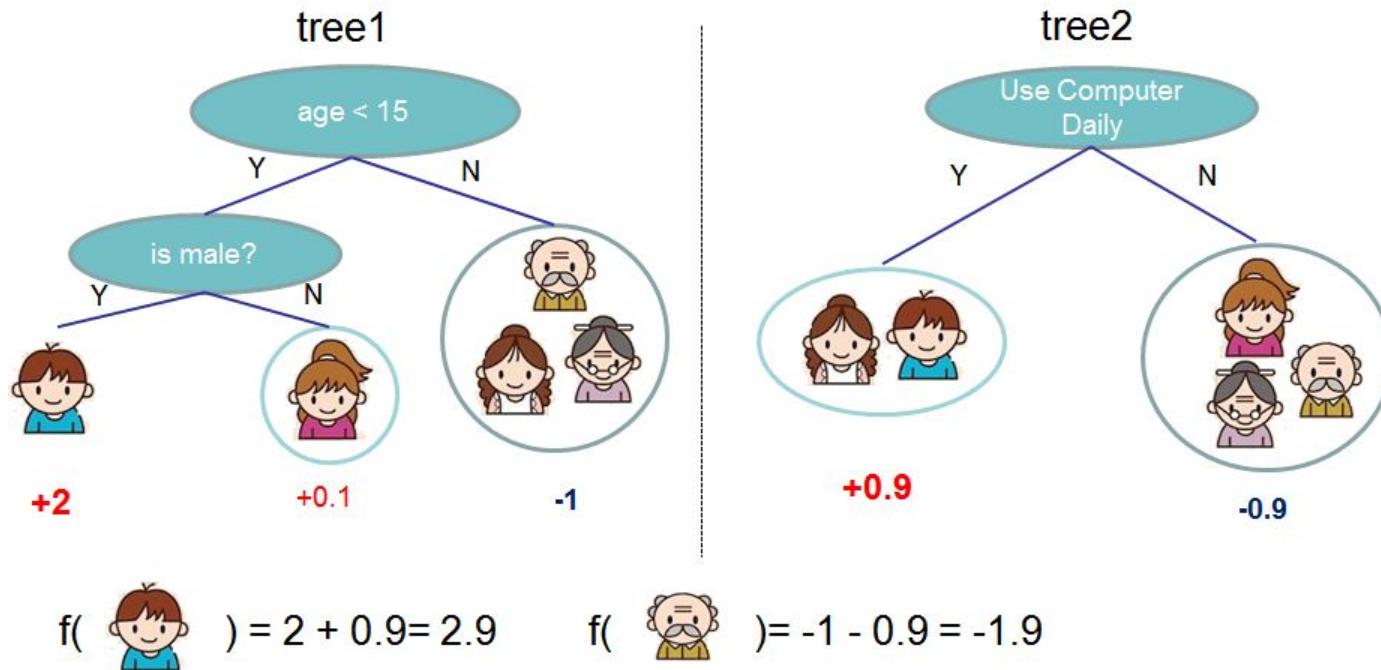
Similar to decision tree

Difference is the leaf node contains a score



Tree Gradient Boosting

Multiple trees with different rules. The subsequent tree try to correct the errors from the previous trees



Extreme Gradient Boosting (XGBoost)

Super popular Tree Boosting library

Highly recommended for spreadsheets type of input data

```
model = XGBClassifier(  
    n_jobs=16,  
    n_estimators=400,  
    max_depth=4,  
    objective="binary:logistic",  
    learning_rate=0.07,  
    subsample=0.9,  
    min_child_weight=6,  
    colsample_bytree=.9,  
    scale_pos_weight=0.8,  
    gamma=8,  
    reg_alpha=6,  
    reg_lambda=1.3)
```

Objective <- type of problem you want to solve

Max_depth <- max depth of tree, higher more overfitting

Min_child_weight <- how strong must the leave be, higher less overfitting

Gamma <- when to stop splitting early

Reg_alpha, reg_lambda <- reduce overfitting

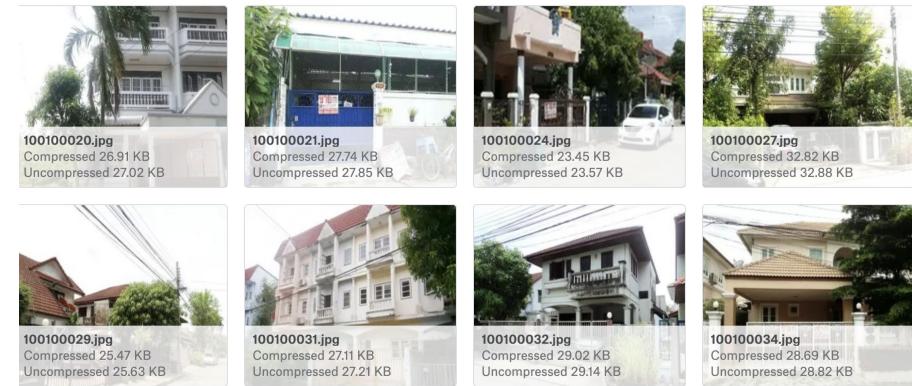
Scale_pos_weight <- weight for class imbalance

