

Lab. – Classification

Convolutional Neural Networks (CNNs)

Yuan-Fu Liao

National Yang Ming Chiao Yung University

Kaggle Competition - The Simpsons Characters Recognition Challenge III – noisy images

- Images of 50 characters extracted from The Simpsons episodes
 - ✓ About 2000 images per character
 - ✓ Pictures are under various size, scenes
 - ✓ not necessarily centered in each image and could sometimes be with or cropped from other characters





Machine Learning 2023@NYCU - Classification

The Simpsons Characters Recognition Challenge III - Noisy Image

22 days to go

abraham_grampa_simpson (2067 files)



About this directory

This file does not have a description yet.



ENQ_DataGen-_0_1004....
14.92 kB



ENQ_DataGen-_0_1007.j...
18.29 kB



ENQ_DataGen-_0_1031.j...
19.23 kB



ENQ_DataGen-_0_1036....
16.45 kB



ENQ_DataGen-_0_1037.j...
16.69 kB



ENQ_DataGen-_0_1039....
14.97 kB



ENQ_DataGen-_0_1042....
16.5 kB



ENQ_DataGen-_0_1047.j...
21.24 kB



ENQ_DataGen-_0_1057.j...
18.21 kB



ENQ_DataGen-_0_1065....
12.14 kB

Data Explorer

1.59 GB

- test-final
- train
 - train
 - abraham_grampa_s...

- ENQ_DataGen-_0_1004....
- ENQ_DataGen-_0_1007.j...
- ENQ_DataGen-_0_1031.j...
- ENQ_DataGen-_0_1036....
- ENQ_DataGen-_0_1037.j...
- ENQ_DataGen-_0_1039....
- ENQ_DataGen-_0_1042....
- ENQ_DataGen-_0_1047.j...
- ENQ_DataGen-_0_1057.j...
- ENQ_DataGen-_0_1065....

20 of 50 Characters in The Simpsons



abraham_grampa_simpson
comic_book_guy

edna_krabappel

apu_nahasapeemapetilon
homer_simpson

bart_simpson
kent_brockman

charles_montgomery_burns
krusty_the_clown

chief_wiggum

lenny_leonard
moe_szyslakned_flanders

lisa_simpson

nelson_muntz

marge_simpson

principal_skinner

mayor_quimby

sideshow_bob
milhouse_van_houten



Demo

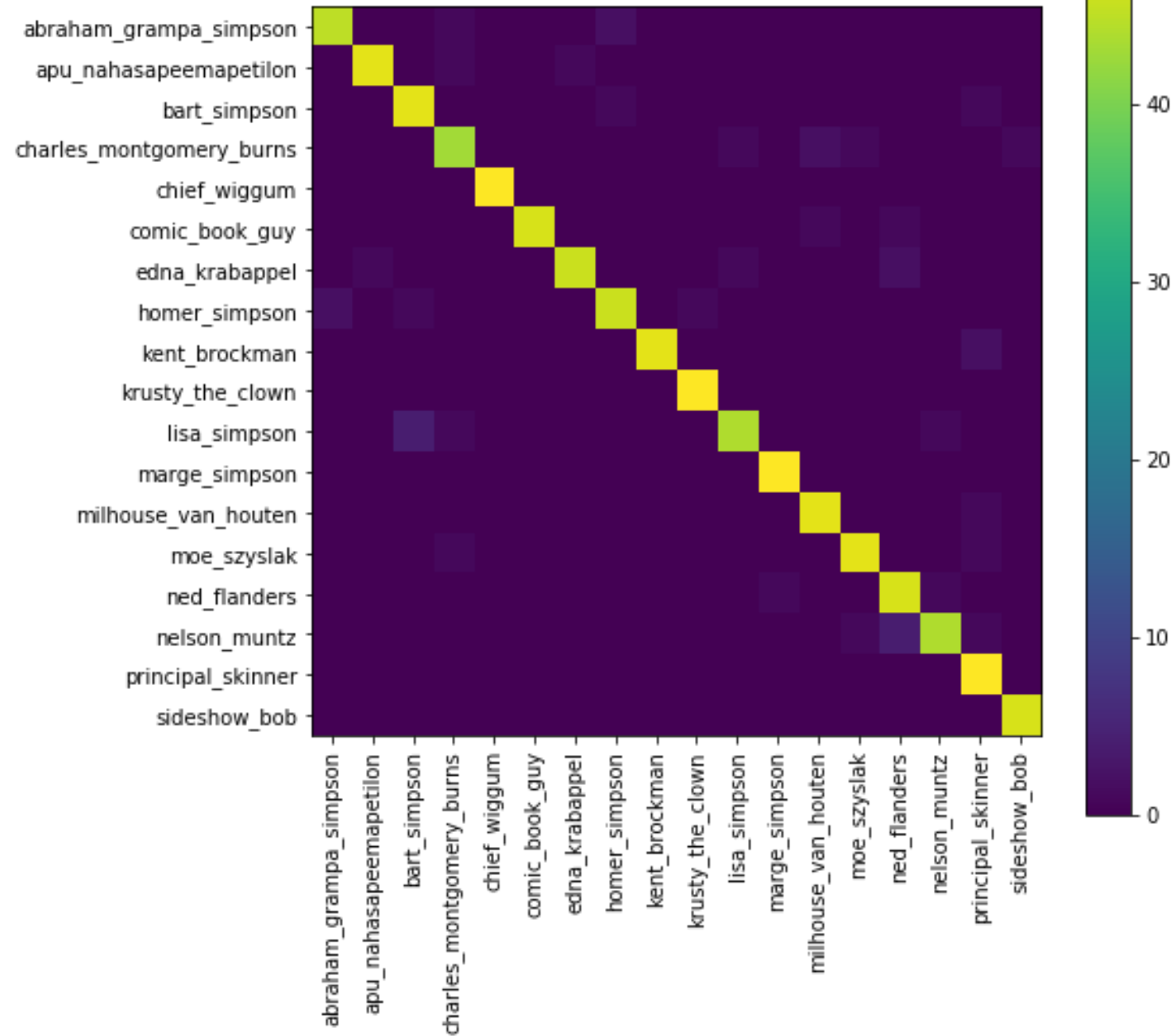


- Task1
 - Predict 50 Simpsons Characters in all pictures



- Task2

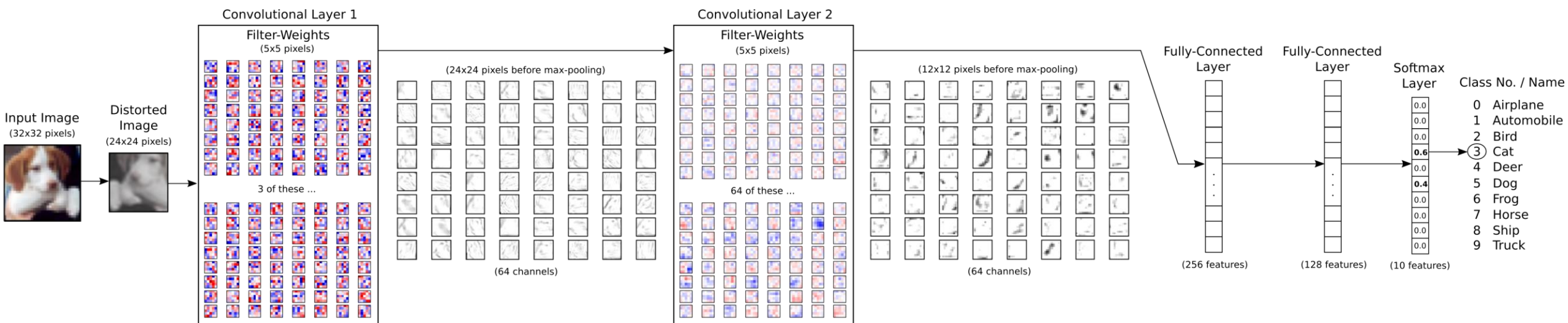
- **Compute the Confusion Matrix (50x50)**



• Task3

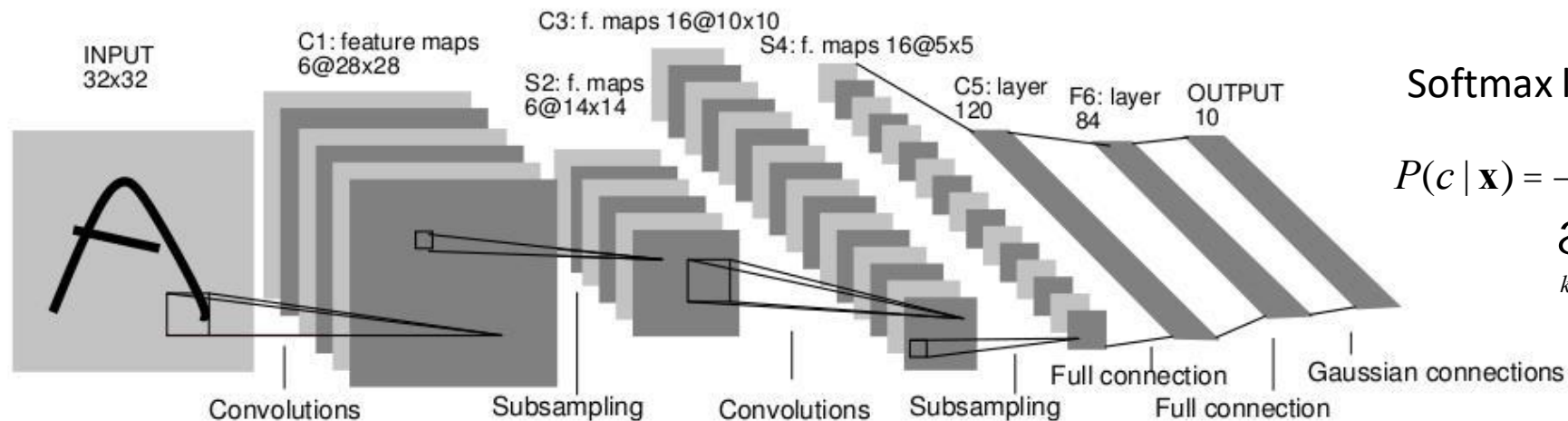
• Visualization and Understanding Convolutional Neural Networks

