

Deep Learning with Tensorflow

http://cvml.ist.ac.at/courses/DLWT_W17/

AlexNet

Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton,
"Imagenet classification with deep convolutional neural networks",
Advances in neural information processing systems, 2012

Djordje Slijepcevic

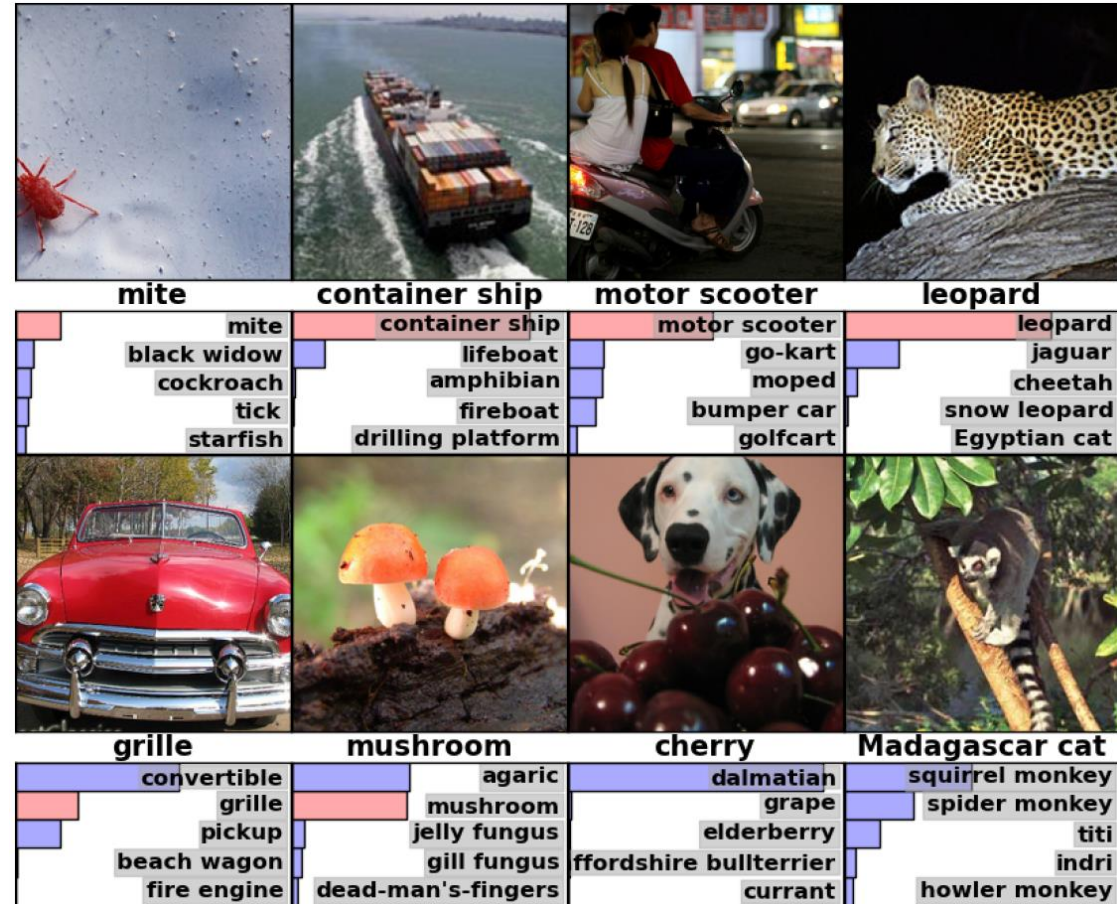
Introduction

- Convolutional Neural Network (CNN)
- Winner of ImageNet Large Scale Visual Recognition Challenge (ILSVRC) 2012
 - first successful CNN application for such a big dataset
 - top-5 test error rate of 15.3% (+10.9% compared to 2nd)
- Relatively simple layout (compared to modern architectures)
 - 5 conv. layers
 - 3 fully connected layers
 - max-pooling layers
 - dropout layers

