# CS 189/289

Today's lecture:
1. From logistic to softmax.
2. Convolutional neural networks
3. Residual neural networks (resnets)
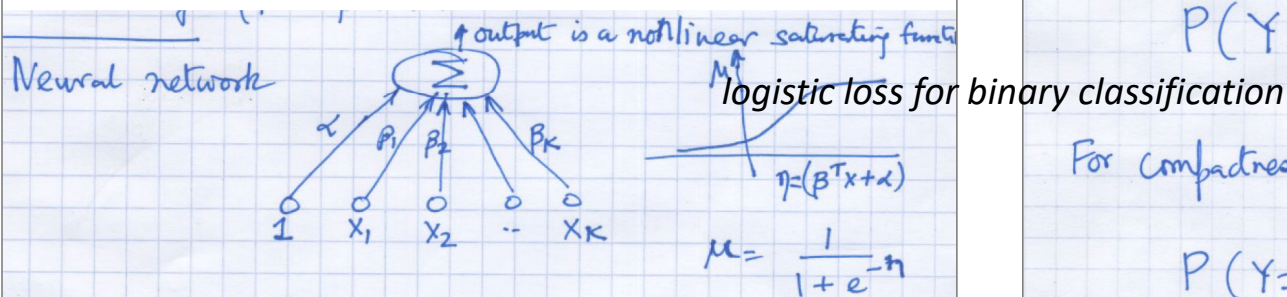
# CS 189/289

Today's lecture:
1. From logistic to softmax.
2. Convolutional neural networks
3. Residual neural networks (resnets)

# Recall: *logistic loss for binary classification*

### Neural networks can be modeled by logistics

Neural network

output is a nonlinear saturating function

$$\eta = (\beta^T x + \alpha)$$

$$\mu = \frac{1}{1 + e^{-\eta}}$$

Standard Trick: Add a $0^{th}$ component to the $\underset{\sim}{x}$ vector, which is fixed to be 1. This is connected with weight $\alpha$

### Modeling the probability distribution

We say that the class label $Y$ is a Bernoulli random variable, with its probability parameter $p$ being as above

$$P(Y = 1 \mid x) = \frac{1}{1 + exp(-\beta^T x)}$$

*logistic loss for binary classification*

For compactness, introduce notation $\mu(x) = \frac{1}{1 + exp(-\beta^T x)}$

$$P(Y = 1 \mid x) = \mu(x)$$

or $\mu(x) = \frac{1}{1 + exp(-\eta(x))}$

As usual we use $y$ to denote values taken by random variables

$$P(y \mid x) = \mu(x)^y (1 - \mu(x))^{1-y}$$

## What if we have more than 2 classes?

# From logistic regression to softmax regression

$$\begin{bmatrix} p(Y = 1|X) \\ p(Y = 0|X) \end{bmatrix} = \begin{bmatrix} \mu \\ 1-\mu \end{bmatrix} = \begin{bmatrix} \dfrac{1}{1 + \exp(-\beta x)} \\ \dfrac{\exp(-\beta x)}{1 + \exp(-\beta x)} \end{bmatrix}$$

# From logistic regression to softmax regression

$$\begin{bmatrix} p(Y = 1|X) \\ p(Y = 0|X) \end{bmatrix} = \begin{bmatrix} \mu \\ \\ 1-\mu \end{bmatrix} = \begin{bmatrix} \dfrac{1}{1 + \exp(-\beta x)} \\ \\ \dfrac{\exp(-\beta x)}{1 + \exp(-\beta x)} \end{bmatrix}$$

Instead we could write this as

$$\frac{1}{e^{\beta_1 x} + e^{\beta_2 x}} \begin{bmatrix} e^{\beta_1 x} \\ \\ e^{\beta_2 x} \end{bmatrix} = \begin{bmatrix} \dfrac{1}{1 + e^{\beta_2 x - \beta_1 x}} \\ \\ \dfrac{1}{1 + e^{-\beta_2 x + \beta_1 x}} \end{bmatrix}$$

equivalent with
$\beta = \beta_1 - \beta_2$

# The softmax function for K-class classification

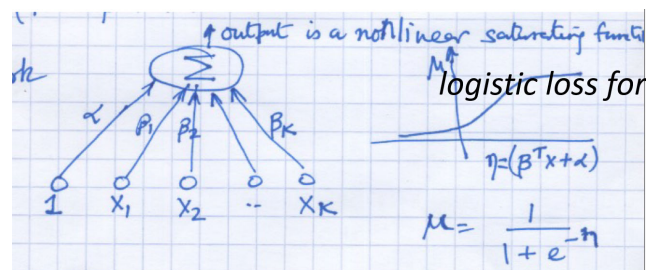$$\begin{bmatrix} p(Y = 1|X) \\ p(Y = 2|X) \\ p(Y = 3|X) \\ \cdots \\ p(Y = K|X) \end{bmatrix} = \frac{1}{\sum_{i=1}^{K} e^{\beta_i X}} \begin{bmatrix} e^{\beta_1 X} \\ e^{\beta_2 X} \\ \\ \\ \\ e^{\beta_K X} \end{bmatrix}$$

- Generalization of logistic regression to more than 2 classes.
- "Softmax regression" or "multinomial logistic regression", parameters $\beta$.
- Use principle of MLE to set $\beta$.
- Needs iterative optimization like gradient descent.
- Can also stick at the top of neural network to get a "softmax" loss.

output is a nonlinear saturating functi

$logistic\ loss\ for$

$\eta = (\beta^T x + \alpha)$

$\mu = \dfrac{1}{1 + e^{-\eta}}$

# The softmax function for K-class classification

$$\begin{bmatrix} p(Y=1|X) \\ p(Y=2|X) \\ p(Y=3|X) \\ \dots \\ p(Y=K|X) \end{bmatrix} = \frac{1}{\sum_{i=1}^{K} e^{\beta_i X}} \begin{bmatrix} e^{\beta_1 X} \\ e^{\beta_2 X} \\ \\ \\ e^{\beta_K X} \end{bmatrix}$$

For class `i`, the logit (log-odds) is defined as:

$$\text{logit}_i = \log \left( \frac{P(y=i|x)}{P(y\neq i|x)} \right)$$

For class `i`, the softmax function is defined as:

$$P(y = i|x) = \frac{e^{\text{logit}_i}}{\sum_{j=1}^{K} e^{\text{logit}_j}}$$



output is a nonlinear saturating function

logistic loss for

$\eta = (\beta^T x + \alpha)$

$\mu = \dfrac{1}{1 + e^{-\eta}}$

# CS 189/289

Today's lecture outline:

1. From logistic to softmax.
2. Convolutional neural networks
3. Residual neural networks (resnets)

# Recall: *fully connected* neural networks



$$O_i = g\left(\sum_j W_{ij}\, g\left(\sum_K W_{jK}\, x_K\right)\right)$$

$O_1$  $O_2$

$O_i$

$W_{11}$  $W_{23}$

$W_{ij}$

$V_1$  $V_2$  $V_3$

$V_j$

$W_{12}$  $W_{35}$  $W_{jK}$

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$  $X_K$

# Recall: feed forward, *fully connected* neural networks



**Computing $\delta$ for neuron in intermediate layer**

Layer $\ell$   $d^{(\ell)}$ nodes

Layer $\ell-1$

Layer 1   $d^{(1)}$ nodes.

Layer 0 (INPUT)

$x_1 \quad x_2 \quad x_3 \quad \cdots \quad x_n$

$$x_j^{(\ell)} = g\left( \sum_{i=0}^{d^{(\ell-1)}} w_{ij}^{(\ell)} \, x_i^{(\ell-1)} \right)$$

$$= g\left( s_j^{(\ell)} \right)$$

$s_j^{(\ell)}$ is the weighted input to node $j$ in layer $\ell$

**INDUCTION STE**

$$\delta_i^{(\ell-1)} = \frac{\partial e(w)}{\partial s_i^{(\ell-1)}}$$

$$= \sum_j \frac{\partial e(w)}{\partial s_j^{(\ell)}} \times \frac{\partial s_j}{\partial x_i^{(\ell-1)}} \times \frac{\partial x_i^{(\ell-1)}}{\partial s_i^{(\ell-1)}}$$

Back-propagation algorithm to compute derivatives of the parameters efficiently.

