



BITS Pilani
Pilani | Dubai | Goa | Hyderabad

DEEP LEARNING

MODULE # 5 : CONVOLUTIONAL NEURAL NETWORK [CNN] ✓

✓
primitive

Conv Layer

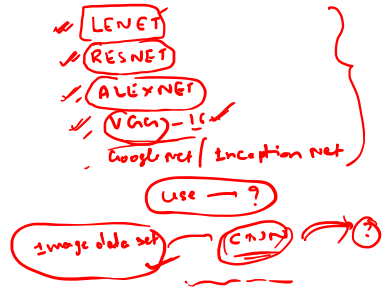
pooling ✓

Fc layer ✓

LENET ✓

- Every convolutional layer includes three parts: convolution, pooling, and nonlinear activation functions.
- Using convolution to extract spatial features.
- Conv filters were 5x5, applied at stride 1.
- Subsampling average pooling layer. Subsampling (Pooling) layers were 2x2 applied at stride 2.
- tanh activation function.
- Using MLP as the last classifier.
- Architecture is [CONV-POOL-CONV-POOL-FC-FC]

Transfer learning ✓



LENET-5

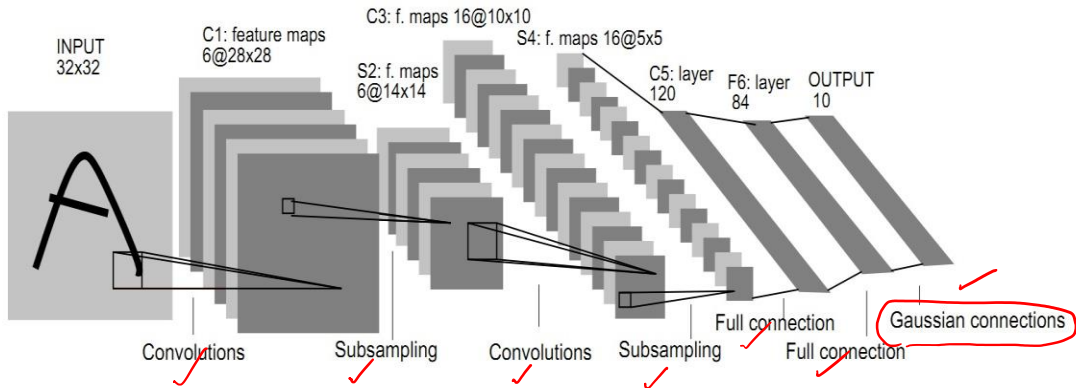
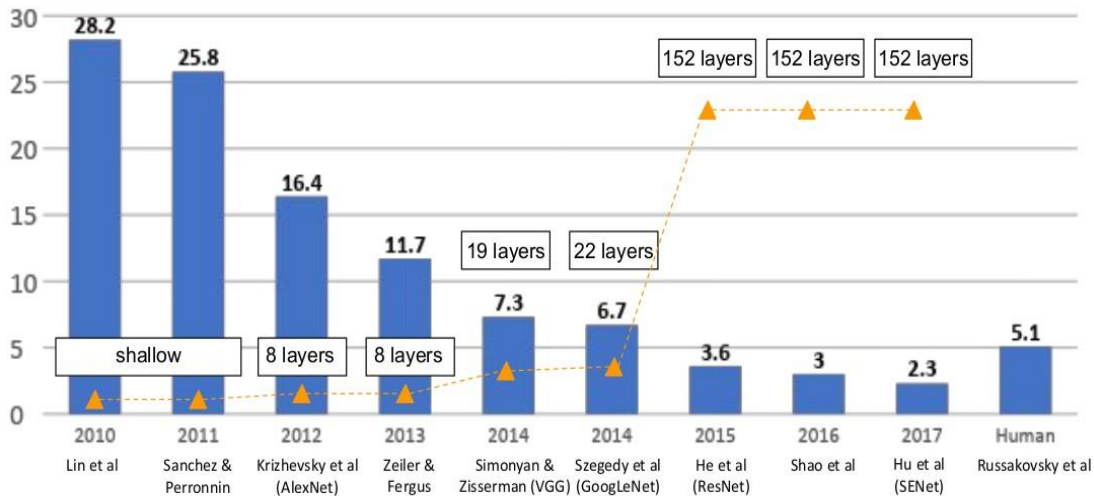