- 畫出每一層filter的權重
- 畫出每一層的feature map



Convolutional Neural Networks

- Neural network with specialized connectivity structure
- Stack multiple stages of feature extractors
- Higher stages compute more global, more invariant features
- Classification layer at the end



Softmax layer:

$$P(c | \mathbf{x}) = \frac{\exp(\mathbf{w}_c \times \mathbf{x})}{\sum_{k=1}^C \exp(\mathbf{w}_k \times \mathbf{x})}$$

Examples

- TRAINING A CLASSIFIER
 - https://pytorch.org/tutorials/beginner/blitz/cifar10_tutorial.html
- TRANSFER LEARNING FOR COMPUTER VISION TUTORIAL
 - <https://pytorch.org/vision/stable/models.html>

Test Image

- 注意測試圖片加入了翻轉變形、改變顏色跟雜訊干擾等等變化
- 因此，在訓練時，就要考慮做資料擴充



Data Augmentation

1. <https://pytorch.org/vision/stable/transforms.html#>
2. https://pytorch.org/vision/master/auto_examples/transforms/plot_transforms_illustrations.html#sphx-glr-auto-examples-transforms-plot-transforms-illustrations-py

CIFAR-10

- CIFAR-10 classification is a common benchmark problem in machine learning.
- The problem is to classify RGB 32x32 pixel images across 10 categories: airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck

