# COMP9444
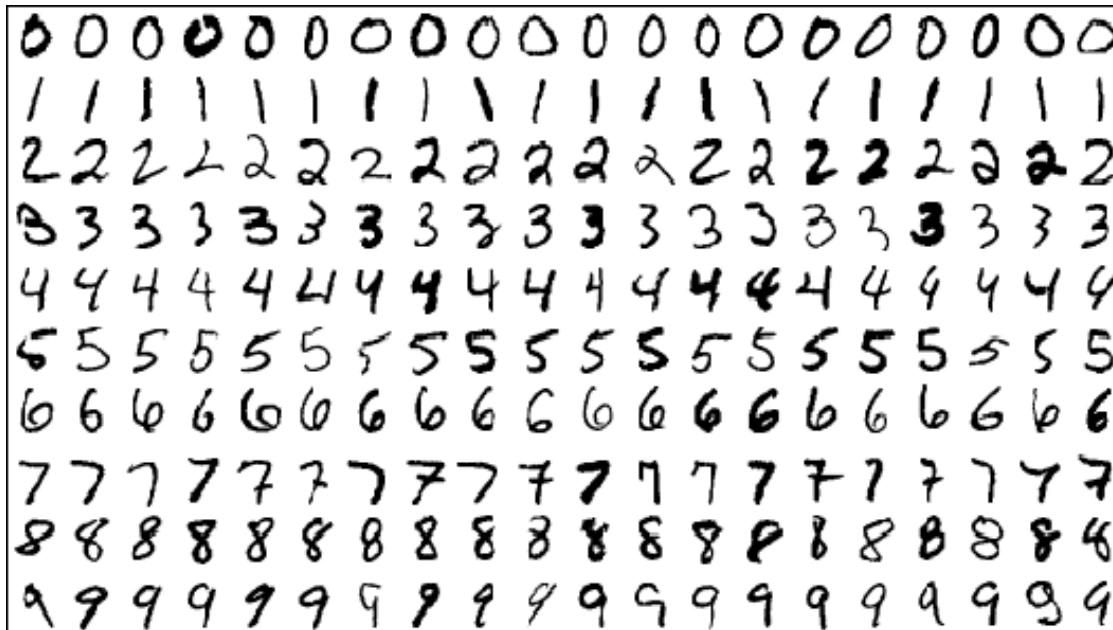# Neural Networks and Deep Learning

# 7. Image Processing

# Outline

- Image Datasets and Tasks

- Convolution in Detail

- AlexNet

- Weight Initialization

- Batch Normalization

- Residual Networks

- Dense Networks

- Style Transfer

# MNIST Handwritten Digit Dataset



- ▪ black and white, resolution $28 \times 28$
- ▪ 60,000 images
- ▪ 10 classes $(0, 1, 2, 3, 4, 5, 6, 7, 8, 9)$

# CIFAR Image Dataset



- color, resolution $32 \times 32$
- 50,000 images
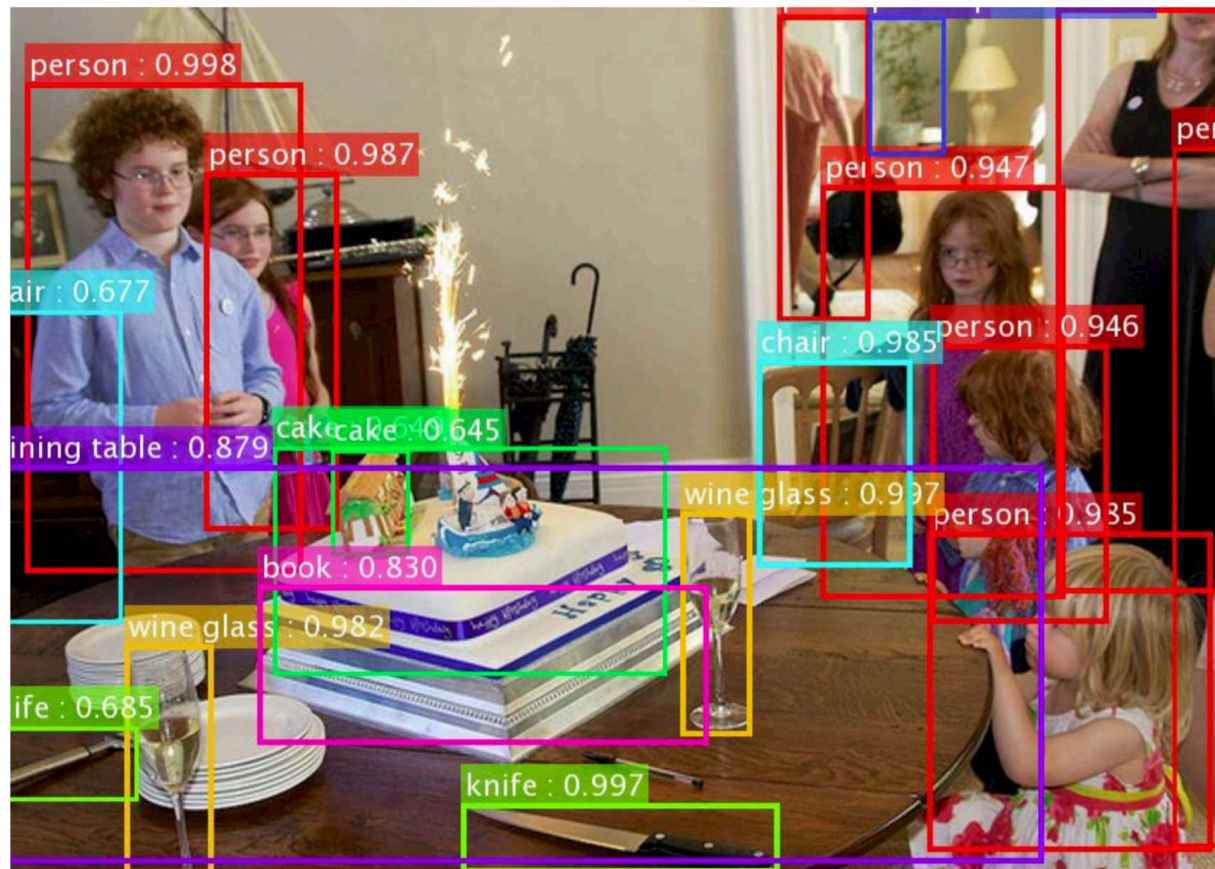- 10 classes