Dataset

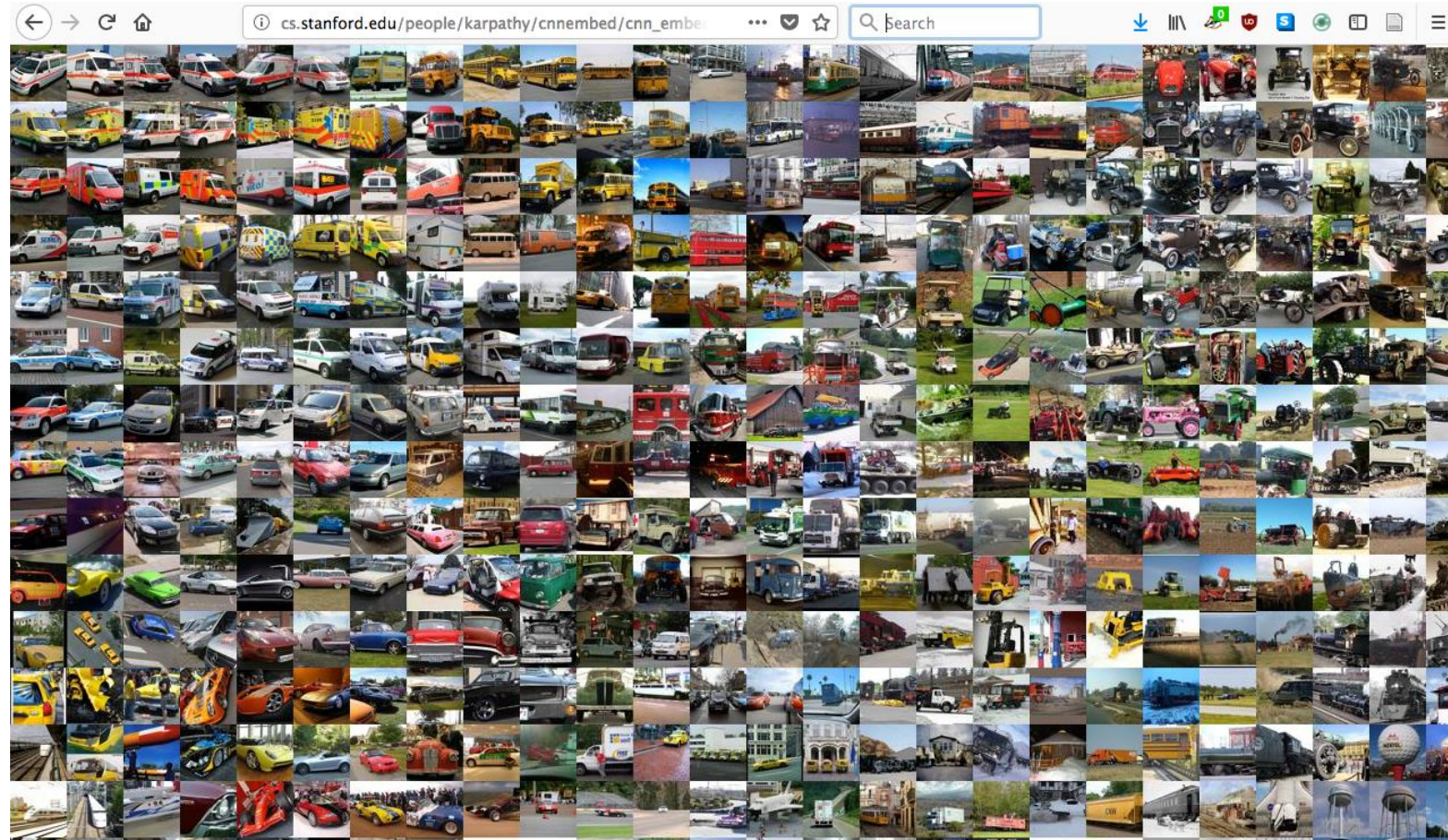
- ImageNet:
 - 15+ million labeled high-resolution images
 - 22000 categories
- ILSVRC uses a subset of ImageNet:
 - ~ 1000 images per category
 - 1000 categories
 - 1.2 million training images | 50000 validation images | 150000 testing images
- AlexNet:
 - images were down-sampled and cropped to 256×256 pixels
 - subtraction of the mean activity over the training set from each pixel

Task

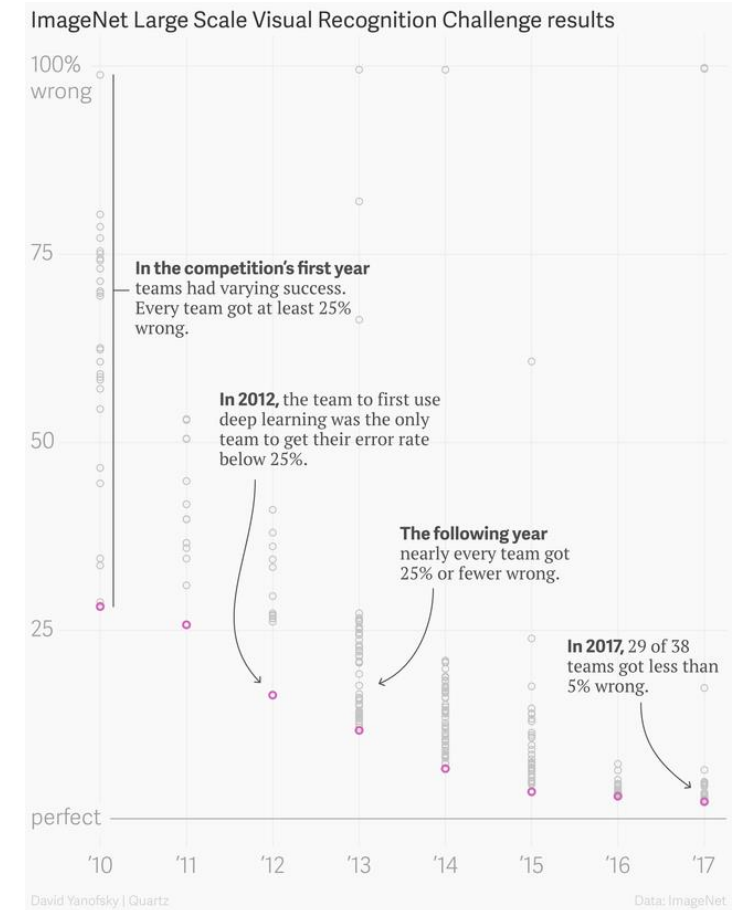


[A. Krizhevsky, I. Sutskever, G.E. Hinton, *ImageNet Classification with Deep Convolutional Neural Networks*, 2012]