Structured vs unstructured data

Trainset_Homedottech_Hackathon.csv (1.18 MB)

20 of 36 columns ▾ Views

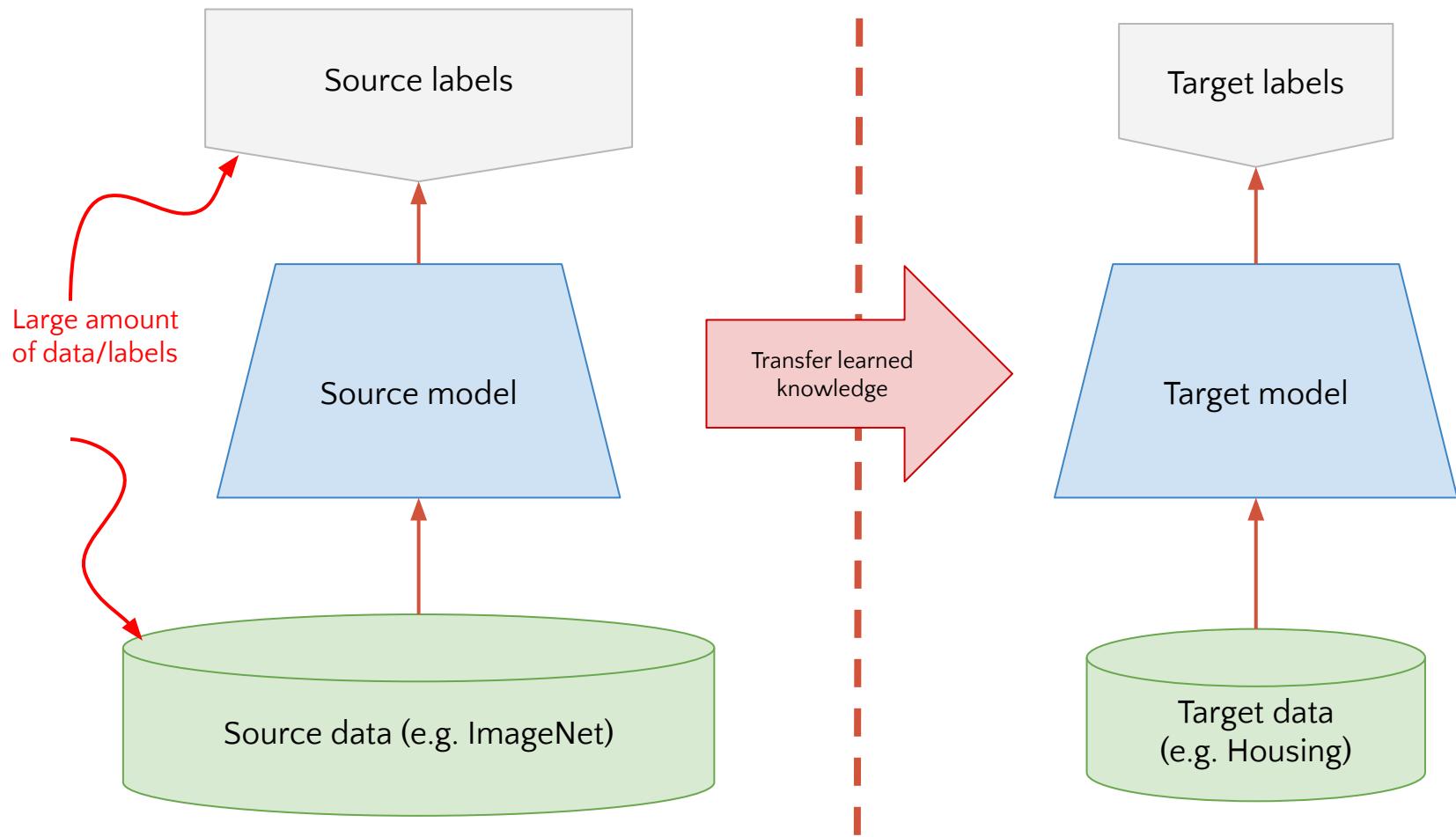
ListingID	ListingInfoID	ListingTypeID	BuildingNameTH	MaxRentPrice
100m	110m	1.00b	1.10b	1 5
1	110053340	1100053070	1	5000000.0
2	100500412	1005004121	1	
3	100100069	1001000691	1	
4	110025306	1100025175	1	
5	110039244	1100039089	1	
6	101400073	1014000731	1	
7	110049848	1100049613	1 ม.นันทวัน วัชรพล	
8	110048432	1100048229	1	
9	110002794	110002715	1	
10	110013078	1100012954	1	
11	110053026	1100052756	2	800000.0
12	110048599	1100048396	1	5500000.0