airplane



automobile



bird



cat



deer



dog



frog



horse



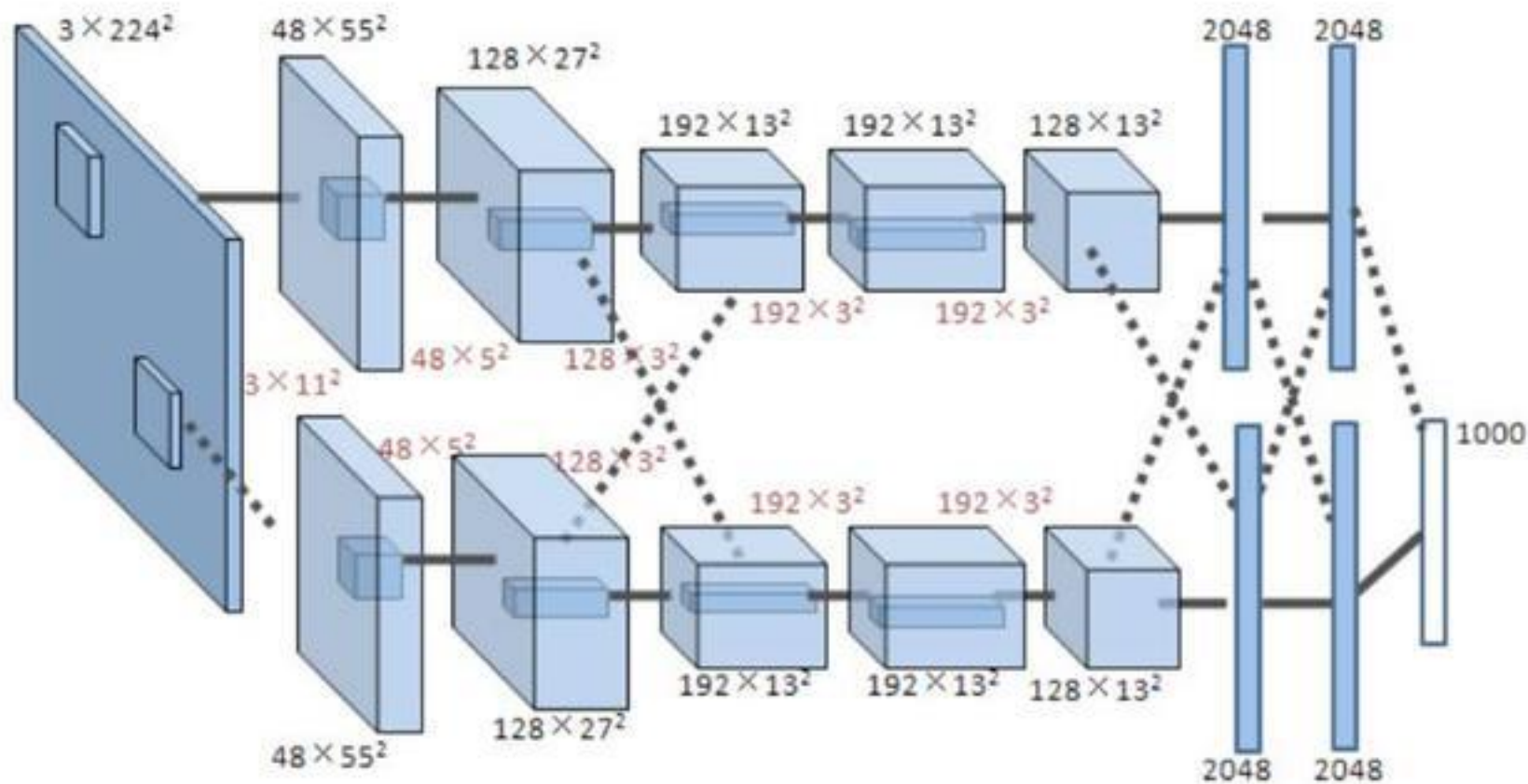
ship



truck



AlexNet



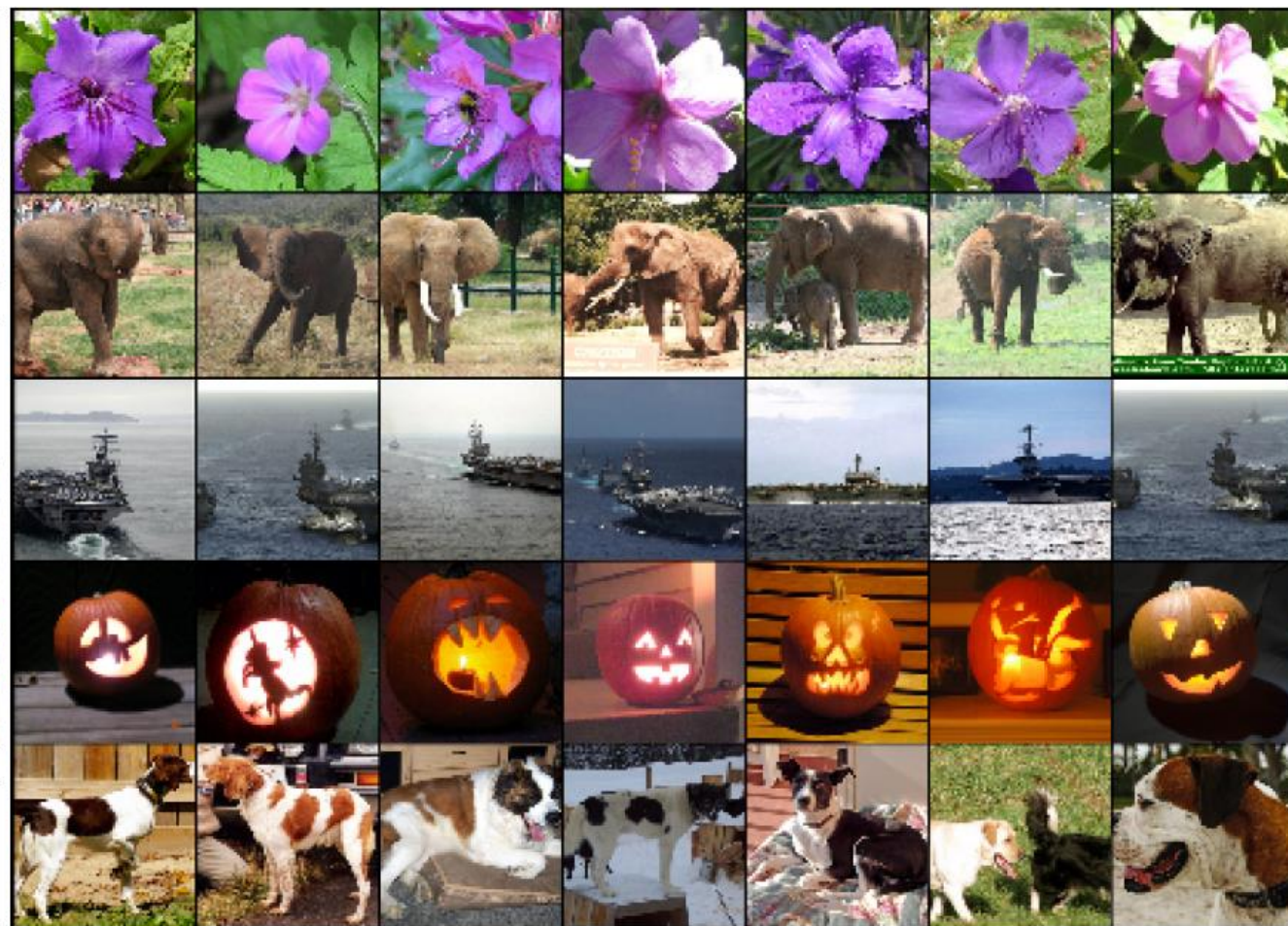
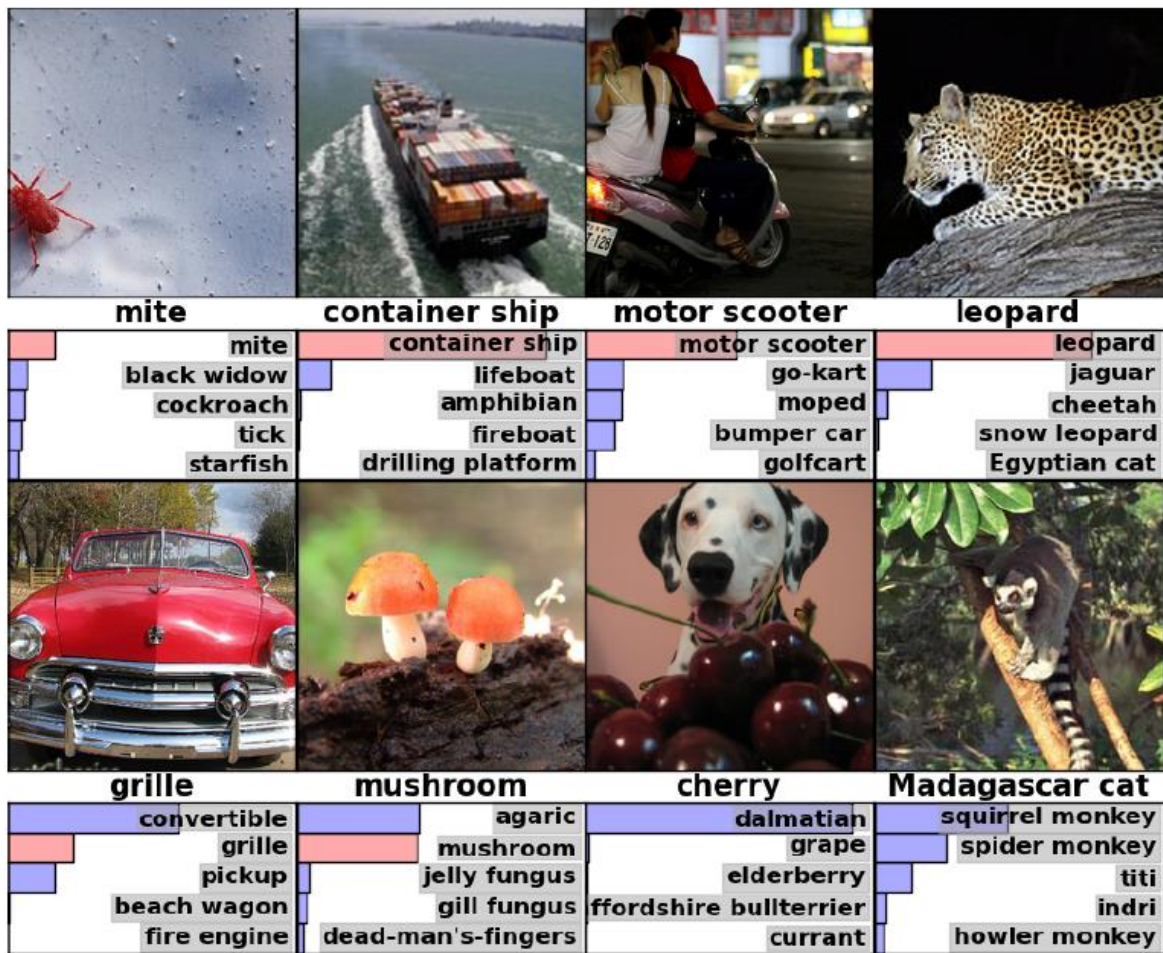


Figure 4: **(Left)** Eight ILSVRC-2010 test images and the five labels considered most probable by our model. The correct label is written under each image, and the probability assigned to the correct label is also shown with a red bar (if it happens to be in the top 5). **(Right)** Five ILSVRC-2010 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.

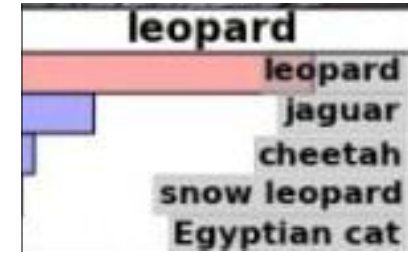
ImageNet Classification with Deep Convolutional Neural Networks

Goa

I



Classifica(on



ImageNet

- Over 15M labeled high resolution images
- Roughly 22K categories
- Collected from web and labeled by Amazon Mechanical Turk



ILSVR C

- Annual competition of image classification at large scale
- 1.2M images in 1K categories
- Classification: make 5 guesses about the image label