Beyond fully connected, feed-forward architectures:

1. *Convolutional*
2. *Residual*
3. *Recurrent* (not "feed-forward").
4. *Attention* and *Transformers*.
5. *Graph*

- As long as we have a feed-forward network, and use only differentiable components, we can apply backprop.

- New architectures have led to break-through successes.

# Pondering fully connected neural networks

- For "fully connected" (FC) layer, $l$, with $n_i(l)$ inputs and $n_o(l)$ outputs, $W_l$ contains $n_i(l) \times n_o(l)$ parameters.
- Adds up quickly to huge #s of parameters.
- Too many parameters can contribute to problems of "overfitting".

# Pondering fully connected neural networks

- For "fully connected" (FC) layer, $l$, with $n_i(l)$ inputs and $n_o(l)$ outputs, $W_l$ contains $n_i(l) \times n_o(l)$ parameters.
- Adds up quickly to huge #s of parameters.
- Too many parameters can contribute to problems of "overfitting".



➢ Strategy to reduce # of free parameters: "bake" in properties that encode problem symmetries.

# Examples of common problem symmetries

*translation invariance*



Predict: is a cat *vs.* not a cat

*translation equivariance*



Predict: which pixels are cat pixels?

# Examples of common problem symmetries

*permutation in**variance***



*permutation **equi**variance*



$$f(x) = f(Perm(x))$$

$$Perm(f(x)) = f(Perm(x))$$

Predict vector output.

# Examples of common problem symmetries

*rotation invariance*

*rotation equivariance*

predict phase (is liquid?)
at room temperature

predict forces (vector)

[from David Rothchild]

# Examples of common problem symmetries

*translation invariance*                    *translation equivariance*

- The *convolution* operation is <u>translation equivariant</u>.
- This operation will form the basis of *convolutional neural networks (CNNs).*
- CNNs also be motivated by the idea of learning re-usable features (next).

'cat'                    'cat'

Predict: is a cat *vs.* not a cat       Predict: which pixels are cat pixels?

# Features sharing across one input example

"Features" (e.g. is there an eye here?) constructed in <u>fully connected layer</u> cannot be shared across the input (e.g. image), because $w$ is not reused across the image.
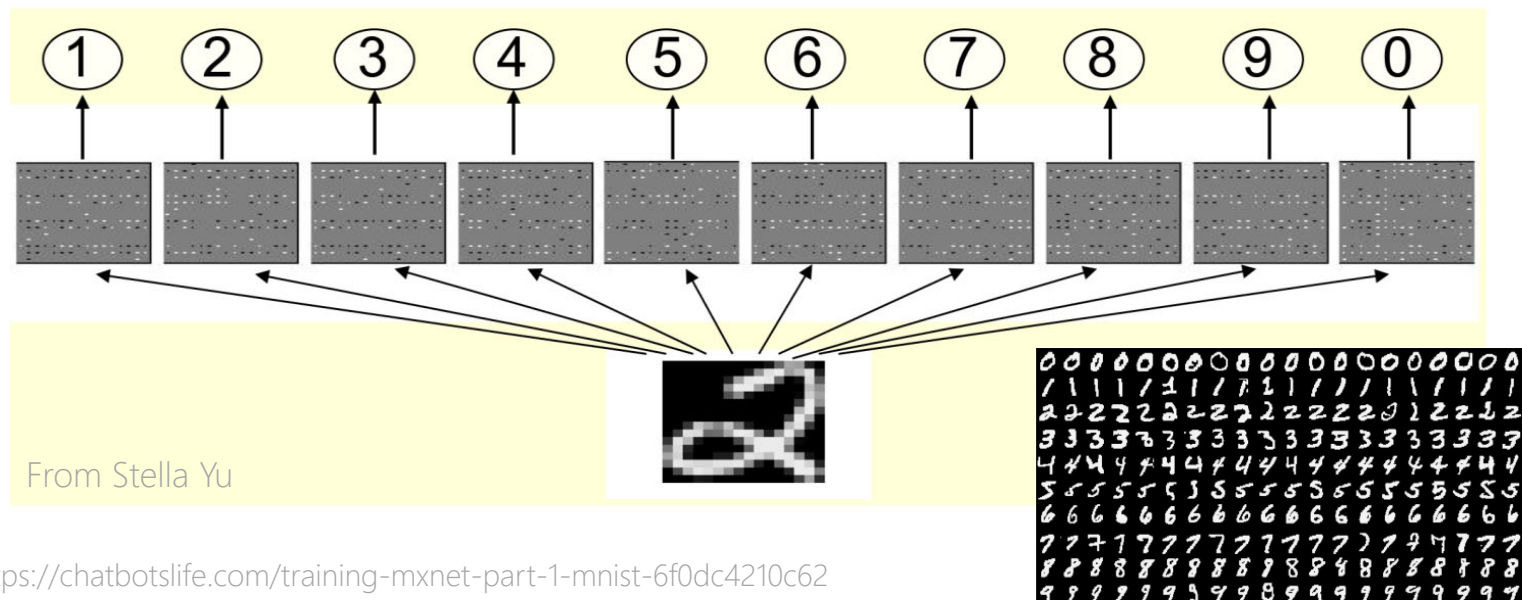


One neuron in FC layer:

$$y = y$$

$w$

$x$

$x$ is the cat matrix flattened to a 1D vector

w operates on the entire image

input layer    hidden layer 1   hidden layer 2   hidden layer 3

output layer

# Features sharing across one input example

"Features" (e.g. is there an eye here?) constructed in <u>fully connected layer</u> cannot be shared across the input (e.g. image), because $w$ is not reused across the image.



With Conv layer:

$w$

$p_1$ $p_2$ $\cdots$ $p_{\#\text{patches}}$

col2im

im2col

use <u>2D</u> image patches as input

Now w takes (overlapping) patches

One neuron in FC layer:

$w$

$y$

$y$

$x$

$x$ is the cat matrix flattened to a <u>1D</u> vector

w operates on the entire image

- ConvNet: learn shared features that are applied to every image patch.
- Also gives us *translational equivariance* for each filter ($w$) response.

# Fully Connected (FC): no feature sharing

- Uses "global template matching".
- e.g. one $W$ matrix per class (single layer):

Iteration 1 of training:



From Stella Yu
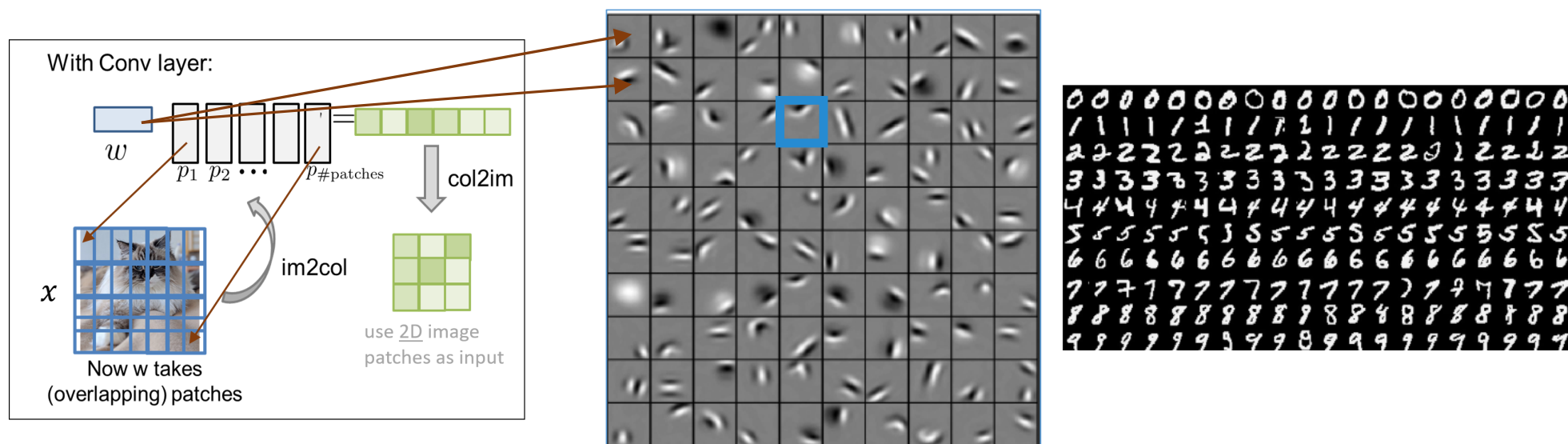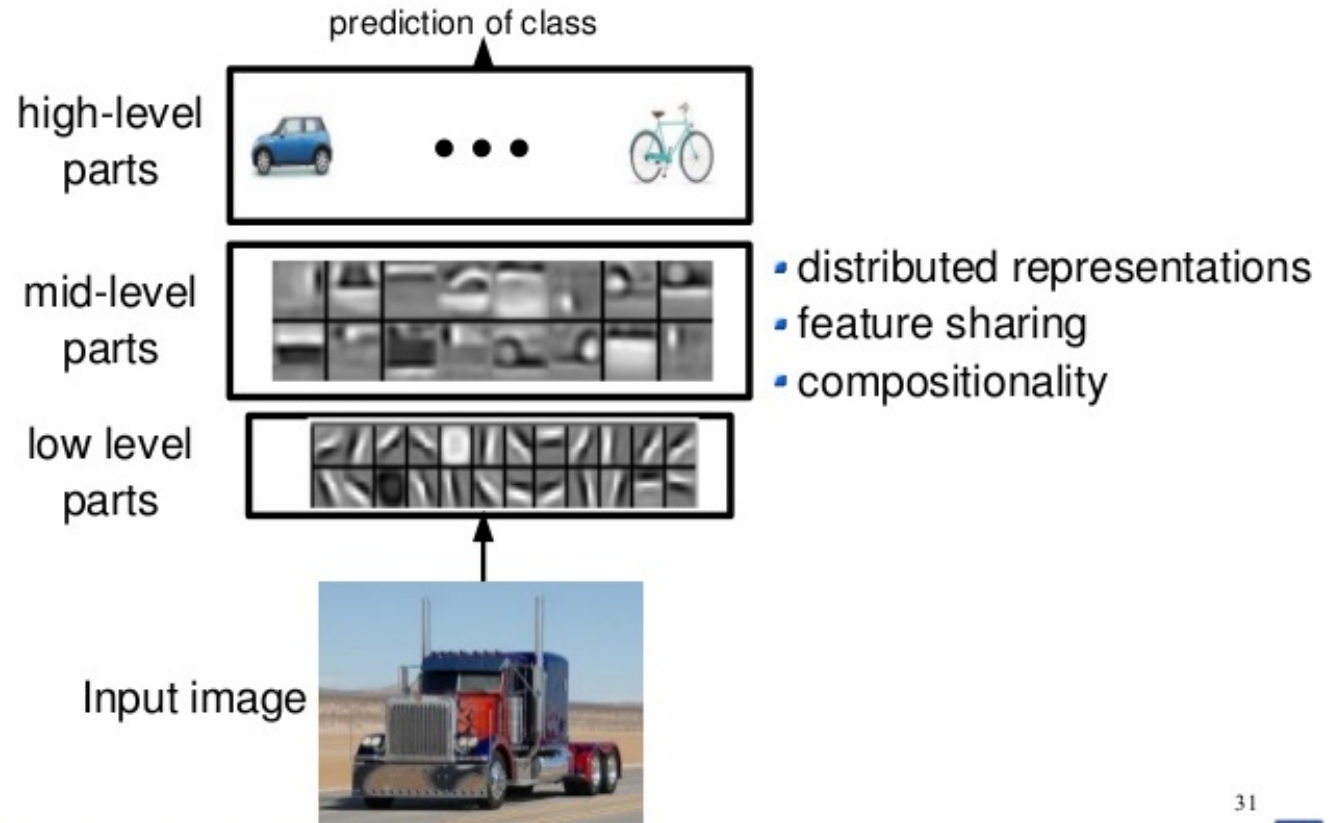
# Fully Connected (FC): no feature sharing

- Uses "global template matching".
- e.g. one $W$ matrix per class (single layer):

Iteration 2 of training:



From Stella Yu

https://chatbotslife.com/training-mxnet-part-1-mnist-6f0dc4210c62

# Fully Connected (FC): no feature sharing

- Uses "global template matching".
- e.g. one $W$ matrix per class (single layer):

Iteration 3 of training:



From Stella Yu

https://chatbotslife.com/training-mxnet-part-1-mnist-6f0dc4210c62

# Fully Connected (FC): no feature sharing

- Uses "global template matching".
- e.g. one $W$ matrix per class (single layer):

Iteration 7 of training:



From Stella Yu

https://chatbotslife.com/training-mxnet-part-1-mnist-6f0dc4210c62

# What would re-usable features look like?

- What if we could learn "local feature filters"
- Then on the next layer, learn how to combine them?

# Convolutional NNs (CNNs/"Convnets")

Can view CNNs as a way to construct hierarchical features, each of which get combined at the next level.



prediction of class

high-level parts

mid-level parts

low level parts

Input image

- distributed representations
- feature sharing
- compositionality

Lee et al. "Convolutional DBN's ..." ICML 2009

Ranzato

https://www.slideshare.net/milkers/lecture-06-marco-aurelio-ranzato-deep-learning

# Convolutional NNs (CNNs/"Convnets")

Can view CNNs as a way to construct hierarchical features, each of which get combined at the next level.



Layer 3

Layer 2

Layer 1

# Convolutional NNs (CNNs/"Convnets")

With Conv layer:

$w$

$p_1$ $p_2$ $\cdots$ $p_{\#patches}$

col2im

im2col

Now w takes
(overlapping) patches

use 2D image
patches as input

Layer 3

Layer 2

Layer 1

"1D" Conv.

# (2D) Convolution



*Convolve* one learned "filter", $W$ with the input to get convolution output $\{v_{ij}\}$:

For each position, $i, j$:

1. Element-wise product of $W$ with image patch centered on $i, j$ (e.g. $3 \times 3$).
2. Sum up the results to get one $v_{ij}$.

*$W$ called filter/template/kernel*

# Convolutional NNs (CNNs/"Convnets")

- We will actually use multiple feature maps, $\{W_k\}_{k=1}^K$
- "Depth" of output "volume" is $K$:

$$h_{ij} = \sigma(v_{ij}).$$

Convnet
Filter

One
Feature
Map

All Feature Maps

**Non-linearity** to get hidden node in a hidden layer in CNN

# Formally: 1D convolution

$b \in \mathbb{R}^7$ $\quad a \in \mathbb{R}^3$ $\quad a * b \in \mathbb{R}^5$



$t = 8$
$t = 5$

- For n-dim convolution, we use an n-dim filter.
- So 1D convolution has a 1D filter.

If $a$ and $b$ are two arrays,

$$(a * b) = (b * a)$$
$$a * (b * c) = (a * b) * c$$

$$(a * b)_t = \sum a_\tau b_{t-\tau}$$

$t$'th element of the convolution

$\tau \in [0, 1, 2, \ldots]$ ← arbitrary

- $\tau$ *is the index of the filter element ('-' means flip filter first)*
- Invalid indices, e.g., $t = 1,2,3$ and $\tau = 3$, are boundaries; don't compute those $t^{th}$ entries, or else *pad out* e.g. with zeros/mirroring input.
- No padding, size of output is $D - K + 1$ for $D$ length input, $K$ length filter.