# ImageNet LSVRC Dataset



■ color, resolution $227 \times 227$

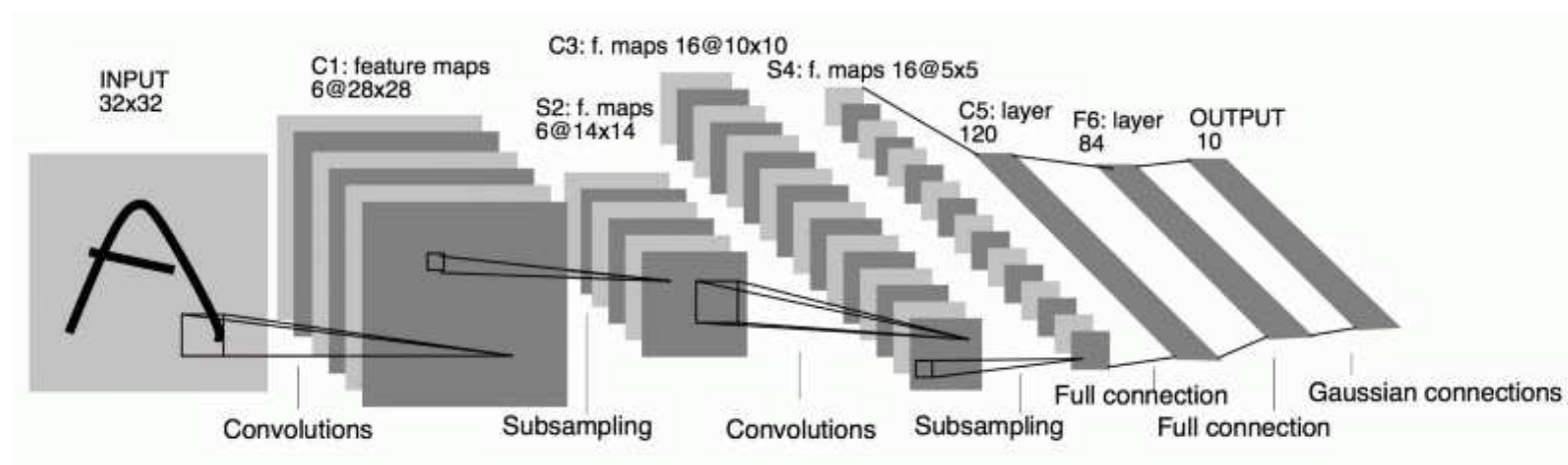■ 1.2 million images

■ 1000 classes

# Image Processing Tasks

- image classification

- object detection

- object segmentation

- style transfer

- generating images

- generating art

- image captioning

# Object Detection

# LeNet trained on MNIST



The $5 \times 5$ window of the first convolution layer extracts from the original $32 \times 32$ image a $28 \times 28$ array of features. Subsampling then halves this size to $14 \times 14$. The second Convolution layer uses another $5 \times 5$ window to extract a $10 \times 10$ array of features, which the second subsampling layer reduces to $5 \times 5$. These activations then pass through two fully connected layers into the 10 output units corresponding to the digits '0' to '9'.