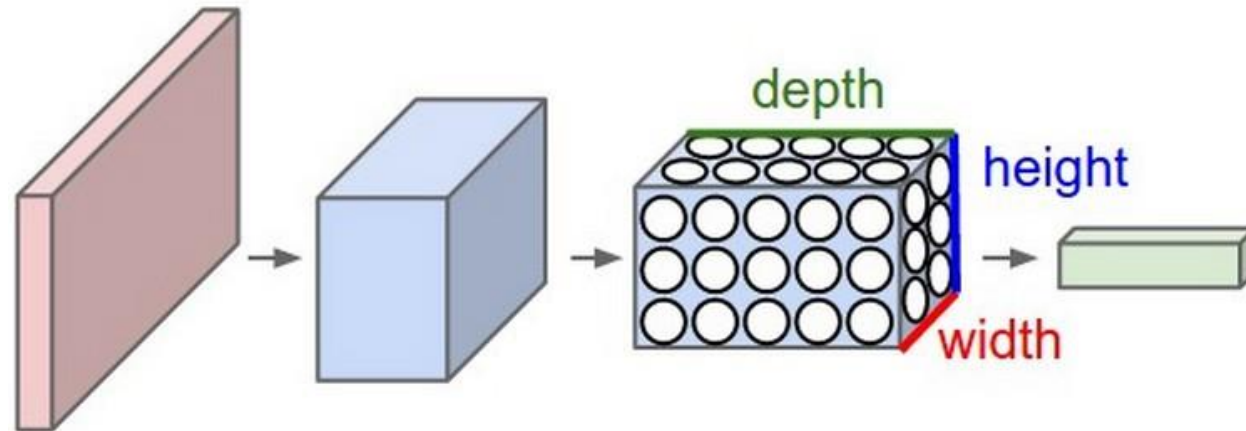
EntleBucher



Appenzeller

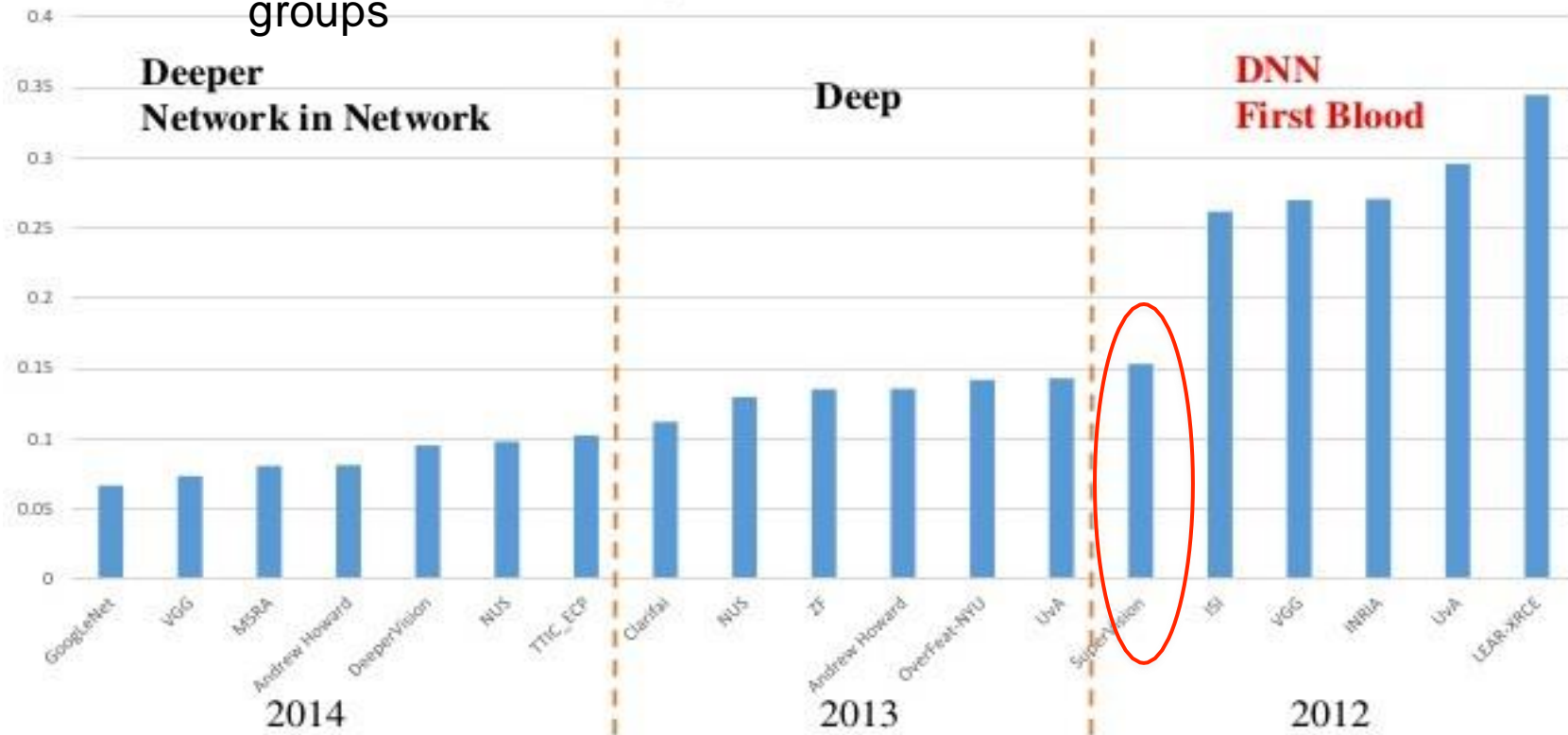
Convolutional Neural Networks

- Model with a large learning capacity
- Prior knowledge to compensate all data we do not have



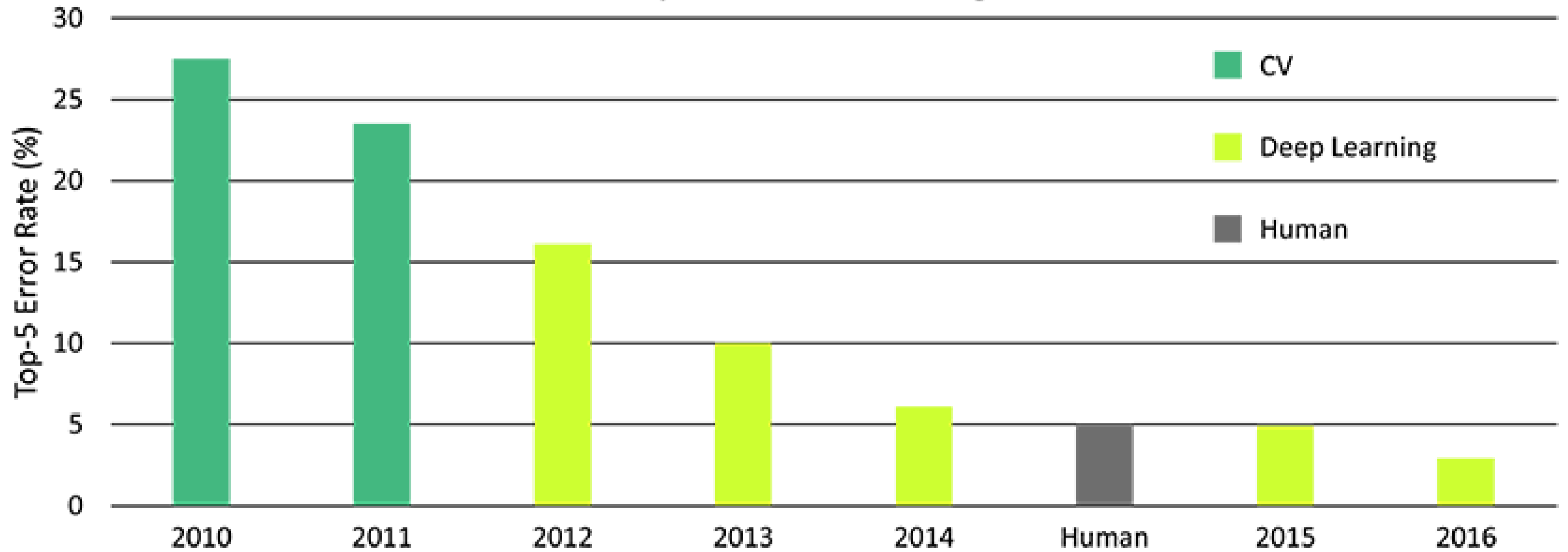
ILSVRC

ImageNet Classification error throughout years and groups

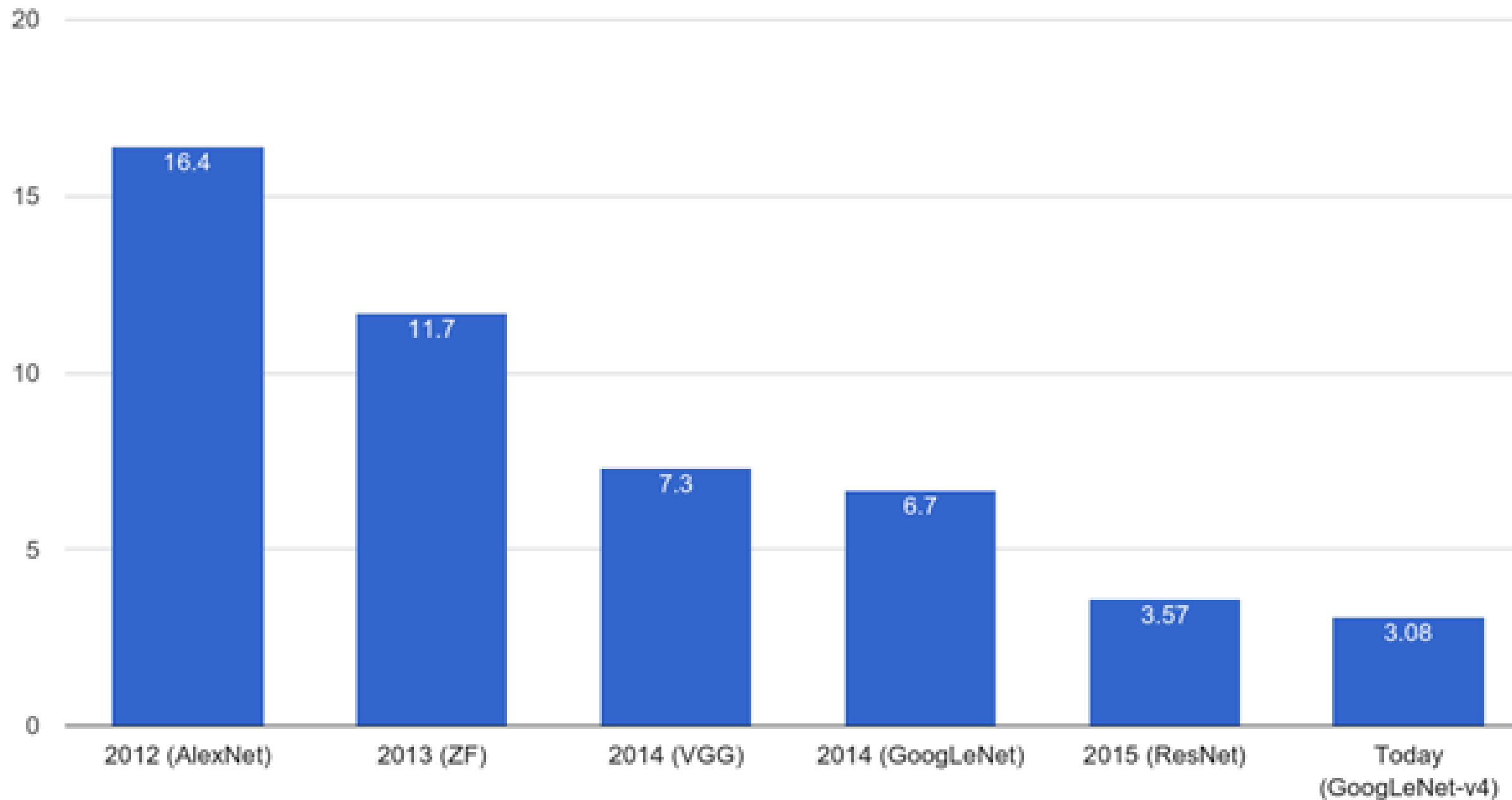


Li Fei-Fei: ImageNet Large Scale Visual Recognition Challenge, 2014 <http://image-net.org/>

ILSVRC Top 5 Error on ImageNet



ImageNet Classification Error (Top 5)