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

IMAGENET LARGE SCALE VISUAL RECOGNITION



IMAGENET DATA

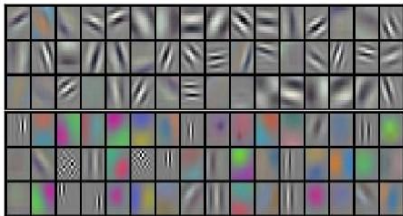


Figure 3: 96 convolutional kernels of size $11 \times 11 \times 3$ learned by the first convolutional layer on the $224 \times 224 \times 3$ input images. The top 48 kernels were learned on GPU 1 while the bottom 48 kernels were learned on GPU

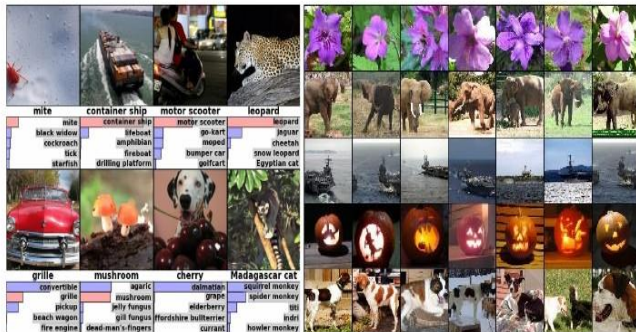


Figure 4: (Left) Eight ILSVRC-2010 test images and the five labels considered most probable by our model. The correct label is written under each image, and the probability assigned to the correct label is also shown with a red bar (if it happens to be in the top 5). (Right) Five ILSVRC-2010 test images in the first column. The remaining columns show the six training images that produce feature vectors in the last hidden layer with the smallest Euclidean distance from the feature vector for the test image.

IMAGENET DATA

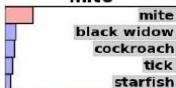


mite

container ship

motor scooter

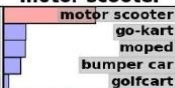
leopard



mite
black widow
cockroach
tick
starfish



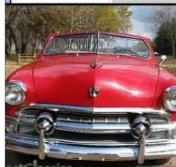
container ship
lifeboat
amphibian
fireboat
drilling platform



motor scooter
go-kart
moped
bumper car
golfcart



leopard
jaguar
cheetah
snow leopard
Egyptian cat



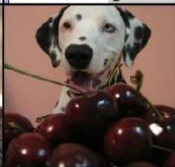
grille

convertible
grille
pickup
beach wagon
fire engine



mushroom

agaric
mushroom
jelly fungus
gill fungus
dead-man's-fingers



cherry

dalmatian
grape
elderberry
ffordshire bullterrier
currant



Madagascar cat

squirrel monkey
spider monkey
titi
indri
howler monkey

The ImageNet set that was used has ~1.2 million images and 1000 classes

Accuracy is measured as top-5 performance: Correct prediction if the true label matches one of the top 5 predictions of the model

ALEXNET ✓

- AlexNet competed in the ImageNet Large Scale Visual Recognition Challenge on September 30, 2012.[3] The network achieved a top-5 error of 15.3%, more than 10.8 percentage points lower than that of the runner up.
- Details/Retrospectives:
 - first use of ReLU ✓
 - used Local Response Normalisation (LRN) layers (not common anymore) ✓
 - Approx. 60 million parameters are learned. ✓
 - heavy data augmentation ✓
 - dropout 0.5 ✓
 - batch size 128 ✓
 - SGD Momentum 0.9 ✓
 - Learning rate $1e-2$, reduced by 10 manually ✓
 - L2 weight decay $5e-4$ ✓
 - 7 CNN ensemble: 18.2% → 15.4% ✓