Dataset

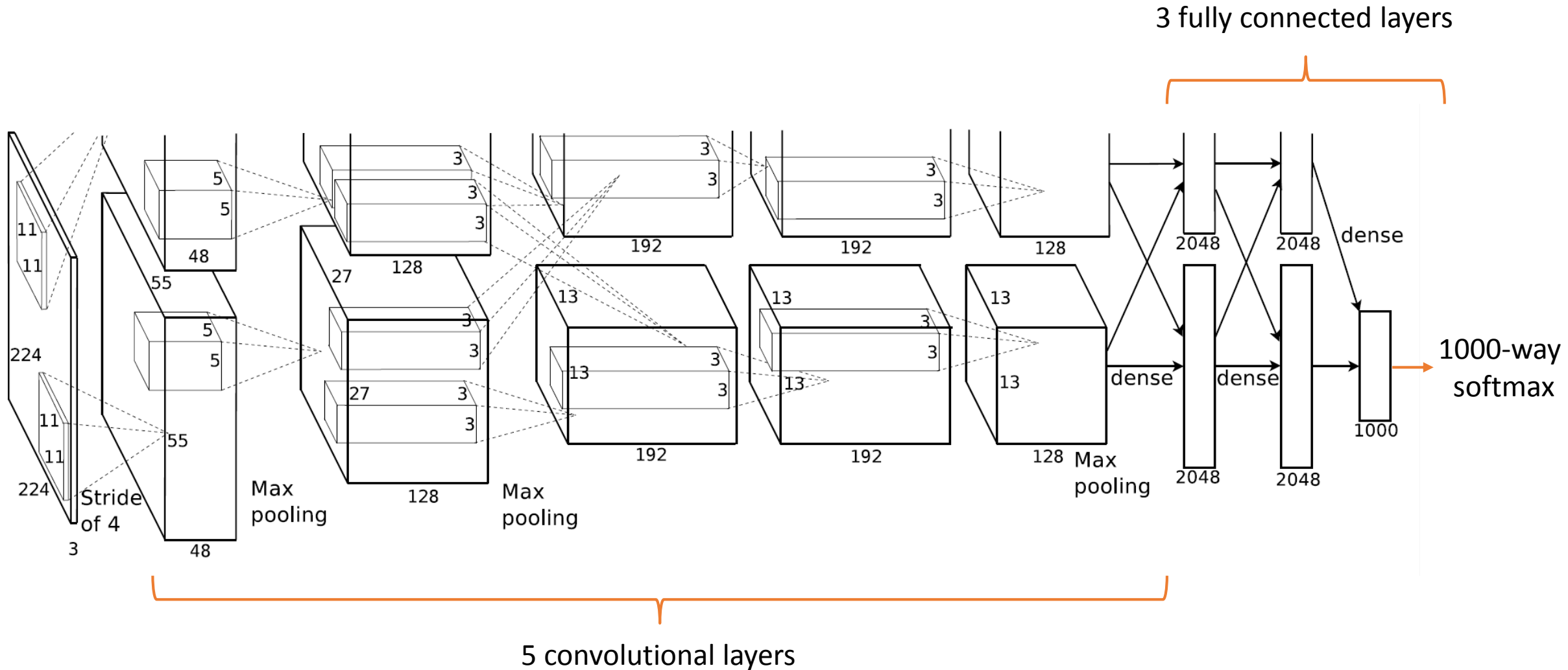


[<http://cs.stanford.edu/people/karpathy/cnnembed/>, 30.11.2017]



[<https://qz.com/1034972/the-data-that-changed-the-direction-of-ai-research-and-possibly-the-world/>, 30.11.2017]

Architecture



Activation function

- Traditionally, saturating nonlinearities:

- hyperbolic tangent function: $f(x) = \tanh(x) = 2 * \frac{1}{1+e^{-2x}} - 1$

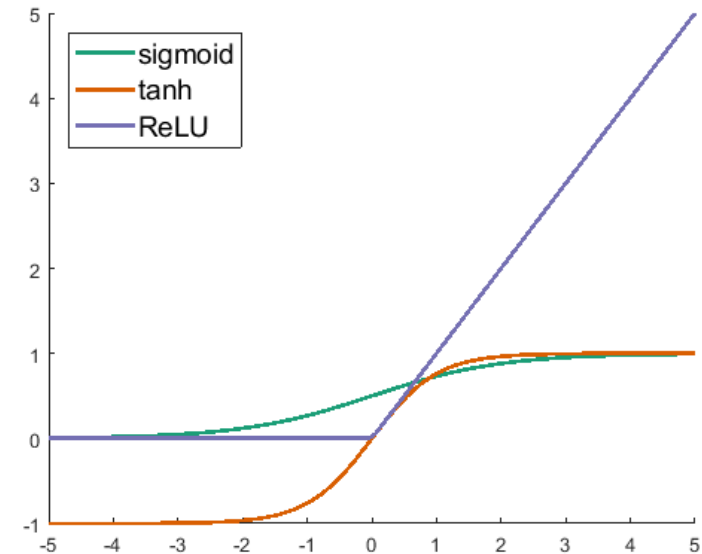
- sigmoid function: $f(x) = \frac{1}{1+e^{-x}}$

→ slow to train

- Non-saturating nonlinearity:

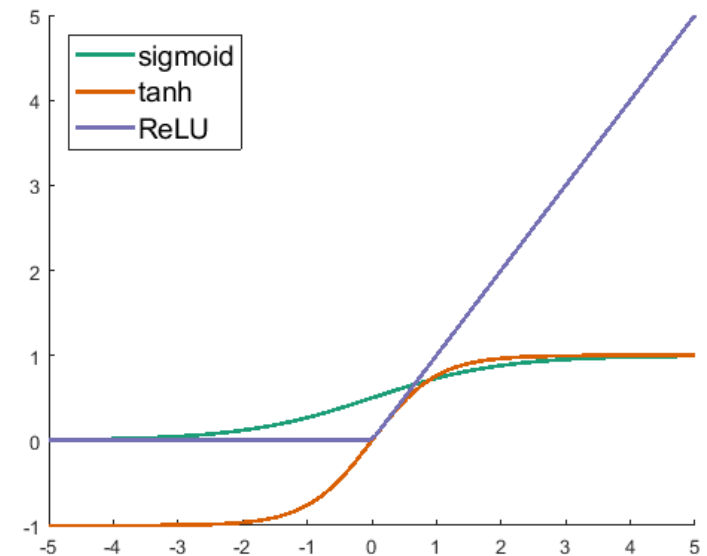
- Rectified Linear Unit (ReLU): $f(x) = \max(0, x)$