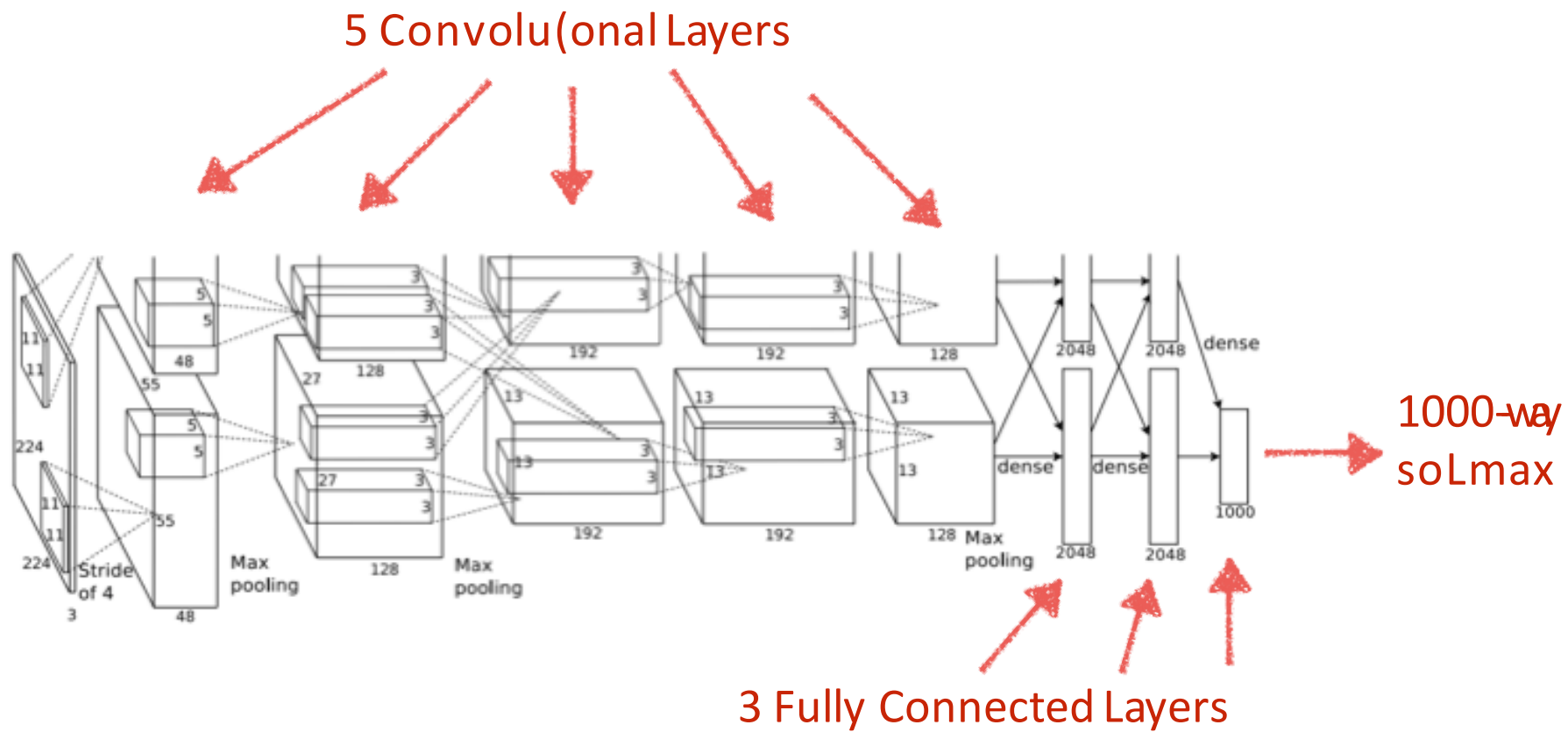


SuperVision (SV)

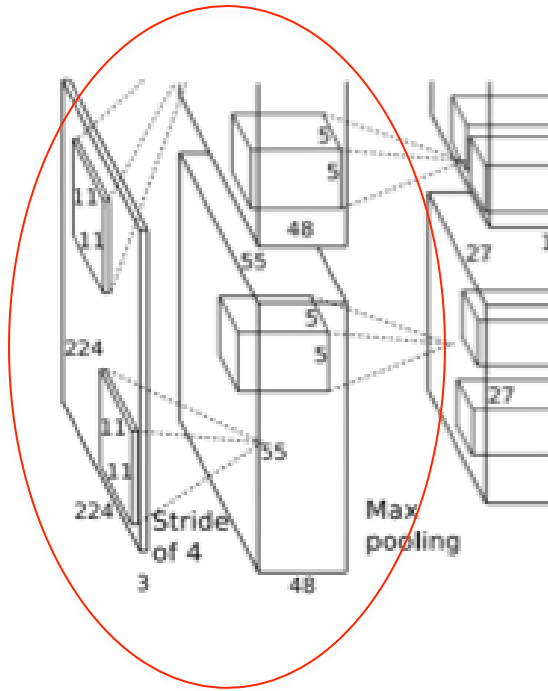
Image classification with deep convolutional neural networks

- 7 hidden “weight” layers
 - 650K neurons
 - 60M parameters
 - 630M connections
-
- Rectified Linear Units, overlapping pooling, dropout trick
 - Randomly extracted 224x224 patches for more data

Architecture

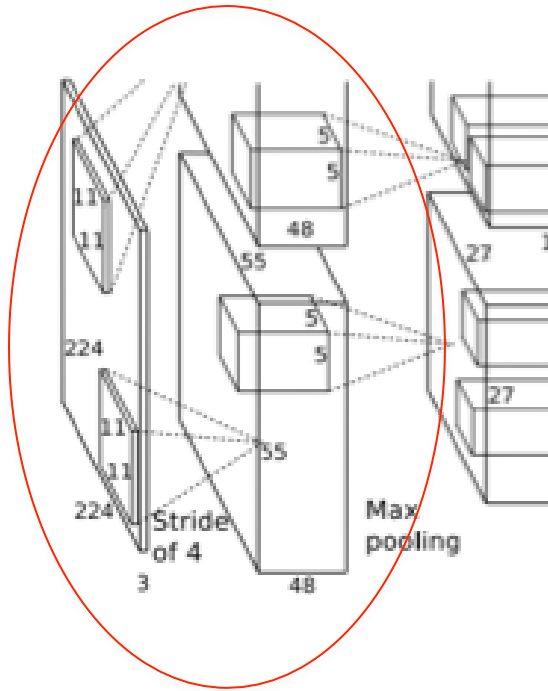


Layer 1 (Convolutional)



- Images: 227x227x3
- F (receptive field size): 11
- S (stride) = 4
- Conv layer output:
55x55x96

Layer 1 (Convolutional)



- $55 \times 55 \times 96 = 290,400$ neurons
- each has $11 \times 11 \times 3 = 363$ weights and 1 bias
- $290400 \times 364 = 105,705,600$ parameters on the first layer of the AlexNet alone!

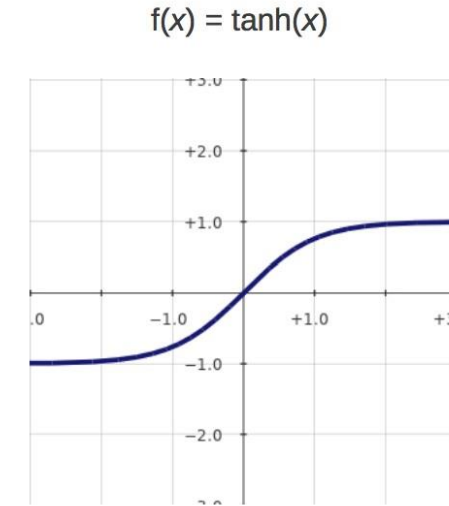
Architecture

ReLU Nonlinearity

- Standard way to model a

neuron $f(x) = \tanh(x)$ or $f(x) = (1 + e^{-x})^{-1}$

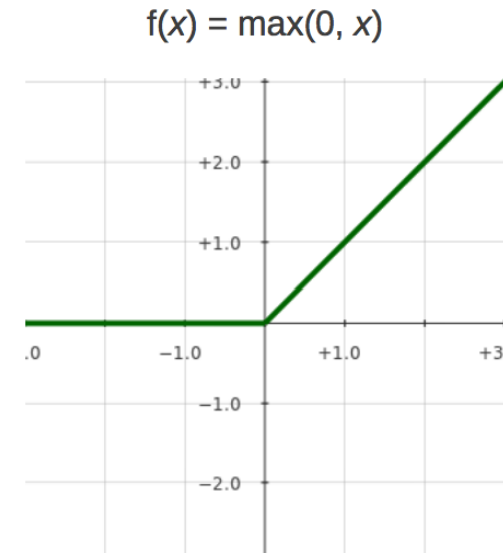
Very slow to
train



- Non-saturating nonlinearity (ReLU)