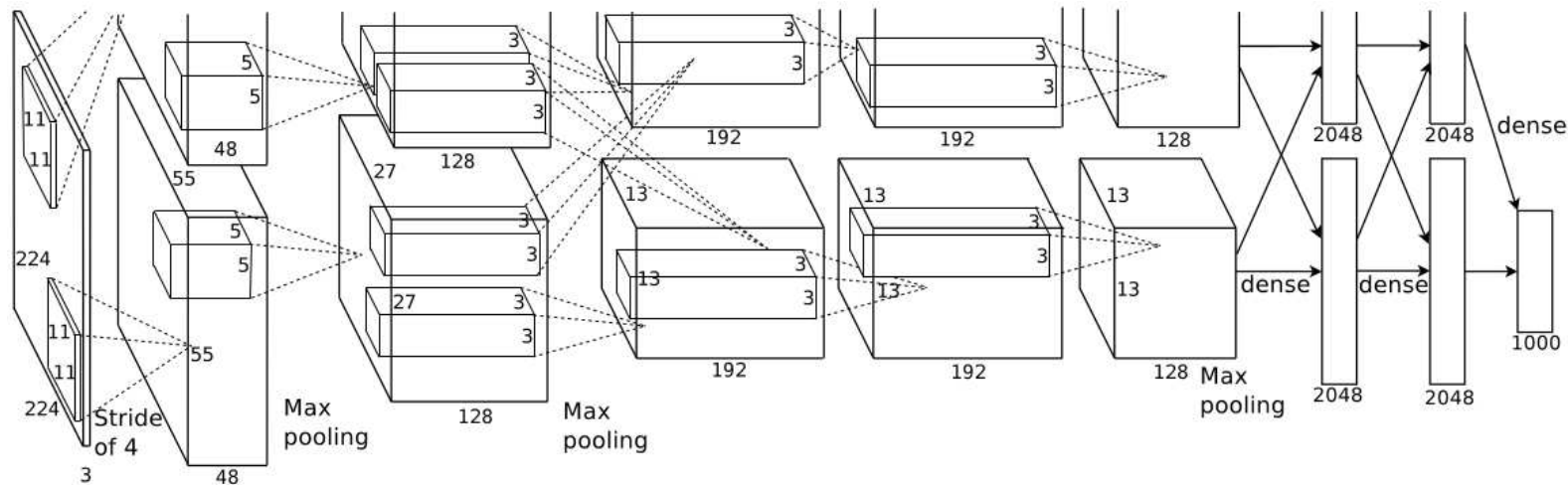
# ImageNet Architectures

■ AlexNet, 8 layers (2012)

■ VGG, 19 layers (2014)

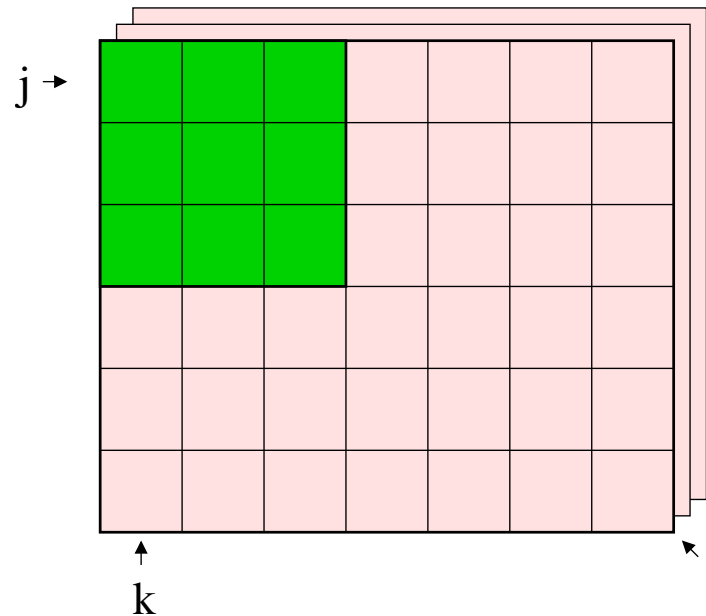■ GoogleNet, 22 layers (2014)

■ ResNets, 152 layers (2015)

# AlexNet Architecture



- 5 convolutional layers + 3 fully connected layers
- max pooling with overlapping stride
- softmax with 1000 classes
- 2 parallel GPUs which interact only at certain layers

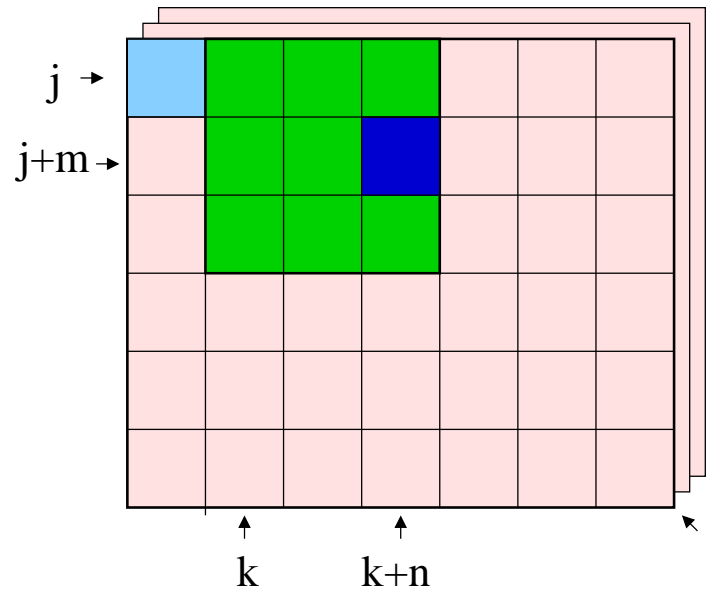# Convolutional Neural Networks



Assume the original image is $J \times K$, with $L$ channels.

We apply an $M \times N$ "filter" to these inputs to compute one hidden unit in the convolution layer. In this example $J = 6, K = 7, L = 3, M = 3, N = 3$.