extract the features

FC layer

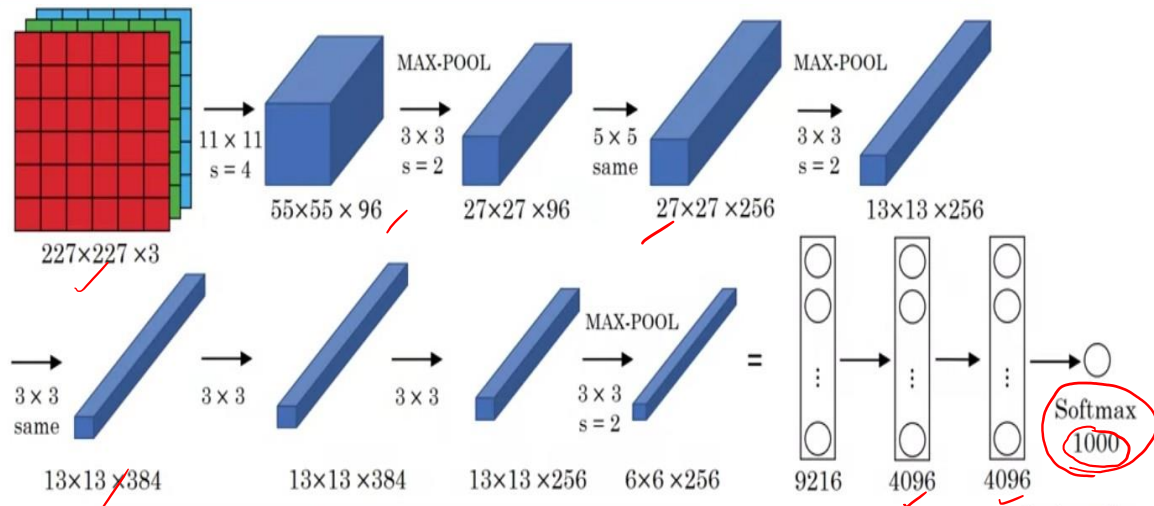
former prop }
Back prop }

ReLU

ALEXNET ARCHITECTURE

- 227x227x3 INPUT
- 55x55x96 CONV1 : 96 11x11 filters at stride 4, pad 0 27x27x96
- MAX POOL1 : 3x3 filters at stride 2 27x27x96
- NORM1 : Normalization layer 27x27x256
- CONV2 : 256 5x5 filters at stride 1, pad 2 13x13x256
- MAX POOL2 : 3x3 filters at stride 2 13x13x256
- NORM2 : Normalization layer 13x13x384
- CONV3 : 384 3x3 filters at stride 1, pad 1 13x13x384
- CONV4 : 384 3x3 filters at stride 1, pad 1 13x13x256
- CONV5 : 256 3x3 filters at stride 1, pad 1 6x6x256
- MAX POOL3 : 3x3 filters at stride 2
- 4096 FC6 : 4096 neurons
- 4096 FC7 : 4096 neurons
- 1000 FC8 : 1000 neurons (class scores)

ALEXNET



[Krizhevsky et al., 2012. ImageNet classification with deep convolutional neural networks]