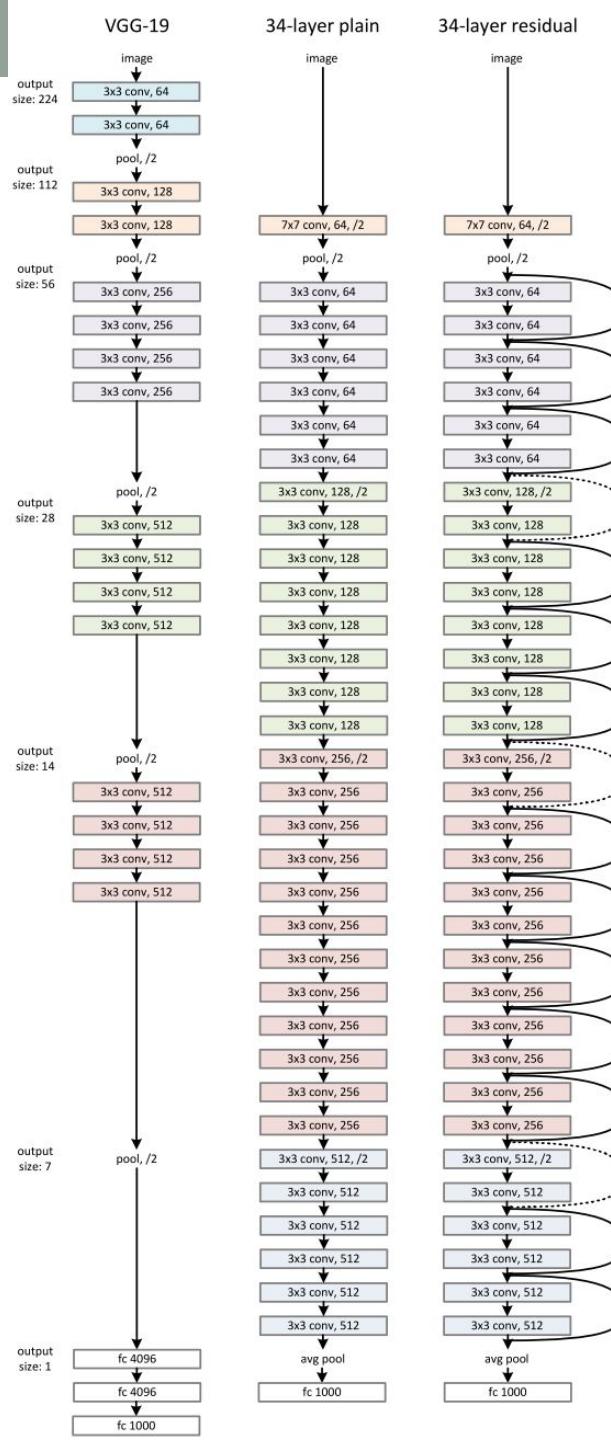
Let's play with XGBoost first

<https://colab.research.google.com/drive/1mVlNEEx8g3nzt8P72JKkduVmRsHoOPlan>

Transfer learning for housing price prediction



Resnet-34



<https://www.kaggle.com/pytorch/resnet34>

Conv layers

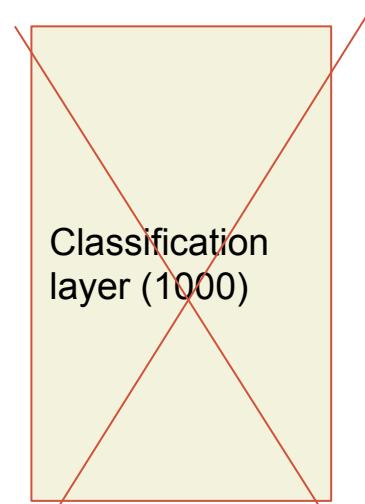
avg pool

Classification
layer (1000)

Conv layers

avg pool

Classification
layer (1000)





124.jpg
Saved 23.45 KB
Resized 25.57 KB

100000027.jpg
Compressed 32.30 KB
Uncompressed 38.40 KB

122.jpg
Saved 29.02 KB
Resized 29.14 KB

100000034.jpg
Compressed 28.69 KB
Uncompressed 28.82 KB

Conv layers

avg pool

Dense
5

Cheap or expensive house?

Conv layers

avg pool

Dense
5

Cheap or expensive house?

Extract these features

+

XGBoost

price



Summary

You can transfer knowledge from other dataset to help reduce the amount of training data.

More sophisticated fine tuning techniques exist:

Chain-taw, ULMfit, etc. See NVIDIA workshop for more details

<https://youtu.be/l8oqxp0up34?t=4410>