# Convolutional Neural Networks



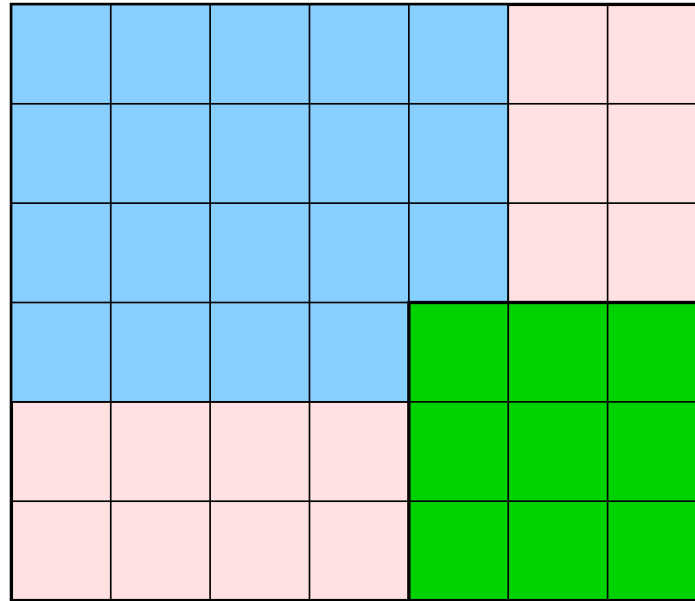$$Z_{j,k}^i = g\left(b^i + \sum_l \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} K_{l,m,n}^i V_{j+m,k+n}^l\right)$$

The same weights are applied to the next $M \times N$ block of inputs, to compute the next hidden unit in the convolution layer ("weight sharing").
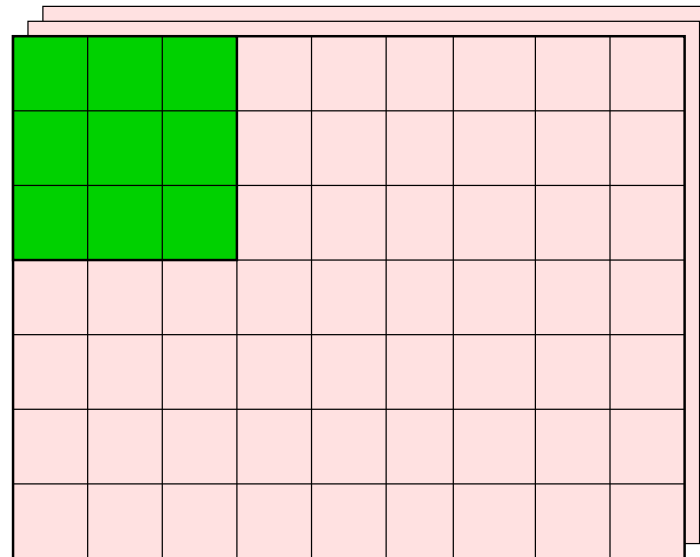
# Convolutional Neural Networks



If the original image size is $J \times K$ and the filter is size $M \times N$, the convolution layer will be $(J + 1 - M) \times (K + 1 - N)$

# Stride



Assume the original image is $J \times K$, with $L$ channels.

We again apply an $M \times N$ filter, but this time with a "stride" of $s > 1$.

In this example $J = 7, K = 9, L = 3, M = 3, N = 3, s = 2$.

# Stride



$$Z^i_{j,k} = g\Big(b^i + \sum_l \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} K^i_{l,m,n} V^l_{j+m,k+n}\Big)$$

The same formula is used, but $j$ and $k$ are now incremented by $s$ each time.

The number of free parameters is $1 + L \times M \times N$

# Stride Dimensions



$j$ takes on the values $0, s, 2s, \ldots, (J-M)$

$k$ takes on the values $0, s, 2s, \ldots, (K-N)$

The next layer is $(1 + (J-M)/s)$ by $(1 + (K-N)/s)$

# Stride with Zero Padding



When combined with zero padding of width $P$,

$j$ takes on the values $0, s, 2s, \ldots, (J + 2P - M)$

$k$ takes on the values $0, s, 2s, \ldots, (K + 2P - N)$

The next layer is $(1 + (J + 2P - M)/s)$ by $(1 + (K + 2P - N)/s)$

# Example: AlexNet Conv Layer 1

For example, in the first convolutional layer of AlexNet,
$J = K = 224$, $P = 2$, $M = N = 11$, $s = 4$.

The width of the next layer is

$$1 + (J + 2P - M)/s = 1 + (224 + 2 \times 2 - 11)/4 = 55$$

Question: If there are 96 filters in this layer, compute the number of:

weights per neuron?

neurons?

connections?

independent parameters?

# Example: AlexNet Conv Layer 1

For example, in the first convolutional layer of AlexNet,
$J = K = 224$, $P = 2$, $M = N = 11$, $s = 4$.