Cross-correlation: $(a \otimes b)_t = \sum_\tau a_\tau b_{t+\tau}$

# 1D convolution



Method 1:  flip-and-filter

$$(a * b)_t = \sum_{\tau} a_\tau b_{t-\tau}$$

# 1D convolution

Method 2: translate-and-scale

$$(a * b)_t = \sum_{\tau} a_\tau b_{t-\tau} \ =$$

# 1D convolution

Method 3
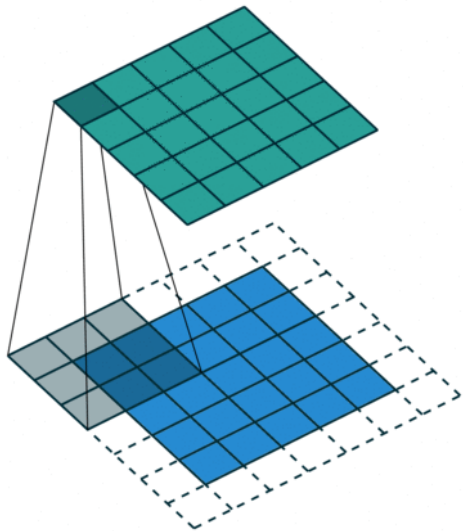
Convolution can also be viewed as matrix multiplication:

$$(a * b)_t = \sum_{\tau} a_\tau b_{t-\tau} = (2, -1, 1) * (1, 1, 2) = \overset{W_k}{\begin{pmatrix} 1 & & \\ 1 & 1 & \\ 2 & 1 & 1 \\ & 2 & 1 \\ & & 2 \end{pmatrix}} \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} =$$

$W_{k'}$ has size **5 × 3**, which means it has 15 entries, yet there are only 3 parameters. Why Convnets to have relatively few parameters!

# From 1D to 2D convolution

$$(A * B)_{ij} = \sum_s \sum_t A_{st} B_{i-s,j-t}$$

Method 1: Flip-and-Filter

# From 1D to 2D convolution

$$(A * B)_{ij} = \sum_s \sum_t A_{st} B_{i-s,j-t}$$

Method 2: Translate-and-Scale

# 2D convolution
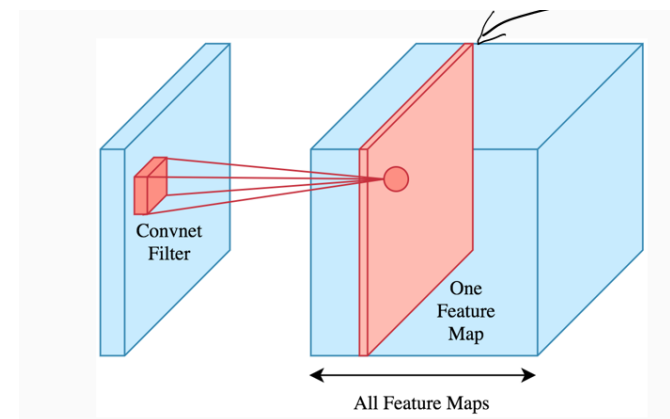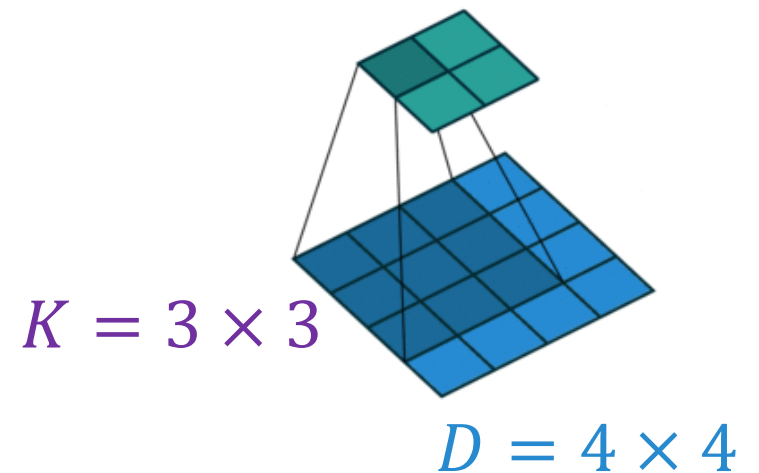
convolution is $2 \times 2$

- Image is $D \times D$.
- $N$ filters each of size $K \times K$.
- No zero-padding.