→ quick to train

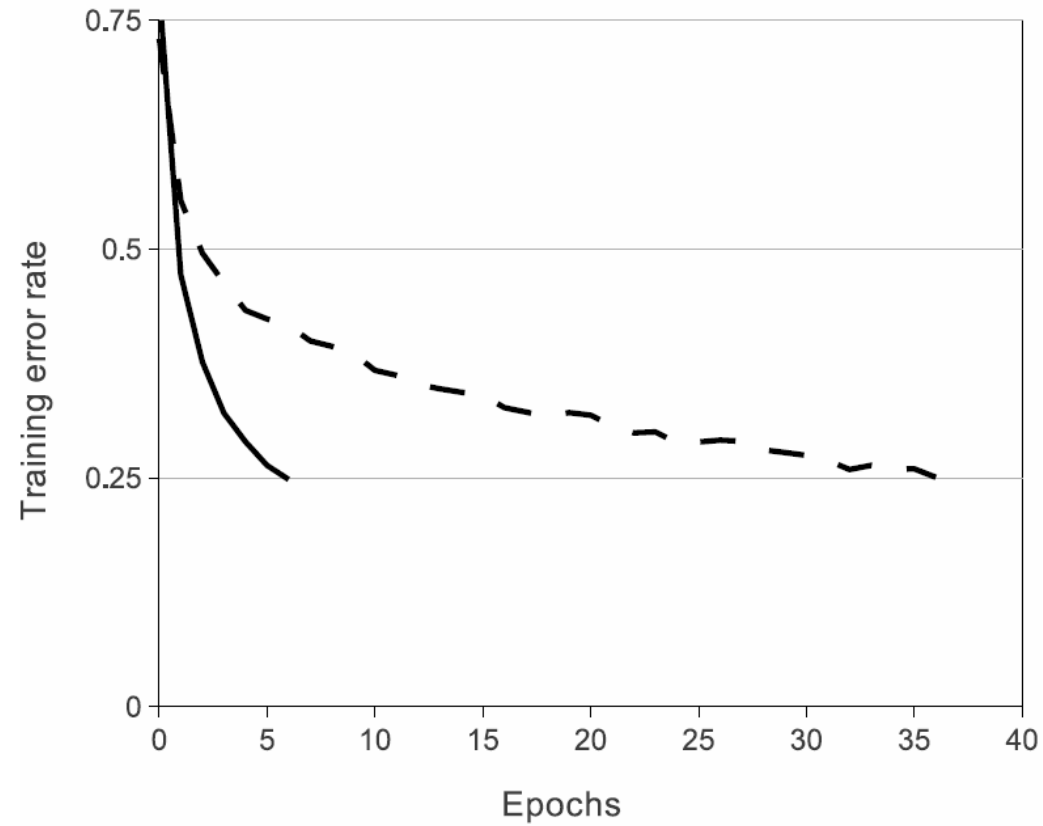


Activation function

- Traditionally, saturating nonlinearities:
 - Saturated neurons facilitate vanishing of gradients
 - exp function is a bit compute expensive
 - slow to train
- Non-saturating nonlinearity:
 - Does not saturate (in the + region)
 - Very computationally efficient
 - quick to train

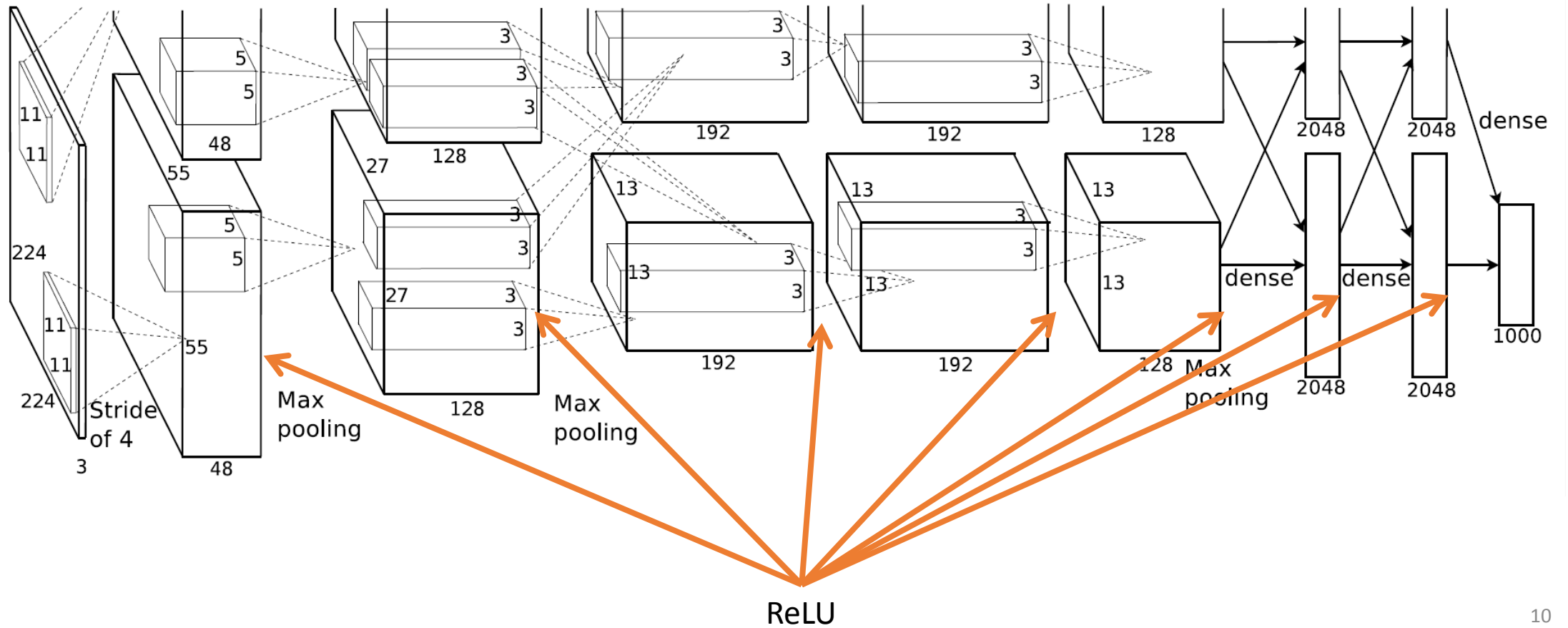


Activation function



- Dataset: CIFAR-10
 - Experiment:
 - CNN (4 layers) + ReLUs (solid line)
 - vs.
 - CNN (4 layers) + tanh (dashed line)
- ReLUs **six times faster**