The width of the next layer is

$$1 + (J - M)/s = 1 + (224 + 2 \times 2 - 11)/4 = 55$$

Question: If there are 96 filters in this layer, compute the number of:

| | | | |
|---|---|---|---|
| weights per neuron? | $1 + 11 \times 11 \times 3$ | $=$ | $364$ |
| neurons? | $55 \times 55 \times 96$ | $=$ | $290,400$ |
| connections? | $55 \times 55 \times 96 \times 364$ | $=$ | $105,705,600$ |
| independent parameters? | $96 \times 364$ | $=$ | $34,944$ |

# Max Pooling



Single depth slice

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

max pool with 2x2 filters
and stride 2

| 6 | 8 |
|---|---|
| 3 | 4 |

y

# Overlapping Pooling

If the previous layer is $J \times K$, and max pooling is applied with width $F$ and stride $s$, the size of the next layer will be

$$(1 + (J - F)/s) \times (1 + (K - F)/s)$$

Question: If max pooling with width 3 and stride 2 is applied to the features of size $55 \times 55$ in the first convolutional layer of AlexNet, what is the size of the next layer?

Answer: $1 + (55 - 3)/2 = 27$.

Question: How many independent parameters does this add to the model?

Answer: None! (no weights to be learned, just computing max)

# AlexNet Details



- 650K neurons

- 630M connections

- 60M parameters

- more parameters that images → danger of overfitting

# Enhancements

- Rectified Linear Units (ReLUs)

- overlapping pooling (width $= 3$, stride $= 2$)

- stochastic gradient descent with momentum and weight decay

- data augmentation to reduce overfitting

- 50% dropout in the fully connected layers

# Data Augmentation

- ten patches of size $224 \times 224$ are cropped from each of the original $227 \times 227$ images (using zero padding)

- the horizontal reflection of each patch is also included.

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

- also include changes in intensity to RGB channels

# Convolution Kernels



- filters on GPU-1 (upper) are color agnostic

- filters on GPU-2 (lower) are color specific

- these resemble Gabor filters

# Dealing with Deep Networks

- \> 10 layers

  ▶ weight initialization

  ▶ batch nomalization

- \> 30 layers

  ▶ skip connections

- \> 100 layers

  ▶ identity skip connections

# Statistics

The mean and variance of a set of $n$ samples $x_1, \ldots, x_n$ are given by

$$\text{Mean}[x] = \frac{1}{n} \sum_{k=1}^{n} x_k$$

$$\text{Var}[x] = \frac{1}{n} \sum_{k=1}^{n} (x_k - \text{Mean}[x])^2 = \left( \frac{1}{n} \sum_{k=1}^{n} x_k^2 \right) - \text{Mean}[x]^2$$

If $w_k, x_k$ are independent and $y = \sum_{k=1}^{n} w_k x_k$ then

$$\text{Var}[y] = n \, \text{Var}[w] \, \text{Var}[x]$$

# Weight Initialization

Consider one layer $(i)$ of a deep neural network with weights $w_{jk}^{(i)}$ connecting the activations $\{x_k^{(i)}\}_{1 \le k \le n_i}$ at the previous layer to $\{x_j^{(i+1)}\}_{1 \le j \le n_{i+1}}$ at the next layer, where $g()$ is the transfer function and

$$x_j^{(i+1)} = g(\text{sum}_j^{(i)}) = g\left(\sum_{k=1}^{n_i} w_{jk}^{(i)} x_k^{(i)}\right)$$

Then

$$\text{Var}[\text{sum}^{(i)}] = n_i \text{Var}[w^{(i)}] \text{Var}[x^{(i)}]$$

$$\text{Var}[x^{(i+1)}] \simeq G_0 \, n_i \text{Var}[w^{(i)}] \text{Var}[x^{(i)}]$$

Where $G_0$ is a constant whose value is estimated to take account of the transfer function.

If some layers are not fully connected, we replace $n_i$ with the average number $n_i^{\text{in}}$ of incoming connections to each node at layer $i+1$.

# Weight Initialization

If the nework has $D$ layers, with input $x = x^{(1)}$ and output $z = x^{(D+1)}$, then

$$\text{Var}[z] \simeq \left( \prod_{i=1}^{D} G_0 \, n_i^{\text{in}} \, \text{Var}[w^{(i)}] \right) \text{Var}[x]$$

When we apply gradient descent through backpropagation, the differentials will follow a similar pattern:

$$\text{Var}[\frac{\partial}{\partial x}] \simeq \left( \prod_{i=1}^{D} G_1 \, n_i^{\text{out}} \, \text{Var}[w^{(i)}] \right) \text{Var}[\frac{\partial}{\partial z}]$$

In this equation, $n_i^{\text{out}}$ is the average number of outgoing connections for each node at layer $i$, and $G_1$ is meant to estimate the average value of the derivative of the transfer function.