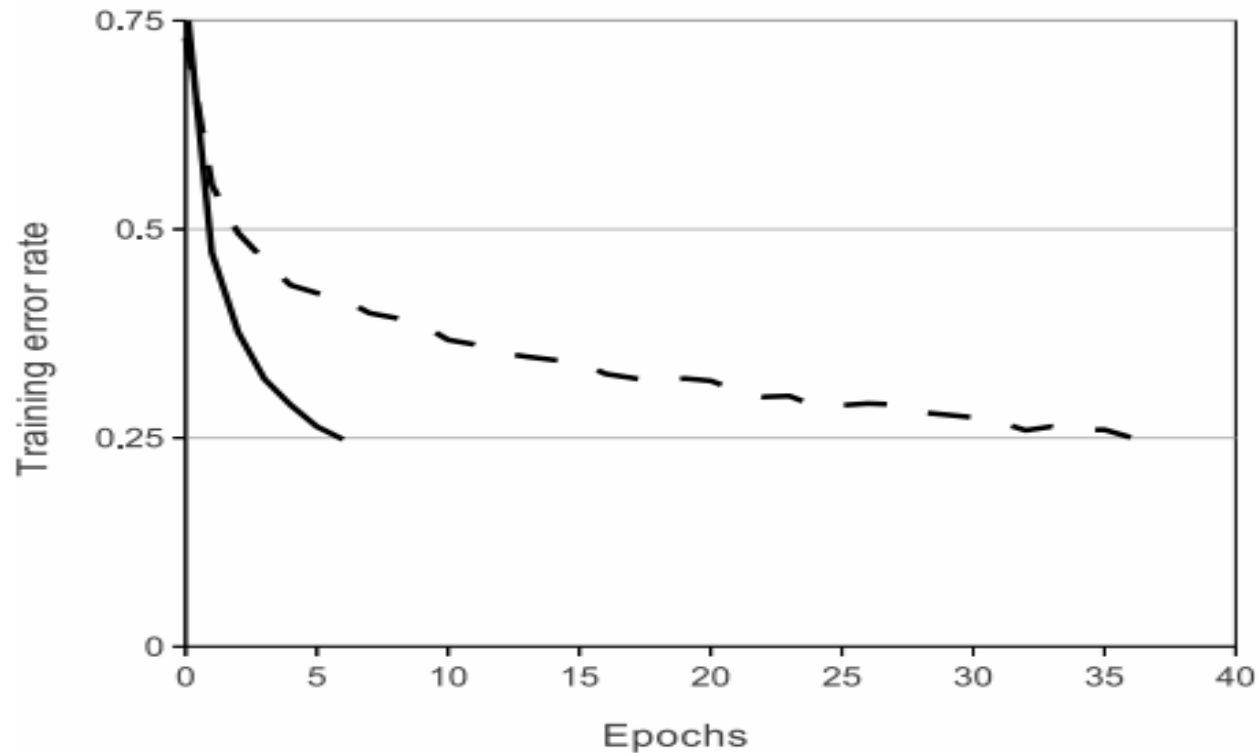
$$f(x) = \max(0, x)$$

Quick to train



Architecture

ReLU
Nonlinearity



A 4 layer CNN with ReLUs (solid line) converges **six times faster** than an equivalent network with tanh neurons (dashed line) on CIFAR-10 dataset