Activation function

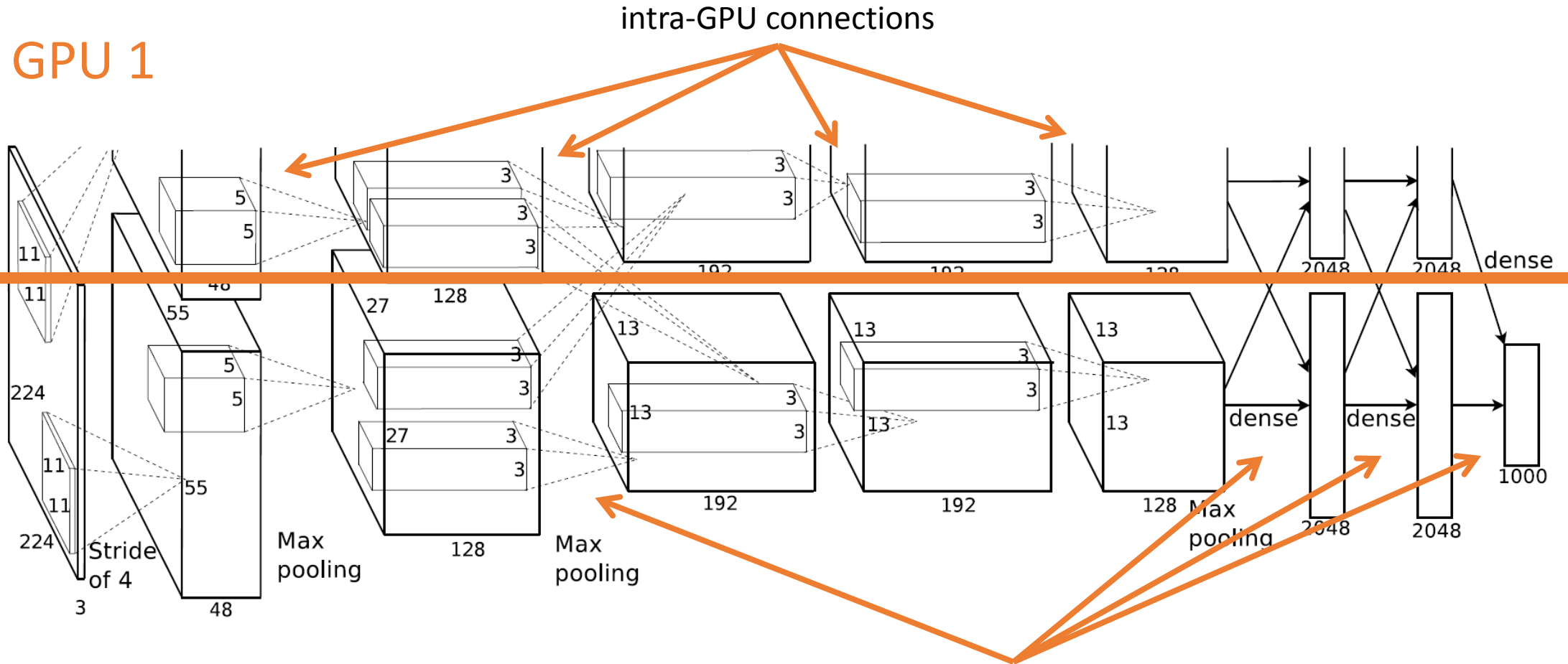


Training on Multiple GPUs

- Half of the neurons of an certain layer are on each GPU
- GPUs communicate only in certain layers
- Improvement (as compared with a net with half as many kernels in each convolutional layer trained on one GPU):
 - top-1 error rate by **1.7%**
 - top-5 error rate by **1.2%**

Training on Multiple GPUs

GPU 1



GPU 2

Inter-GPU connections

Local Response Normalization

- ReLUs do not require input normalization to prevent them from saturating
- However, Local Response Normalization aids generalization

Activity of a neuron by applying kernel i at position (x,y)

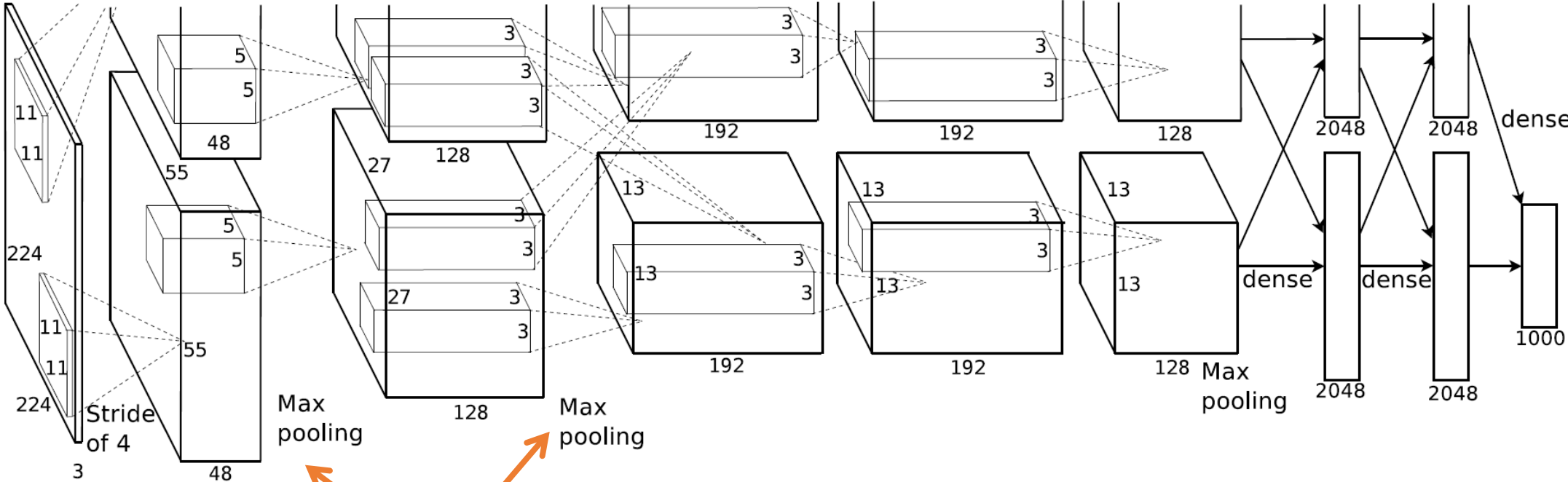
$$b_{x,y}^i = a_{x,y}^i / \left(k + \alpha \sum_{j=\max(0, i-\frac{n}{2})}^{\min(N-1, i+\frac{n}{2})} (a_{x,y}^j)^2 \right)^\beta$$