For Rectified Linear Units, we can assume $G_0 = G_1 = \frac{1}{2}$
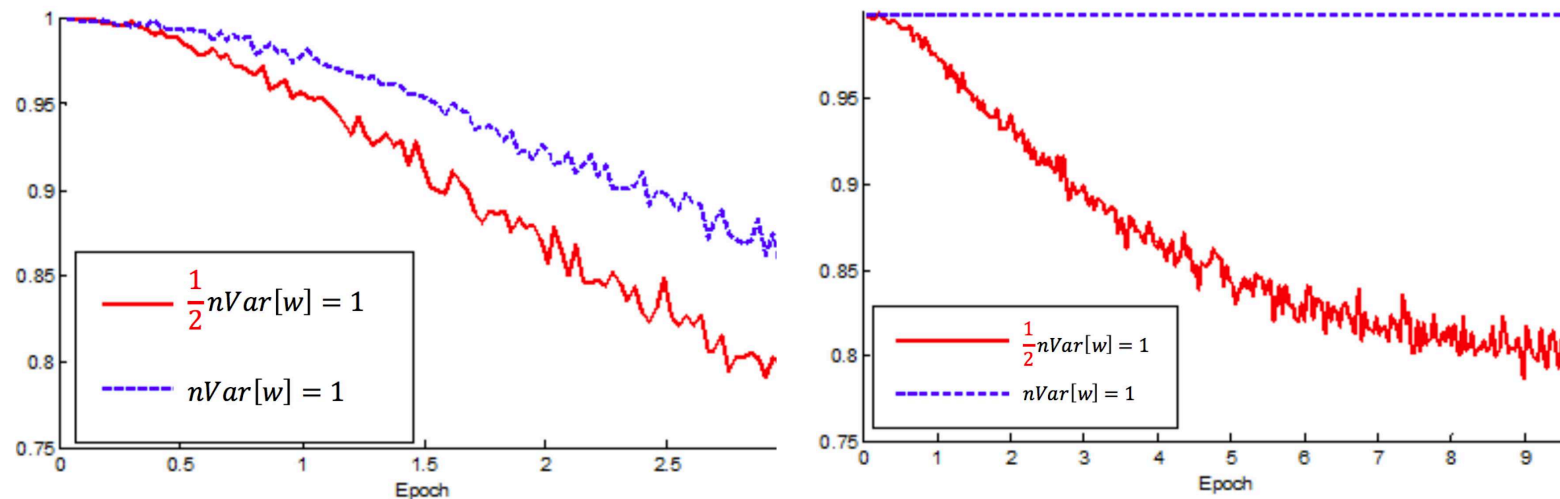
# Weight Initialization

In order to have healthy forward and backward propagation, each term in the product must be approximately equal to 1. Any deviation from this could cause the activations to either vanish or saturate, and the differentials to either decay or explode exponentially.

$$\text{Var}[z] \simeq \left( \prod_{i=1}^{D} G_0 \, n_i^{\text{in}} \text{Var}[w^{(i)}] \right) \text{Var}[x]$$

$$\text{Var}[\frac{\partial}{\partial x}] \simeq \left( \prod_{i=1}^{D} G_1 \, n_i^{\text{out}} \text{Var}[w^{(i)}] \right) \text{Var}[\frac{\partial}{\partial z}]$$

We therefore choose the initial weights $\{w_{jk}^{(i)}\}$ in each layer $(i)$ such that

$$G_1 n_i^{\text{out}} \text{Var}[w^{(i)}] = 1$$

# Weight Initialization



- 22-layer ReLU network (left), $G_1 = \frac{1}{2}$ converges faster than $G_1 = 1$

- 30-layer ReLU network (right), $G_0 = \frac{1}{2}$ is successful while $G_1 = 1$ fails to learn at all

# Batch Normalization

We can normalize the activations $x_k^{(i)}$ of node $k$ in layer $(i)$ relative to the mean and variance of those activations, calculated over a mini-batch of training items:

$$\hat{x}_k^{(i)} = \frac{x_k^{(i)} - \text{Mean}[x_k^{(i)}]}{\sqrt{\text{Var}[x_k^{(i)}]}}$$

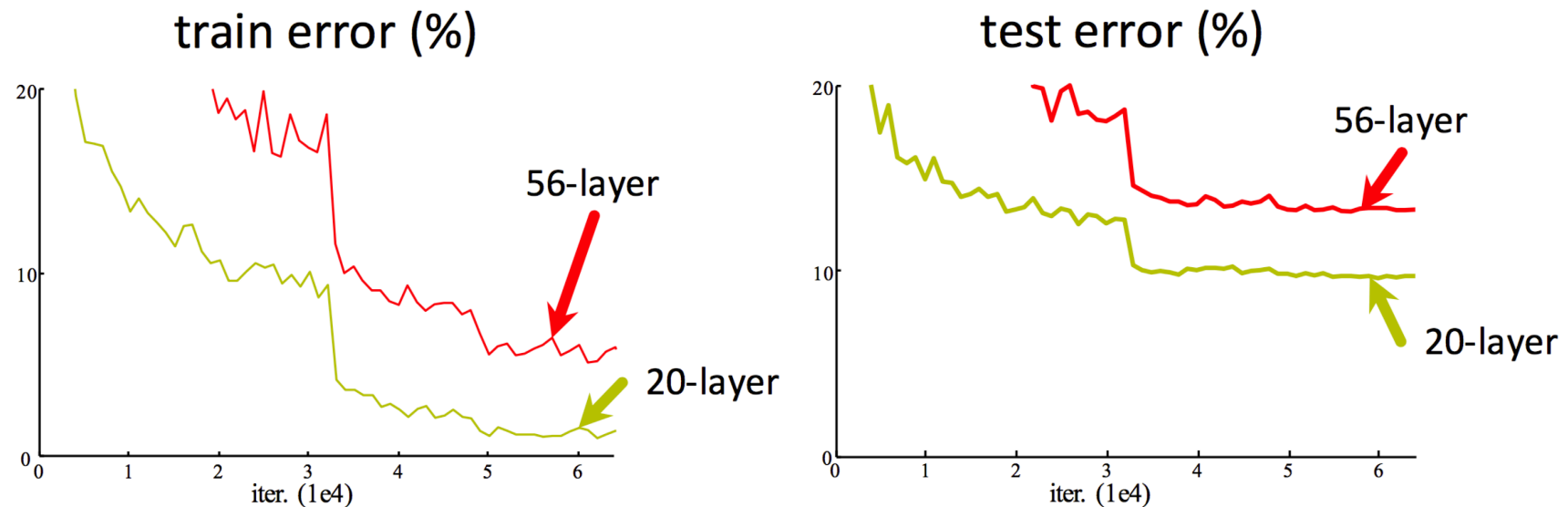These activations can then be shifted and re-scaled to