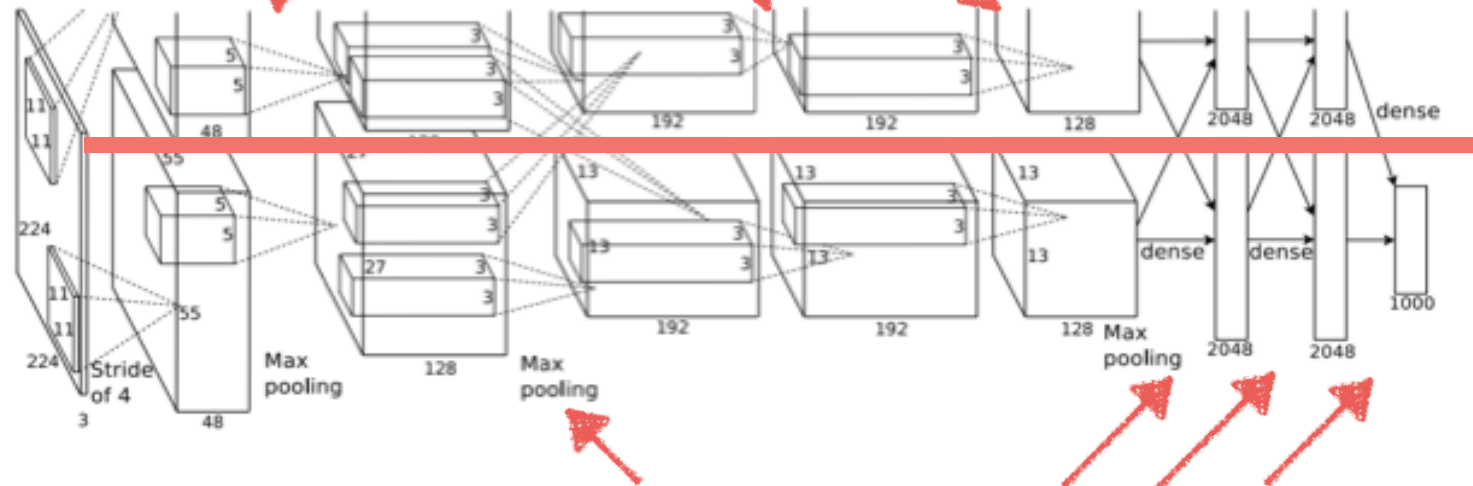
Architecture

Training on Multiple

GPUs

GPU #1

intra-GPU connec(ons



GPU #2

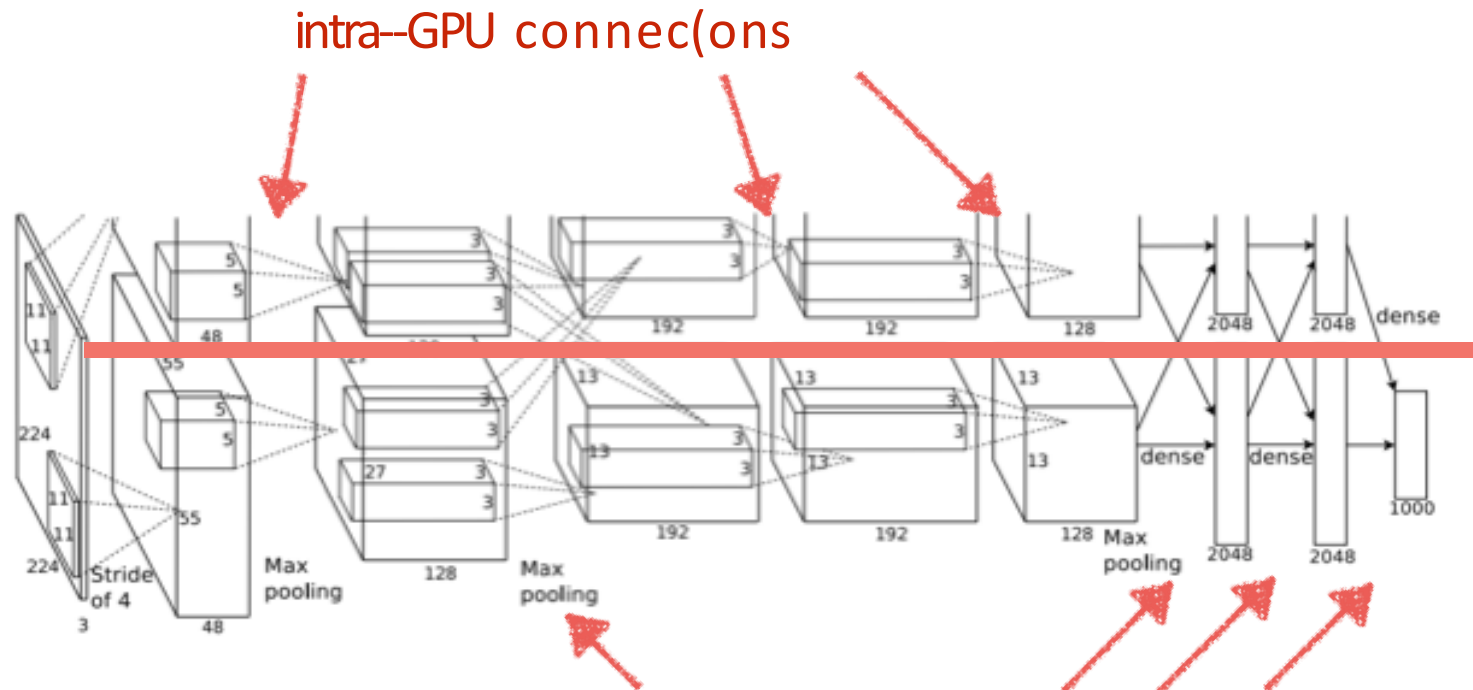
inter-GPU connec(ons

Architecture

Training on Multiple

GPUs

GPU #1



GPU #2

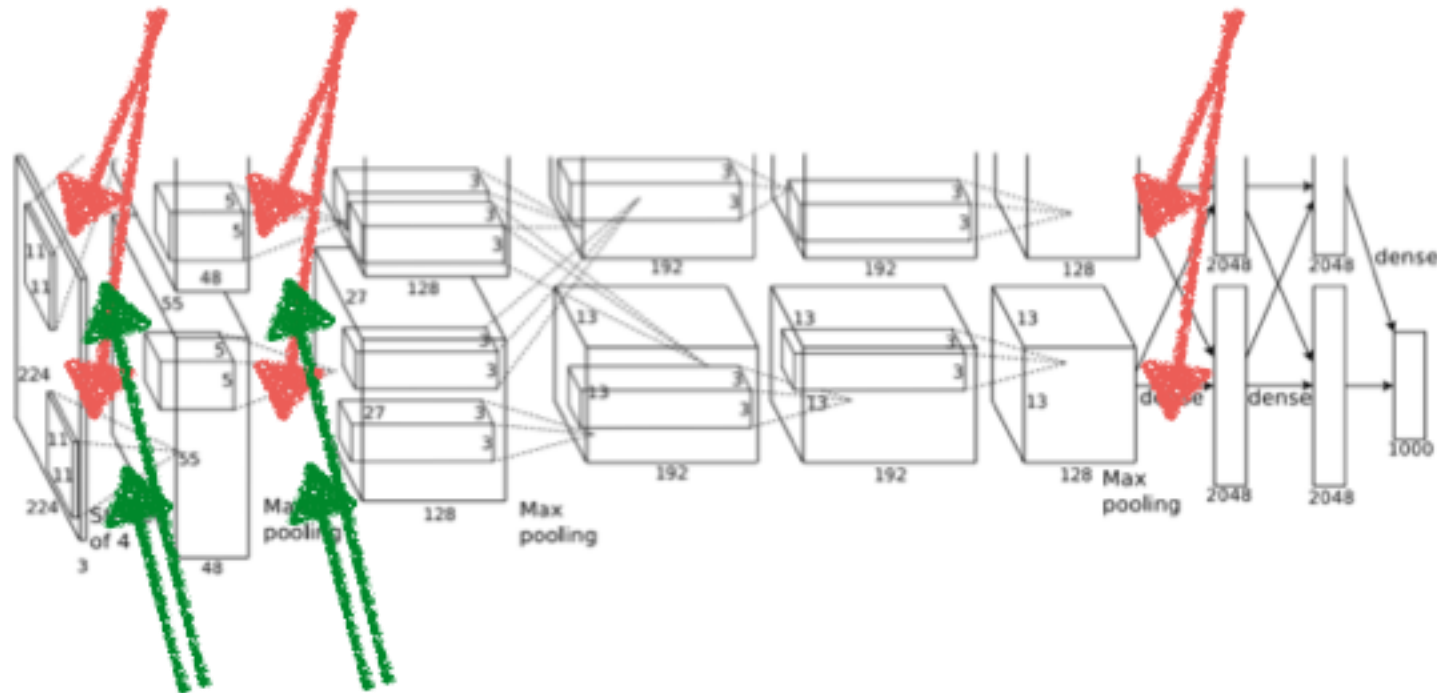
inter-GPU connections

Top-1 and Top-5 error rates decreases by 1.7% & 1.2% respectively, comparing to the net trained with one GPU and half neurons!!

e

Overlapping Pooling

Max-pooling layers



Response normalization layers

Architecture

Local Response Normalization

- No need to input normalization with ReLUs.
- But still the following local normalization scheme helps generalization.

$$b_{x,y}^i = a_{x,y}^i / \left(k + \alpha \sum_{j=\max(0, i-n/2)}^{\min(N-1, i+n/2)} (a_{x,y}^j)^2 \right)^\beta$$

Response-normalized activity

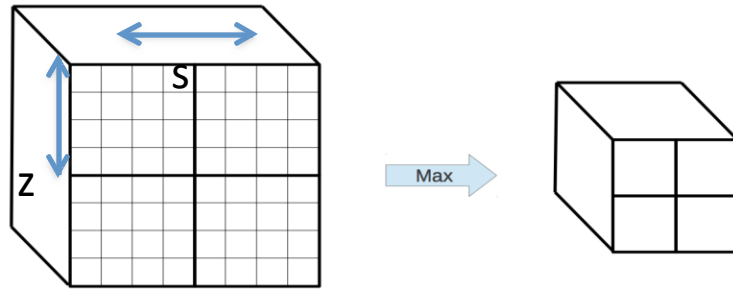
Activity of a neuron computed by applying kernel l at position (x,y) and then applying the ReLU nonlinearity

- Response normalization reduces top-1 and top-5 error rates by 1.4% and 1.2% , respectively.

Architectur

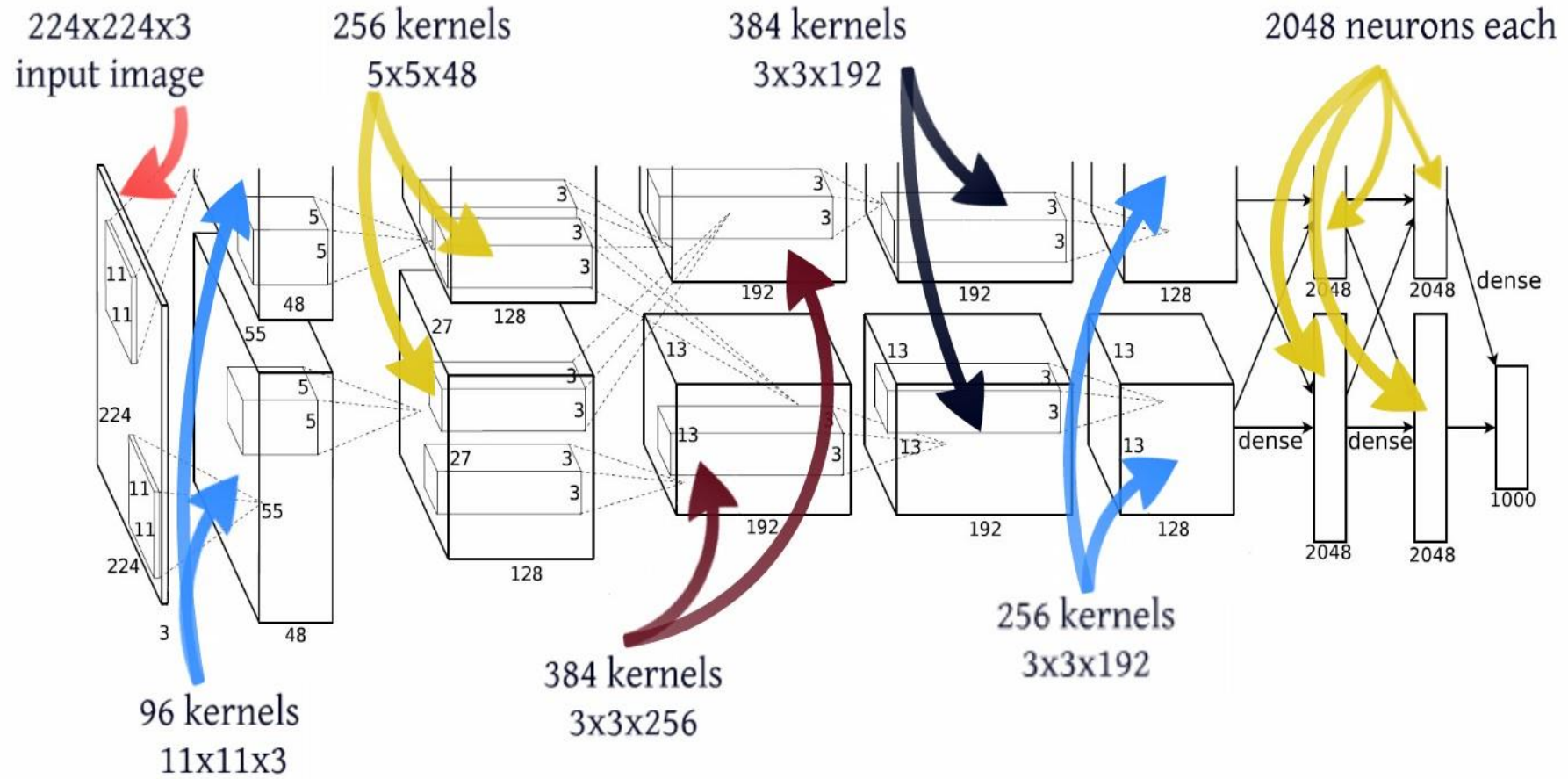
Overlapping Pooling

- Traditional pooling ($s = z$)