$$\begin{aligned} k &= 2 \\ n &= 5 \\ \alpha &= 10^{-4} \\ \beta &= 0.75 \end{aligned}$$

- Improvement:
 - top-1 error rate by **1.4%**
 - top-5 error rate by **1.2%**

sum runs over n "adjacent" kernel maps at the same spatial position

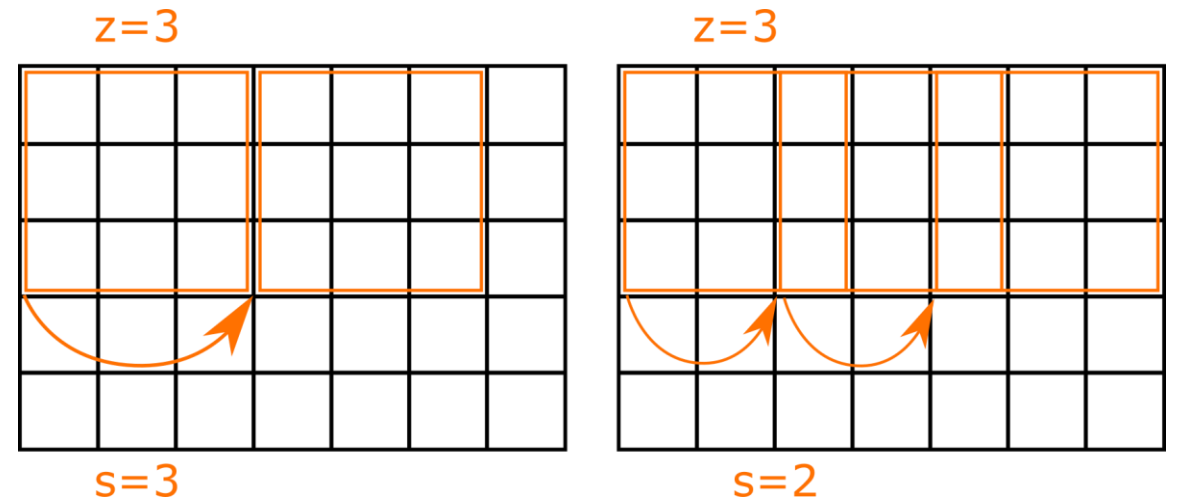
Local Response Normalization



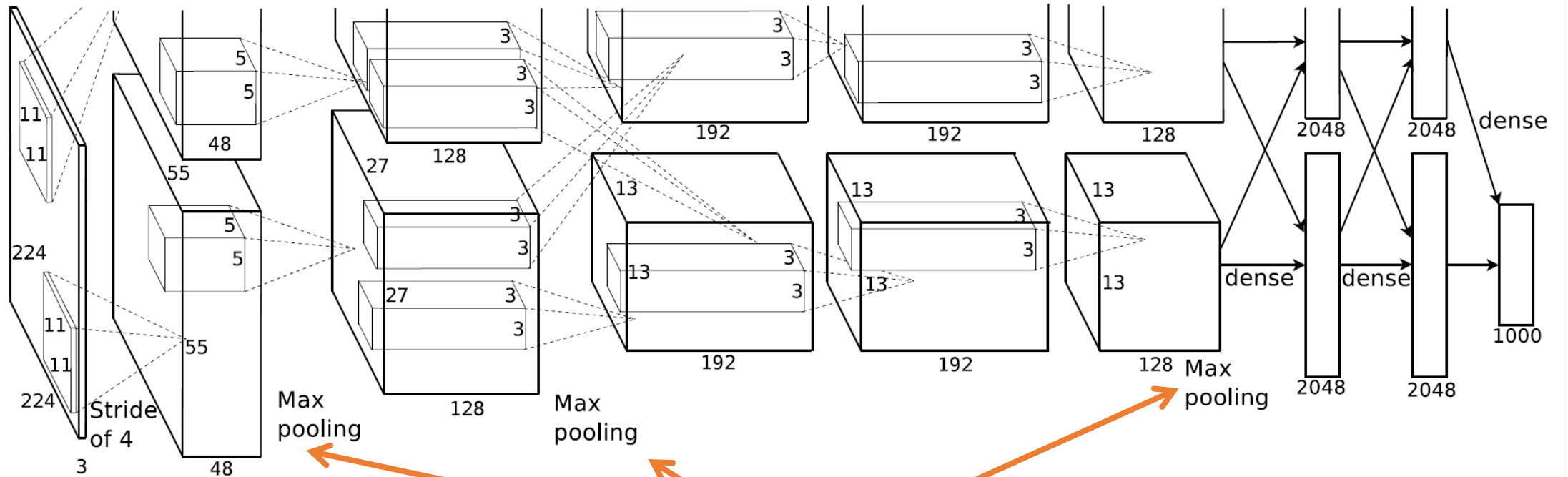
Local Response Normalization

Overlapping Pooling

- Pooling layers summarize the outputs of neighboring neurons in the same kernel map.
- Overlapping pooling $\rightarrow s < z$
- Improvement using MaxPooling:
 - top-1 error rate by **0.4%**
 - top-5 error rates by **0.3%**