$K = 3 \times 3$

Then output from one filter has size:
$$(D - K + 1) \times (D - K + 1)$$

$D = 4 \times 4$

For all $N$ filters,
$$N \times (D - K + 1) \times (D - K + 1)$$



Convnet
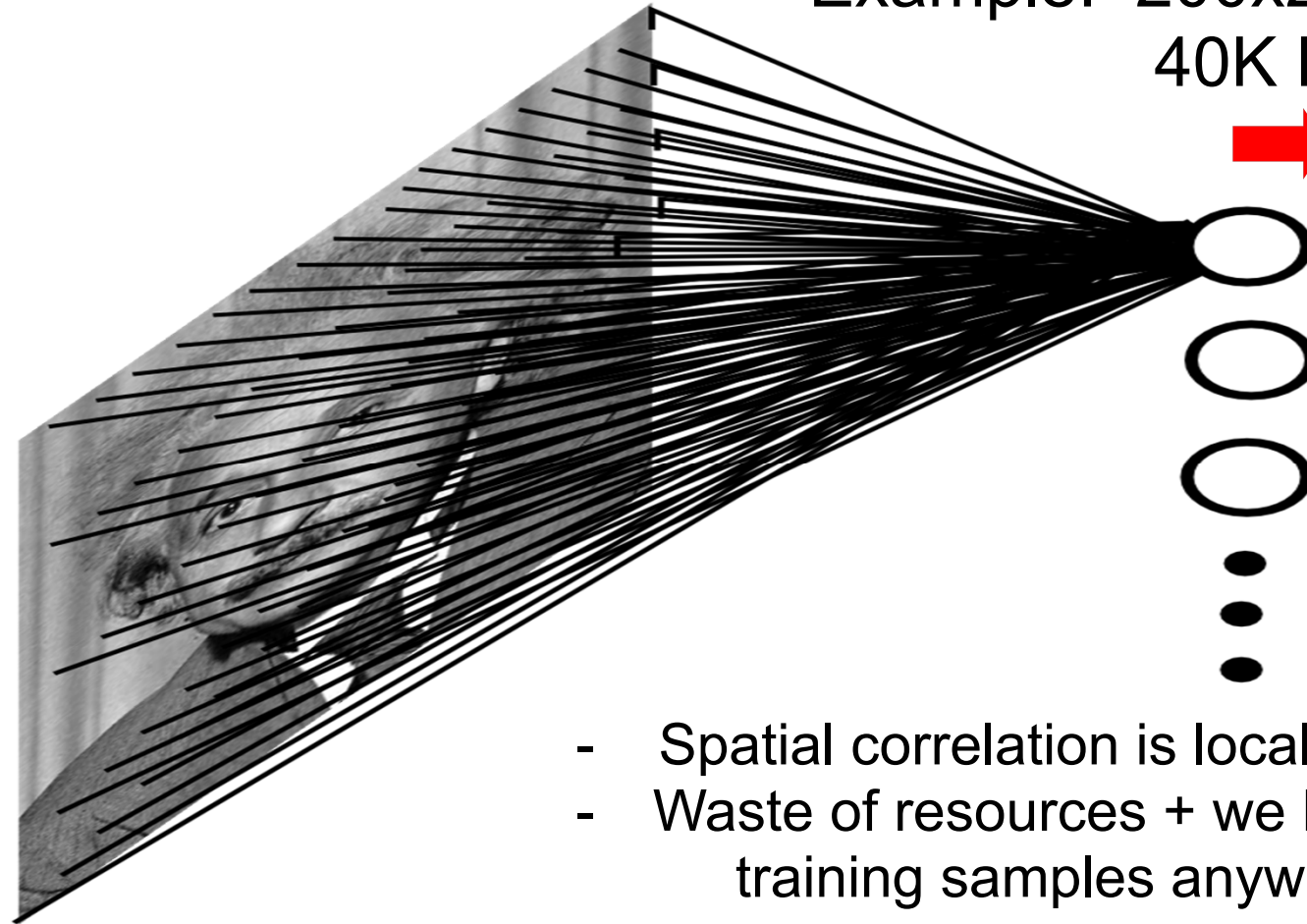Filter

One
Feature
Map

All Feature Maps

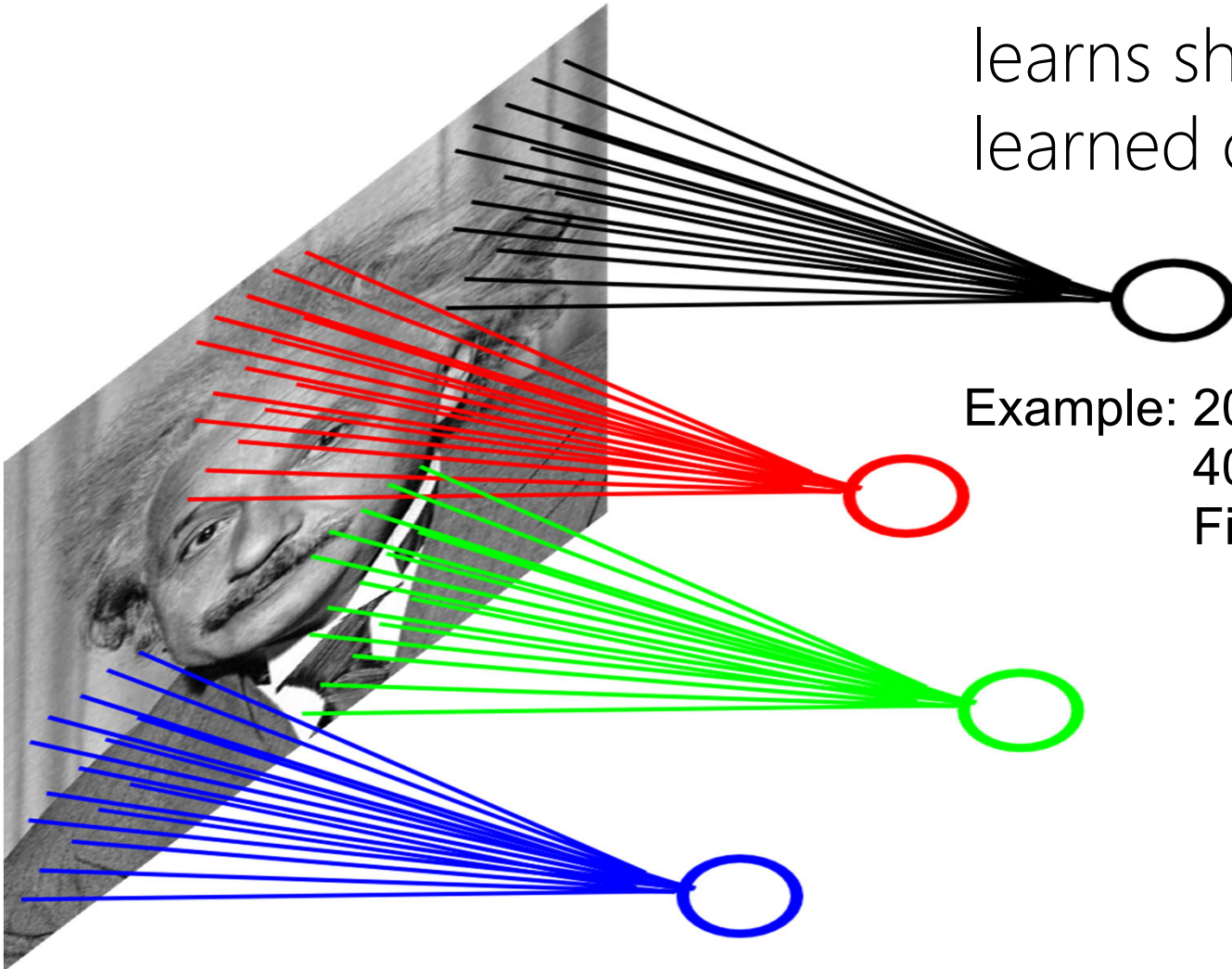# Fully-connected layer (no shared features)

Example: 200x200 image
40K hidden units

➡ **~2B parameters**!!!



- Spatial correlation is local
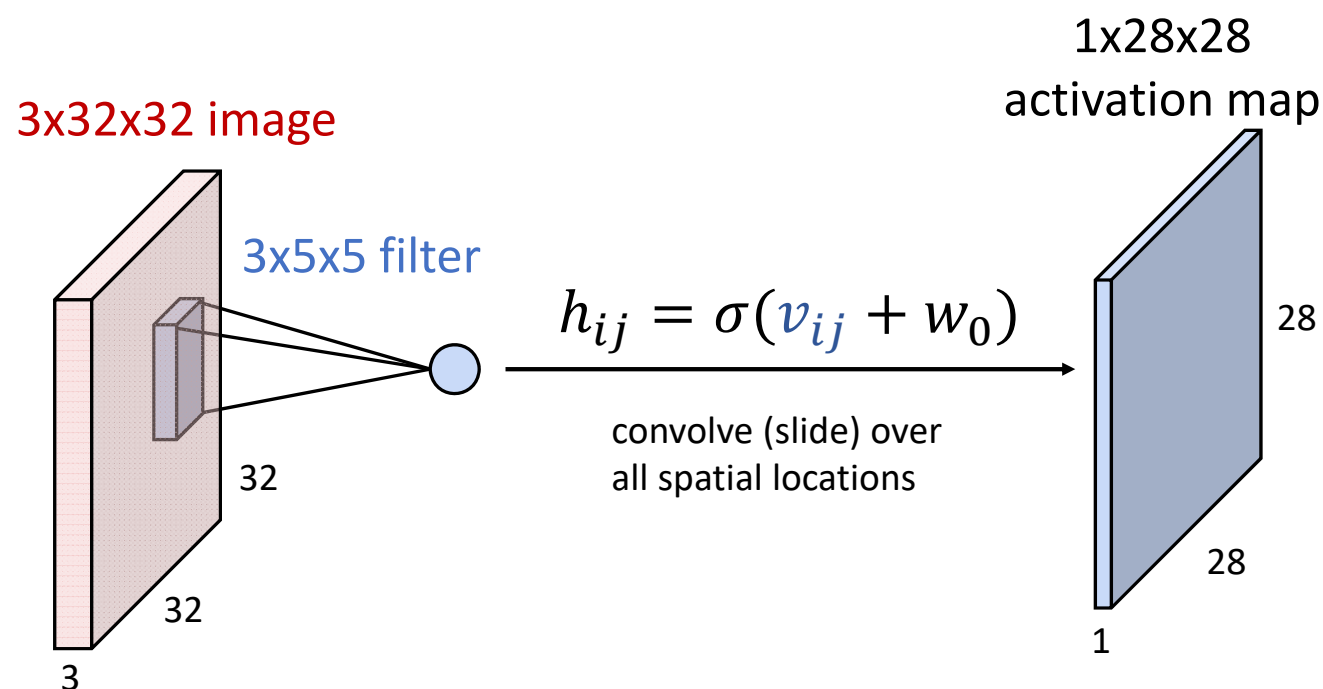- Waste of resources + we have not enough training samples anyway..