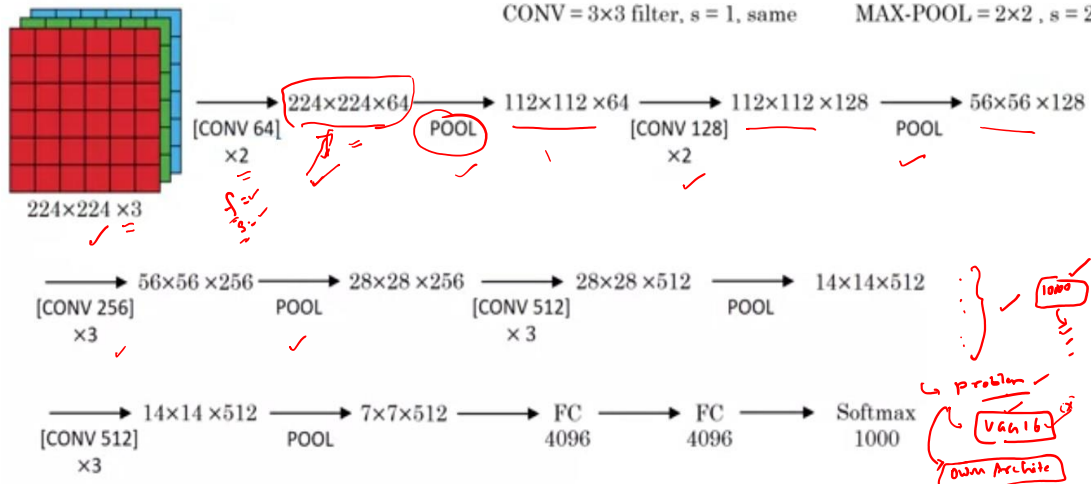
Andrew Ng

- VGG stands for Visual Geometry Group with 16 Layers
- Plug and play in Caffe
- Deeper the better
- Details
 - ILSVRC'14 2nd in classification, 1st in localization
 - Similar training procedure as Krizhevsky 2012
 - Approx. 138 million parameters are learned.
 - No Local Response Normalisation (LRN)
 - Use VGG16 or VGG19 (VGG19 only slightly better, more memory)
 - Use ensembles for best results
 - All convolutions with a 3×3 kernel
 - All max-pooling layers with a 2×2 kernel
 - FC7 features generalize well to other tasks

VGG 16

CONV = 3×3 filter, $s = 1$, same

MAX-POOL = 2×2 , $s = 2$



[Simonyan & Zisserman 2015. Very deep convolutional networks for large-scale image recognition]

Andrew Ng

GOOGLE NET / INCEPTION NET

- **ILSVRC'14** classification winner (6.7% top 5 error) **22** layers with weights
- **Only** 5 million parameters (12x less than AlexNet and 27x less than VGG-16)
- **Inception** Module - convolutional “blocks” – efficient
 - Design a good local network topology (network within a network) and then stack these modules on top of each other.
- **Linear** layers at the end
- **Max** pooling in between, multiple Conv layers between pooling
- **Great** ideas for data augmentation
- **Deeper** networks, with computational efficiency
- **No** FC layers.
- **After** the last convolutional layer, a global average pooling layer is used that spatially averages across each feature map, before final FC layer.

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

Transfer Learning with CNNs

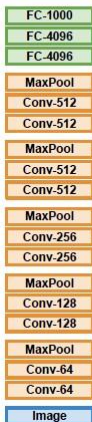
1. Train on Imagenet ✓



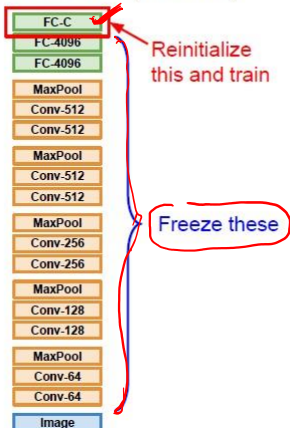
Transfer Learning with CNNs

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

1. Train on Imagenet



2. Small Dataset (C classes)

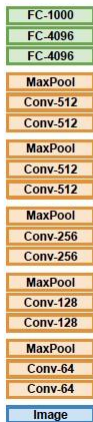


TRANSFER LEARNING

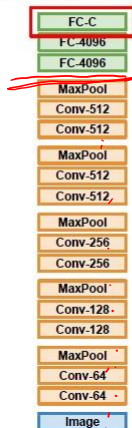
Transfer Learning with CNNs

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

1. Train on Imagenet



2. Small Dataset (C classes)



Reinitialize
this and train

Freeze these

$$\frac{\text{softmax}}{e^{x_1} + e^{x_2} + \dots + e^{x_{1000}}}$$

3. Bigger dataset



Train these

With bigger
dataset, train
more layers

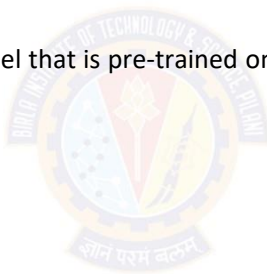
Freeze these

Lower learning rate
when finetuning;
1/10 of original LR
is good starting
point

v 6.16

TRANSFER LEARNING FOR IMAGE DATA

- Use a deep learning model that is pre-trained on large dataset like ImageNet or MS Coco.
- Oxford VGG Model
- Google Inception
- Model Microsoft
- ResNet Model



TRANSFER LEARNING FOR TEXT DATA

- Embedding is the mapping of words to a high-dimensional continuous vector space where different words with similar meanings have similar vector representations.
- Google's word2vec Model
- Stanford's GloVe Model
- FastText
- Gensim



TRANSFER LEARNING - WHEN TO USE?

- You need a lot of a data if you want to train/use CNNs / RNNs.
- Task A and Task B have the same type of input. Eg: Input is images for both tasks.
- We have lot of data for training Task A and relatively low data for training in Task B.
- Low level features obtained from Task A could be more helpful for learning Task B.



Thank You!