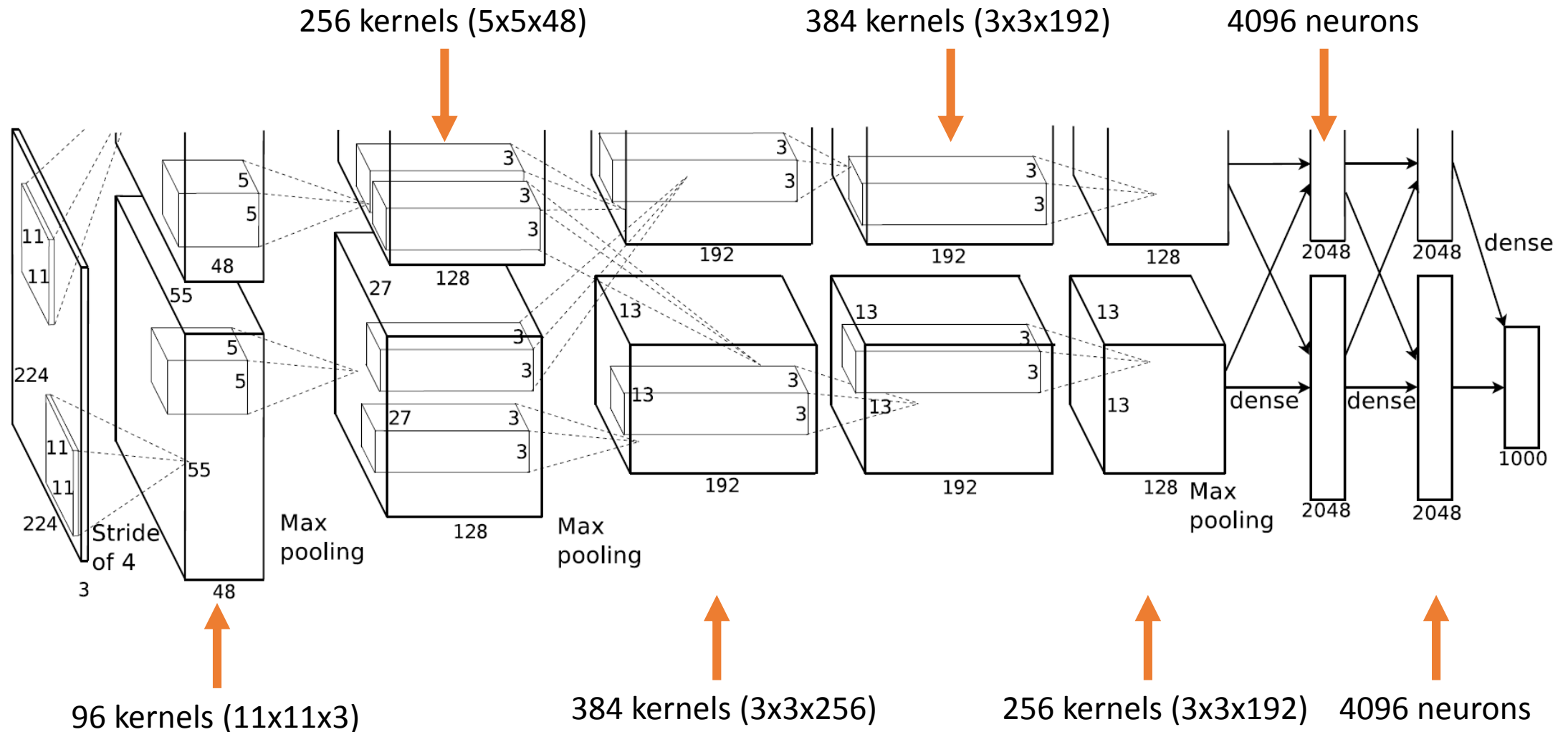


Overlapping Pooling



Overlapping Pooling

Overall Architecture



Reducing Overfitting - Data Augmentation

- 1st : image translations and horizontal reflections
 - random 224x224 patches + horizontal reflections from the 256x256 images
 - Testing: five 224x224 patches + horizontal reflections → averaging the predictions over the ten patches
- 2nd : change the intensity of RGB channels
 - PCA on the set of RGB pixel values throughout the ImageNet training set
 - To each RGB image pixel $I_{xy} = [I_{xy}^R, I_{xy}^G, I_{xy}^B]$ following is added

$$[p_1, p_2, p_3][\alpha_1 \lambda_1, \alpha_2 \lambda_2, \alpha_3 \lambda_3]^T \quad |\alpha_i \sim N(0, 0.1)$$

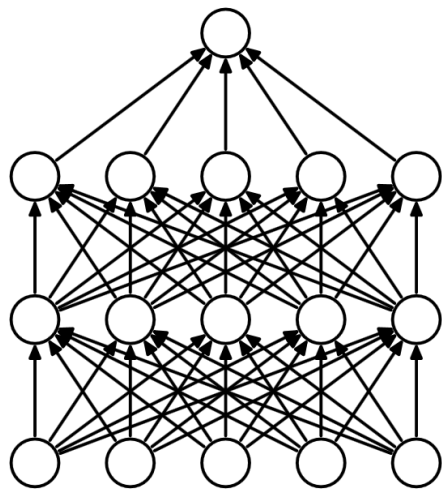
- Improvement:

- top-1 error rate by **1%**

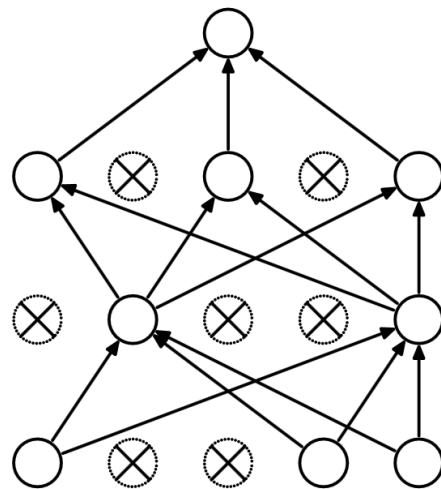
eigenvectors

eigenvalues

Reducing Overfitting - Dropout



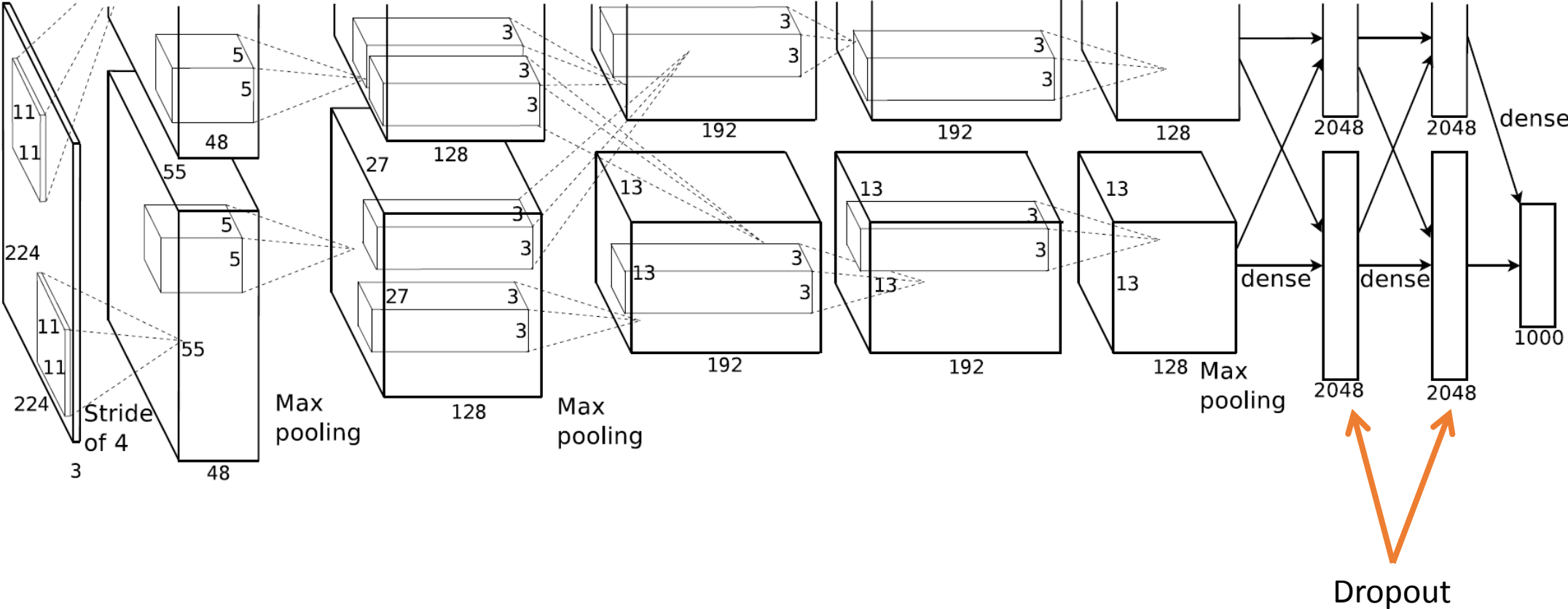
(a) Standard Neural Net



(b) After applying dropout.

- Output of each hidden neuron is set to zero with probability 0.5
- Learning more robust features
- Doubles the number of iterations required to converge
- Applied in the first two fully connected layers

Reducing Overfitting - Dropout



Stochastic Gradient Descent

- Training process
 - Minimizing the cross-entropy loss function:

$$L(w) = \sum_{i=1}^N \sum_{c=1}^{1000} -y_{ic} \log f_c(x_i) + \epsilon \|w\|_2^2$$

indicator that example i has label c

predicted probability of class c for image x

Stochastic Gradient Descent

- SGD with a batch size of 128
- Learning rate initialized at 0.01; divided by 10 if validation error rate stopped improving
- Update rule for weight w :

$$v_{i+1} := 0.9 * v_i - 0.0005 * \epsilon * w_i - \epsilon * \left\langle \frac{\partial L}{\partial w} \Big|_{w_i} \right\rangle_{D_i}$$

$w_{i+1} := w_i + v_{i+1}$

The diagram illustrates the update rule for the velocity v_{i+1} and the weight w_{i+1} . The velocity update rule is $v_{i+1} := 0.9 * v_i - 0.0005 * \epsilon * w_i - \epsilon * \left\langle \frac{\partial L}{\partial w} \Big|_{w_i} \right\rangle_{D_i}$. The weight update rule is $w_{i+1} := w_i + v_{i+1}$. Annotations with orange arrows point to specific terms: 'momentum' points to the coefficient 0.9; 'weight decay' points to the coefficient 0.0005; 'learning rate' points to the coefficient ϵ ; and 'Gradient of Loss' points to the term $\left\langle \frac{\partial L}{\partial w} \Big|_{w_i} \right\rangle_{D_i}$.

- ~ 90 cycles \rightarrow five to six days on two NVIDIA GTX 580 3GB GPUs