Convolutional layer

learns shared features via learned convolution kernels

Example: 200x200 image
40K hidden units
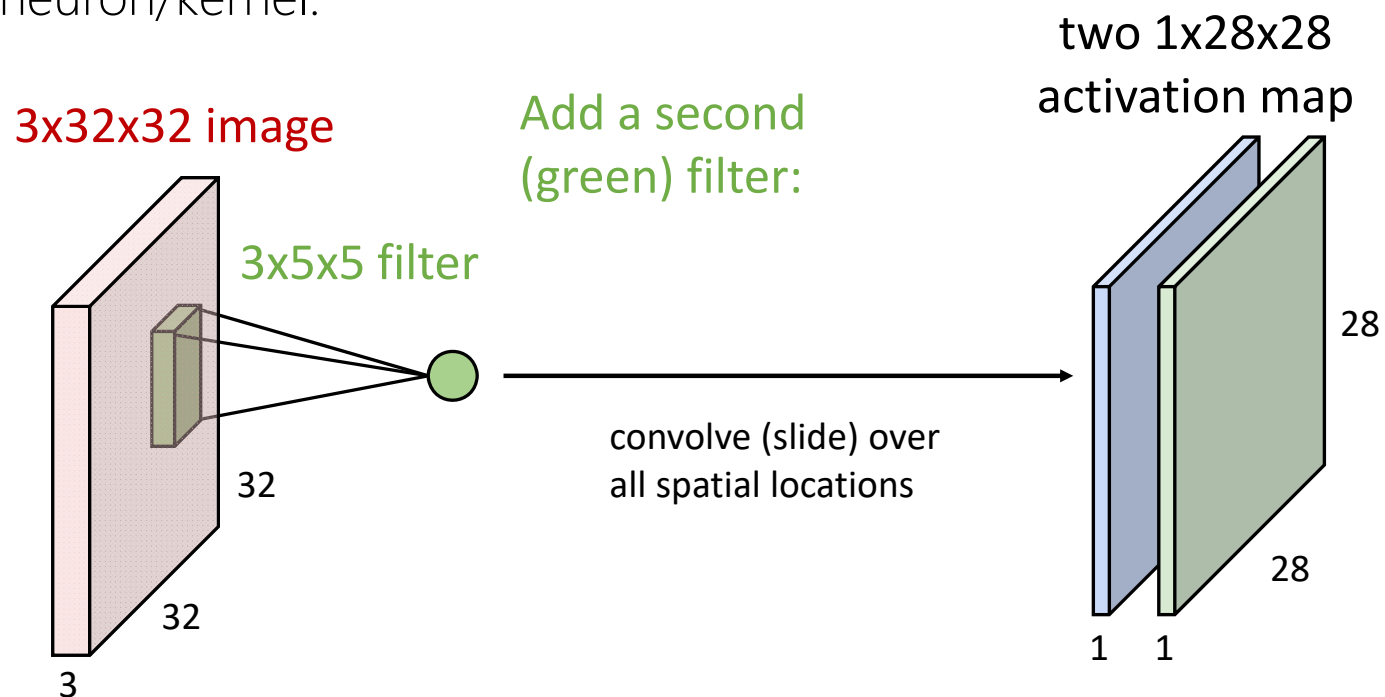Filter size: 10x10
**4M parameters**

# Convolution Layer

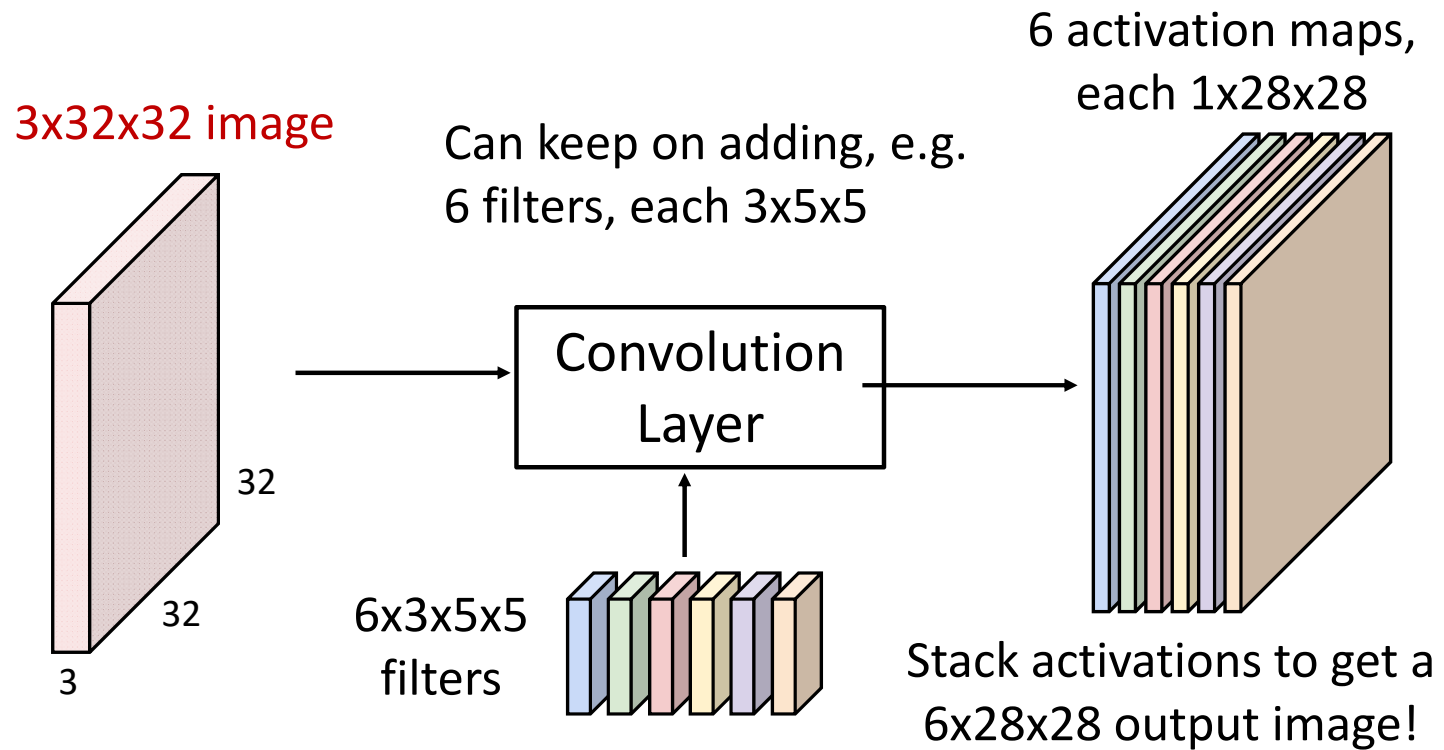One "neuron"/kernel that "looks at" 5x5 region and outputs a sheet of *activation map*