$$y_k^{(i)} = \beta_k^{(i)} + \gamma_k^{(i)} \hat{x}_k^{(i)}$$

$\beta_k^{(i)}, \gamma_k^{(i)}$ are additional parameters, for each node, which are trained by backpropagation along with the other parameters (weights) in the network.

After training is complete, $\text{Mean}[x_k^{(i)}]$ and $\text{Var}[x_k^{(i)}]$ are either pre-computed on the entire training set, or updated using running averages.
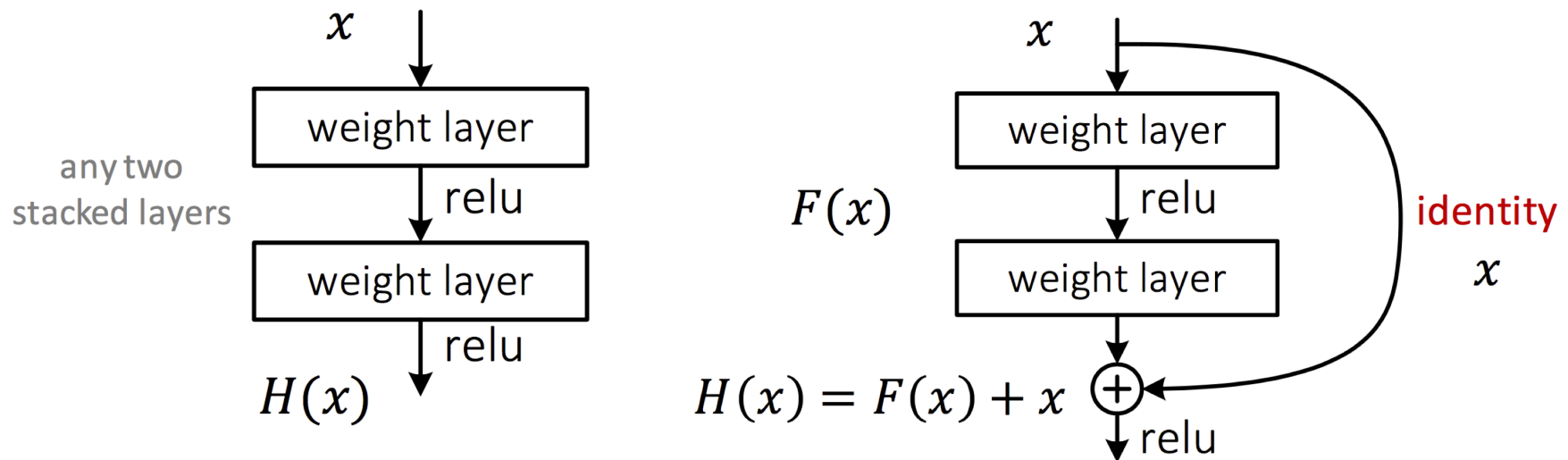
# Going Deeper



If we simply stack additional layers, it can lead to higher training error as well as higher test error.

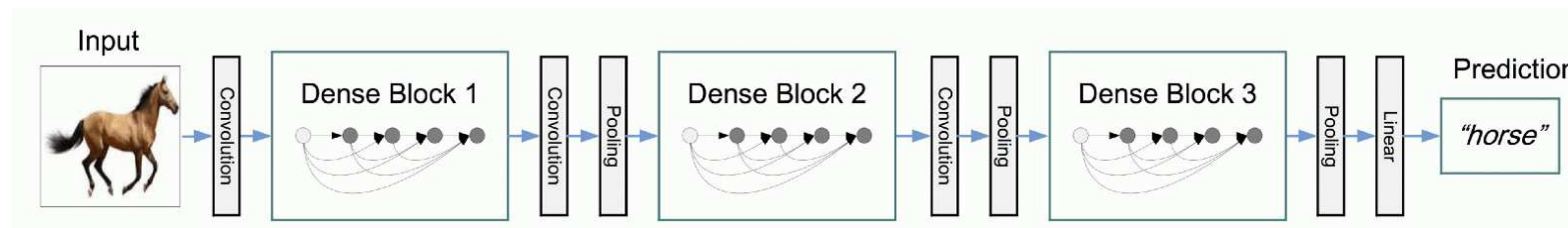# Residual Networks



any two stacked layers

$F(x)$

identity

$x$

$H(x)$

$H(x) = F(x) + x$

relu

Idea: Take any two consecutive stacked layers in a deep network and add a "skip" connection which bipasses these layers and is added to their output.