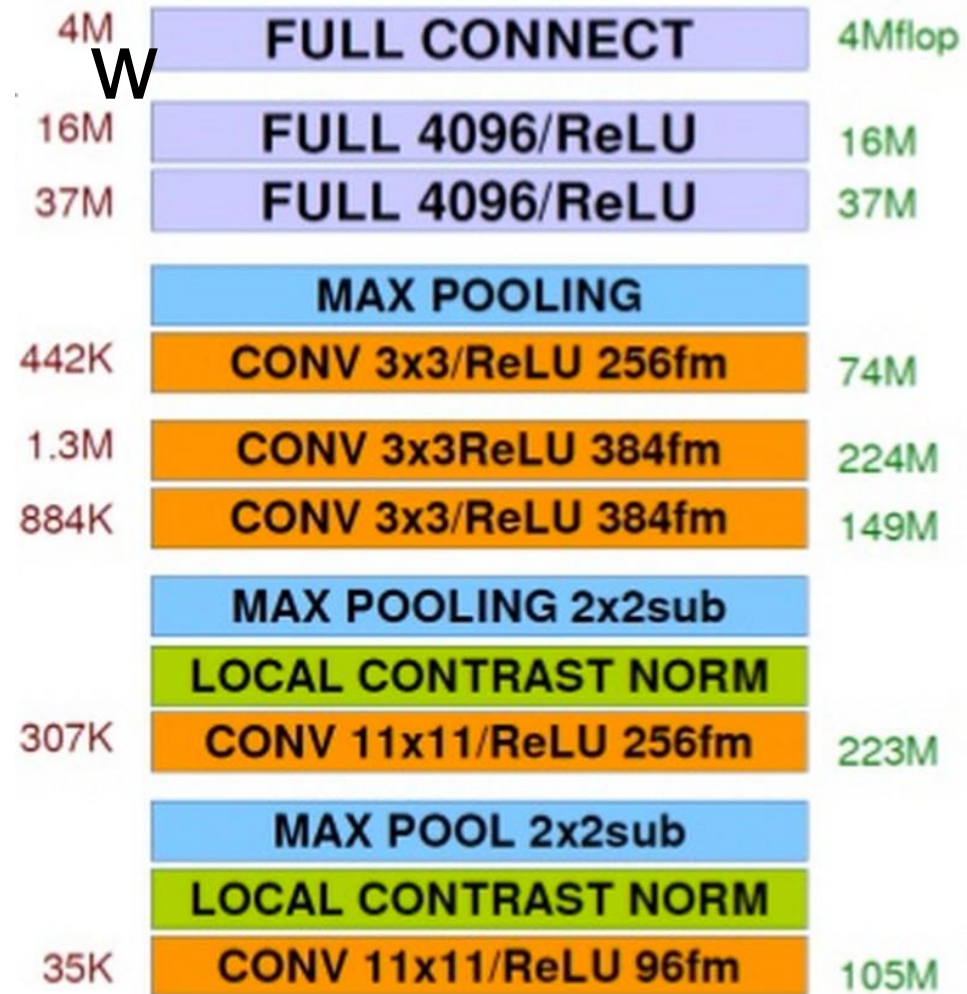


- $s < z \rightarrow$ overlapping pooling
- top-1 and top-5 error rates decrease by 0.4% and 0.3%, respectively, compared to the non-overlapping scheme $s = 2, z = 2$

Architecture



Architecture Overview



Reducing Overfitting

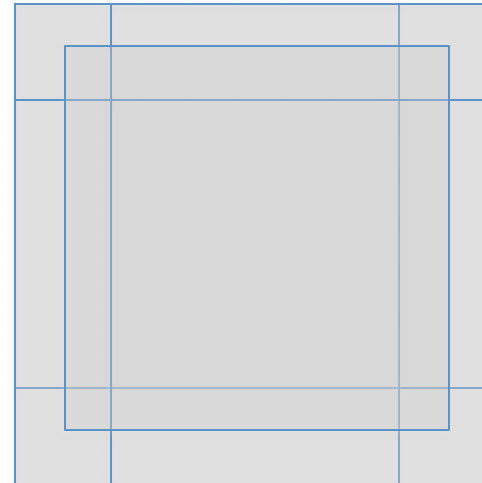
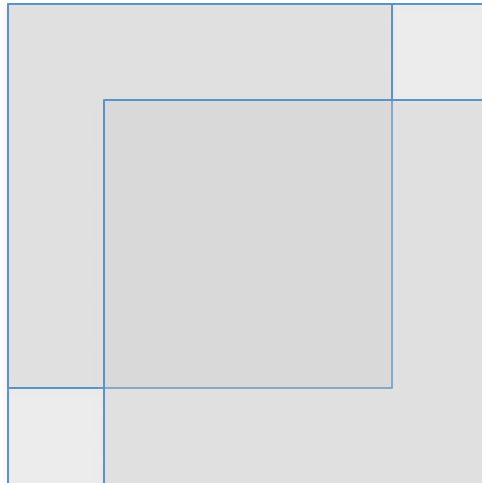
Data Augmentation

- 60 million parameters, 650,000 neurons
→ Overfits a lot.
- Crop 224x224 patches (and their horizontal reflections.)

Reducing Overfitting

Data Augmentation

- At test time, average the predictions on the 10 patches.



Reducing

- Softmax

$$L = \frac{1}{N} \sum_i -\log \left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right) + \lambda \sum_k \sum_l W_{k,l}^2$$

$j = 1 \dots 1000$

$P(y_i | x_i; W)$ Likelihood

- No need to calibrate to average the predictions over 10 patches.

cf. SVM

$$L = \frac{1}{N} \sum_i \sum_{j \neq y_i} \left[\max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + \Delta) + \lambda \sum_k \sum_l W_{k,l}^2 \right]$$

Reducing Overfitting

Data Augmentation

- Change the intensity of RGB channels
-

$$I_{xy} = [I_{xy}, I_{xy}^G, I_{xy}^{BR}]^T$$

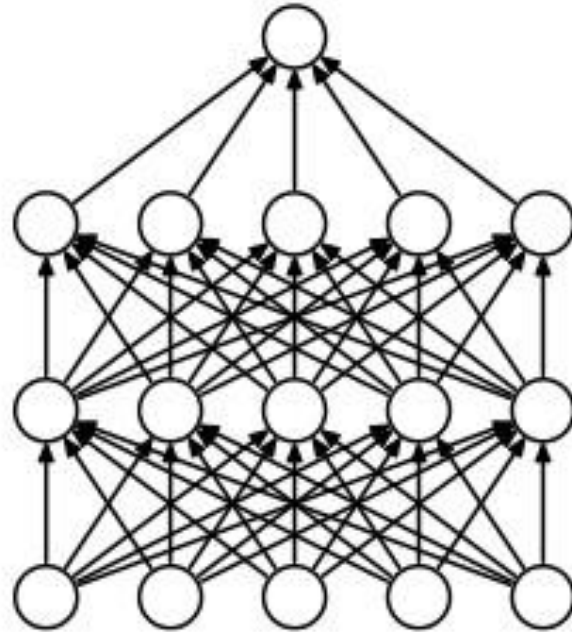
add multiples of principle
components

$$[\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3][\alpha_1 \lambda_1, \alpha_2 \lambda_2, \alpha_3 \lambda_3]^T$$

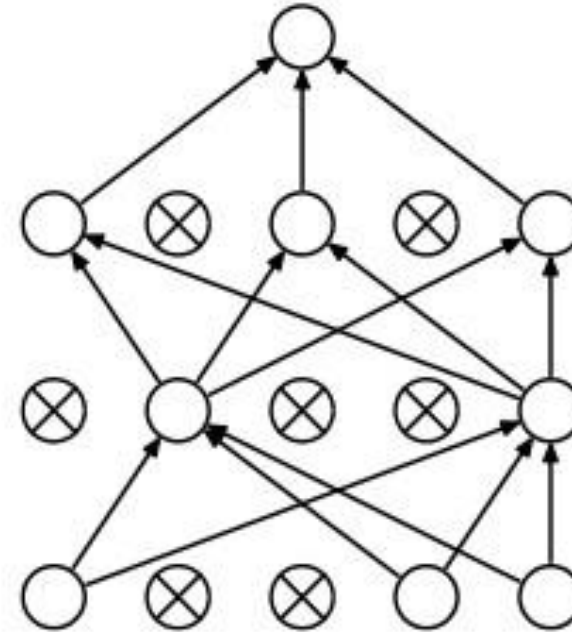
$$\langle_i \sim N(0, 0.1)$$

Reducing Overfitting

Dropout



Standard Neural Net



After applying dropout.

- With probability 0.5
- last two 4096 fully-connected layers.

Stochastic Gradient Descent Learning

Momentum Update

$$\begin{aligned} v_{i+1} &:= \underbrace{0.9}_{\text{momentum(damping parameter)}} \cdot v_i - \underbrace{0.0005}_{\text{weight decay}} \cdot \epsilon \cdot w_i - \underbrace{\epsilon}_{\text{Learning rate (initialized at 0.01)}} \cdot \underbrace{\left\langle \frac{\partial L}{\partial w} \Big|_{w_i} \right\rangle_{D_i}}_{\substack{\text{Gradient of Loss} \\ \text{w.r.t weight} \\ \text{Averaged over batch}}} \\ w_{i+1} &:= w_i + v_{i+1} \end{aligned}$$

Batch size: 128

- The training took **5 to 6 days** on **two NVIDIA GTX 580 3GB GPUs**.

Results: ILSVRC-2010

Model	Top-1	Top-5
<i>Sparse coding [2]</i>	47.1%	28.2%
<i>SIFT + FVs [24]</i>	45.7%	25.7%
CNN	37.5%	17.0%

Table 1: Comparison of results on ILSVRC-2010 test set. In *italics* are best results achieved by others.

Results: ILSVRC-2012

Model	Top-1 (val)	Top-5 (val)	Top-5 (test)
<i>SIFT + FVs [7]</i>	—	—	26.2%
1 CNN	40.7%	18.2%	—
5 CNNs	38.1%	16.4%	16.4%
1 CNN*	39.0%	16.6%	—
7 CNNs*	36.7%	15.4%	15.3%

Table 2: Comparison of error rates on ILSVRC-2012 validation and test sets. In *italics* are best results achieved by others. Models with an asterisk* were “pre-trained” to classify the entire ImageNet 2011 Fall release. See Section 6 for details.

96 Convolutional



- 11 x 11 x 3 size kernels.
- top 48 kernels on GPU 1 : color-agnostic
- bottom 48 kernels on GPU 2 : color-specific.

Why?

Eight ILSVRC-2010 test images



Five ILSVRC-2010 test



The output from the last 4096 fully-connected layer
: 4096 dimensional feature.

Discussion

- Depth is really important.

removing a single convolutional layer degrades the performance.

K. Simonyan, A. Zisserman.

[Very Deep Convolutional Networks for Large-Scale Image Recognition](#). Technical report, 2014.

→ 16-layer model, 19-layer model. 7.3% top-5 test error on ILSVRC-2012