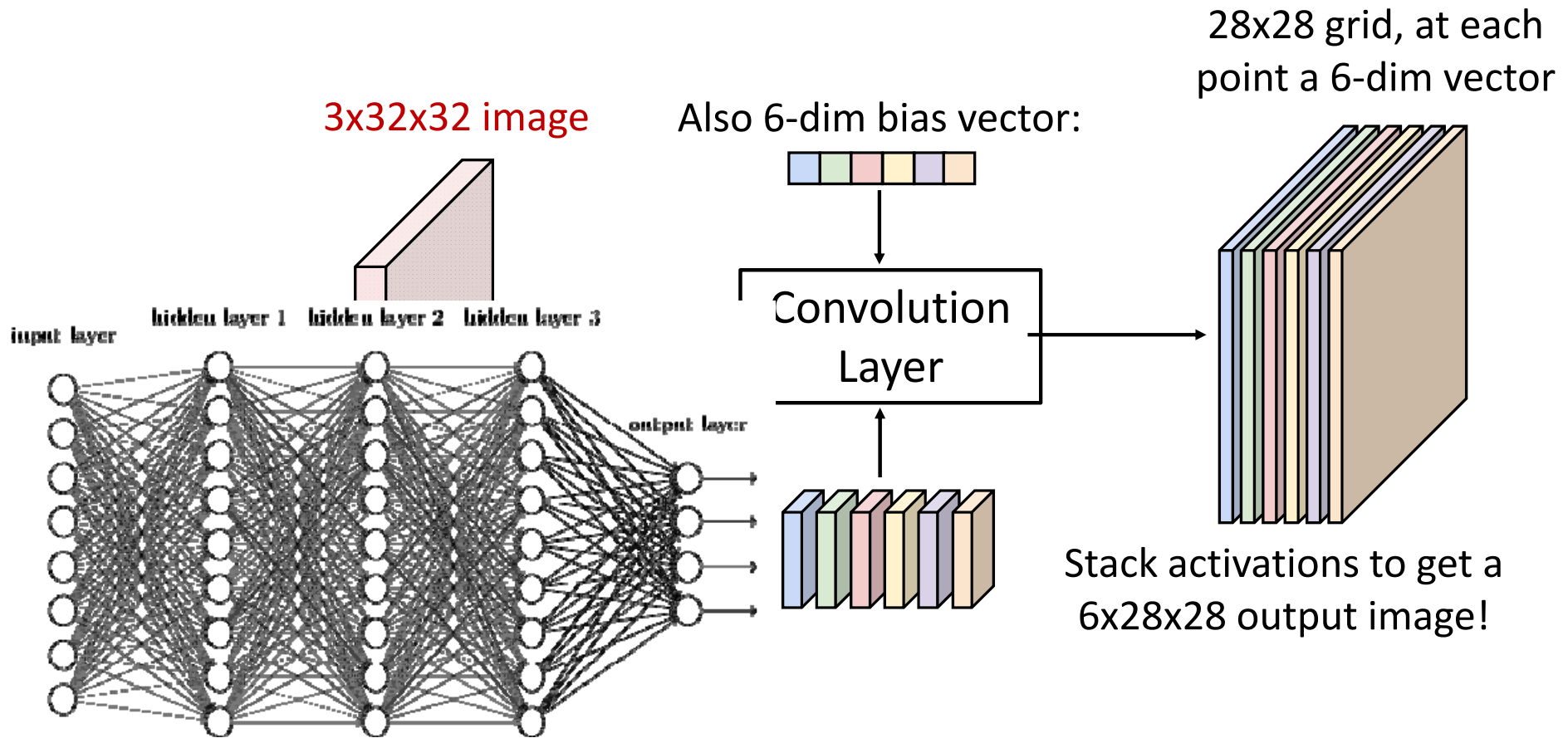**3x32x32 image**

**3x5x5 filter**

32

32

3

$$h_{ij} = \sigma(v_{ij} + w_0)$$

convolve (slide) over
all spatial locations

**1x28x28
activation map**

28

28

1

# Convolution Layer

Add a second neuron/kernel.

**3x32x32 image**

**3x5x5 filter**

32

32

3

**Add a second (green) filter:**

convolve (slide) over all spatial locations

**two 1x28x28 activation map**

28

28

1    1

# Convolution Layer



3x32x32 image

Can keep on adding, e.g.
6 filters, each 3x5x5

6 activation maps,
each 1x28x28

Convolution
Layer

6x3x5x5
filters

32

32

3

Stack activations to get a
6x28x28 output image!

# Convolution Layer

3x32x32 image

Also 6-dim bias vector:

28x28 grid, at each point a 6-dim vector

Convolution Layer

Stack activations to get a 6x28x28 output image!

# Intuition of 2D convolution kernels



*"blurring" filter*

# Intuition of 2D convolution kernels



"oriented edges"

# Intuition of 2D convolution kernels



$$
* \quad
\begin{array}{|c|c|c|}
\hline
0 & -1 & 0 \\
\hline
-1 & 4 & -1 \\
\hline
0 & -1 & 0 \\
\hline
\end{array}
\quad = 
$$

*"sharpen"*

# Intuition of 2D convolution kernels



$$
* \quad
\begin{array}{|c|c|c|}
\hline
w_{11} & w_{12} & w_{13} \\
\hline
w_{21} & w_{22} & w_{23} \\
\hline
w_{31} & w_{32} & w_{33} \\
\hline
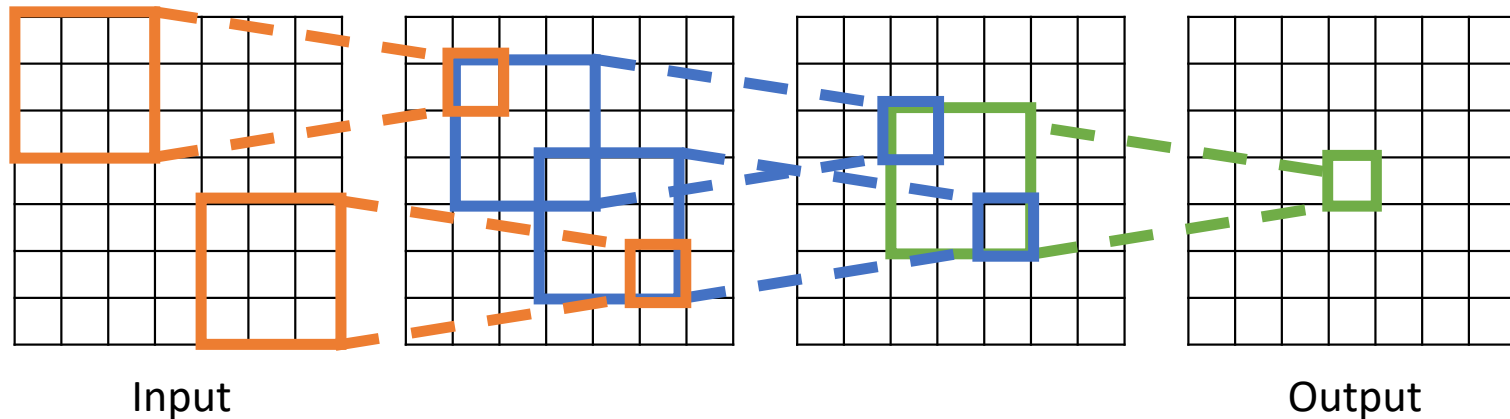\end{array}
\quad = \quad ?
$$

*Gradient descent on loss will decide.*

# Receptive Fields

For convolution with kernel size K, each element in the output depends on a K x K **receptive field** in the input