# Residual Networks

- the preceding layers attempt to do the "whole" job, making $x$ as close as possible to the target output of the entire network

- $F(x)$ is a residual component which corrects the errors from previous layers, or provides additional details which the previous layers were not powerful enough to compute

- with skip connections, both training and test error drop as you add more layers

- with more than 100 layers, need to apply relu before adding the residual instead of afterwards. This is called an identity skip connection.

# Dense Networks



Recently, good results have been achieved using networks with densely connected blocks, within which each layer is connected by shortcut connections to all the preceding layers.

# Texture Synthesis

# Neural Texture Synthesis

1. pretrain CNN on ImageNet (VGG-19)

2. pass input texture through CNN; compute feature map $F_{ik}^l$ for $i^{\text{th}}$ filter at spatial location $k$ in layer (depth) $l$

3. compute the Gram matrix for each pair of features

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l$$

4. feed (initially random) image into CNN

5. compute L2 distance between Gram matrices of original and new image

6. backprop to get gradient on image pixels

7. update image and go to step 5.
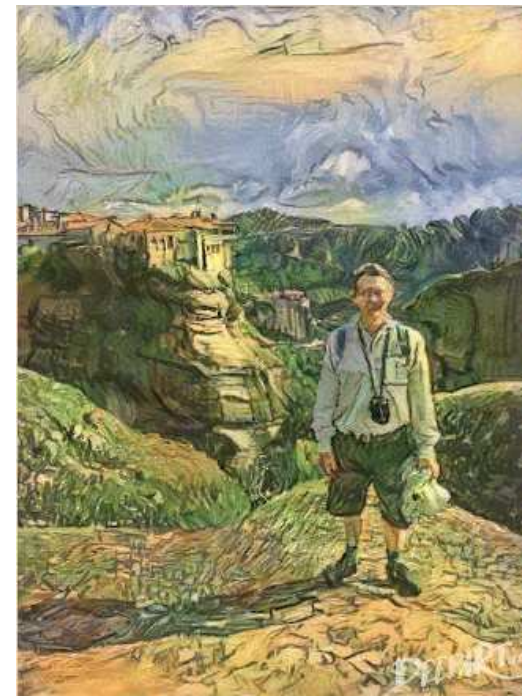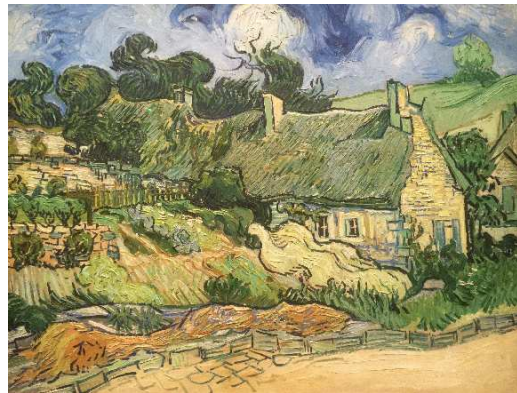
# Neural Texture Synthesis

We can introduce a scaling factor $w_l$ for each layer $l$ in the network, and define the Cost function as

$$E_{\text{style}} = \frac{1}{4} \sum_{l=0}^{L} \frac{w_l}{N_l^2 M_l^2} \sum_{i,j} (G_{ij}^l - A_{ij}^l)^2$$

where $N_l$, $M_l$ are the number of filters, and size of feature maps, in layer $l$, and $G_{ij}^l$, $A_{ij}^l$ are the Gram matrices for the original and synthetic image.

# Neural Style Transfer



content      +      style      →      new image

# Neural Style Transfer

# Neural Style Transfer

For Neural Style Transfer, we minimize a cost function which is

$$E_{\text{total}} = \alpha \; E_{\text{content}} \; + \; \beta \, E_{\text{style}}$$

$$= \frac{\alpha}{2} \sum_{i,k} ||F_{ik}^l(x) - F_{ik}^l(x_c)||^2 + \frac{\beta}{4} \sum_{l=0}^{L} \frac{w_l}{N_l^2 M_l^2} \sum_{i,j} (G_{ij}^l - A_{ij}^l)^2$$

where

$$
\begin{aligned}
x_c, x \quad &= \quad \text{content image, synthetic image} \\
F_{ik}^l \quad &= \quad i^{\text{th}} \text{ filter at position } k \text{ in layer } l \\
N_l, M_l \quad &= \quad \text{number of filters, and size of feature maps, in layer } l \\
w_l \quad &= \quad \text{weighting factor for layer } l \\
G_{ij}^l, A_{ij}^l \quad &= \quad \text{Gram matrices for style image, and synthetic image}
\end{aligned}
$$

# References

- "ImageNet Classification with Deep Convolutional Neural Networks", Krizhevsky et al., 2015.

- "Understanding the difficulty of training deep feedforward neural networks", Glorot & Bengio, 2010.

- "Batch normalization: Accelerating deep network training by reducing internal covariate shift", Ioffe & Szegedy, ICML 2015.

- "Deep Residual Learning for Image Recognition", He et al., 2016.

- "Densely Connected Convolutional Networks", Huang et al., 2016.

- "A Neural Algorithm of Artistic Style", Gatys et al., 2015.