Input          Output

# Receptive Fields
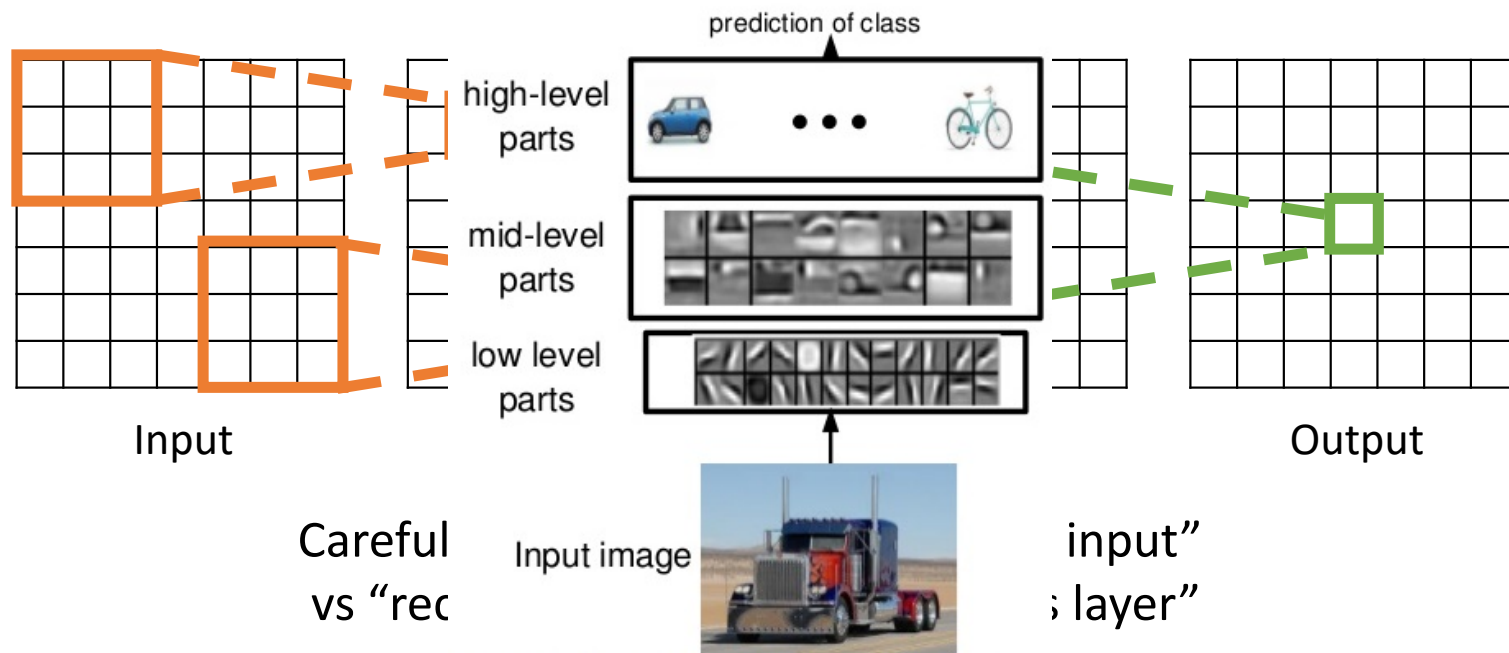
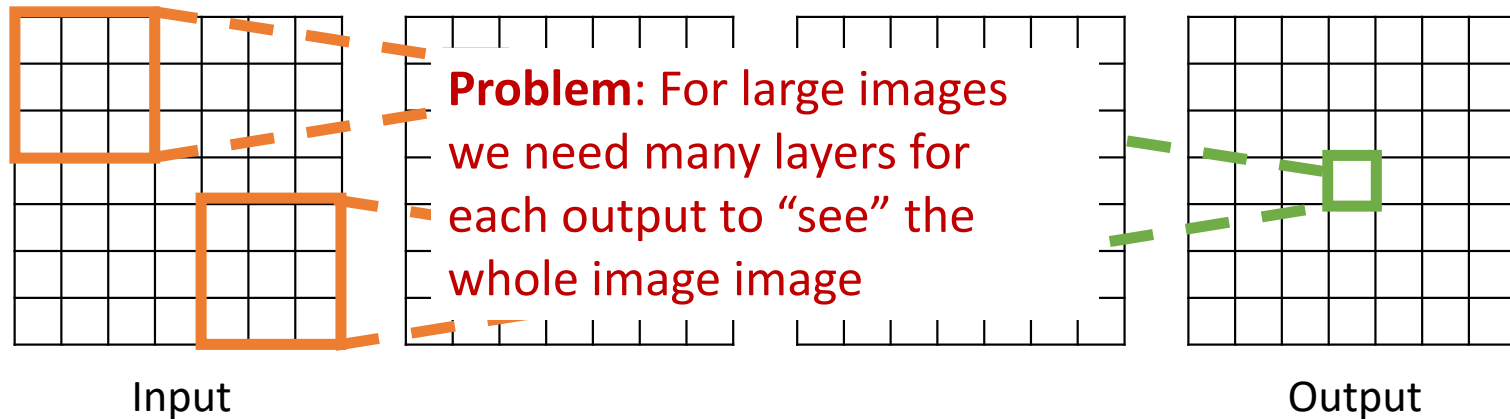Each successive convolution adds $K - 1$ to the receptive field size

With $L$ layers the receptive field size is $1 + L \times (K - 1)$



Input                                                                    Output

Careful – "receptive field wrt to the input"
vs "receptive field wrt the previous layer"

# Receptive Fields

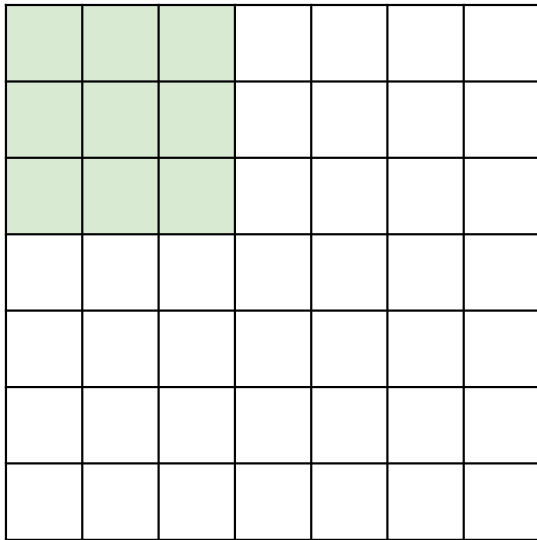Each successive convolution adds $K - 1$ to the receptive field size
With $L$ layers the receptive field size is $1 + L \times (K - 1)$



prediction of class

high-level parts

mid-level parts

low level parts

Input

Input image

Output

Careful ... input"
vs "rec ... layer"

Lee et al. "Convolutional DBN's ..." ICML 2009

# Receptive Fields

Each successive convolution adds $K - 1$ to the receptive field size

With $L$ layers the receptive field size is $1 + L \times (K - 1)$



**Problem**: For large images we need many layers for each output to "see" the whole image image

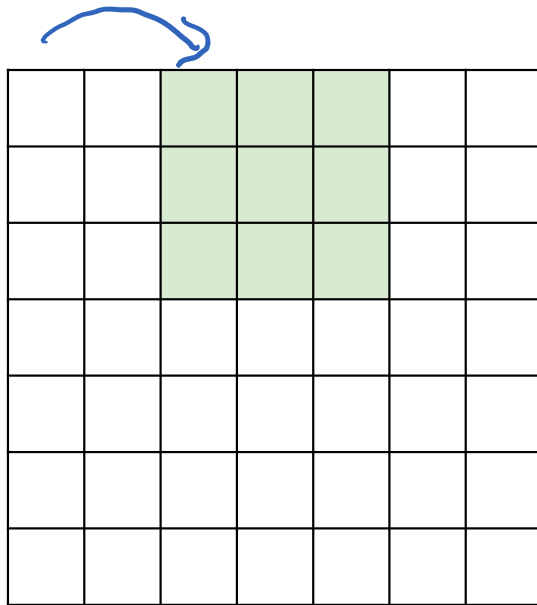Input

Output

Solution: downsample inside the network

1. "Strided" convolution
2. Pooling

# 1. Strided Convolution

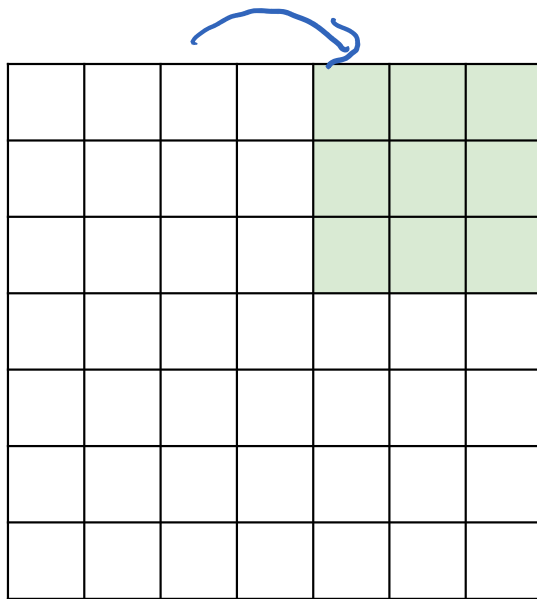Input: 7x7
Filter: 3x3
Stride: 2

# 1. Strided Convolution

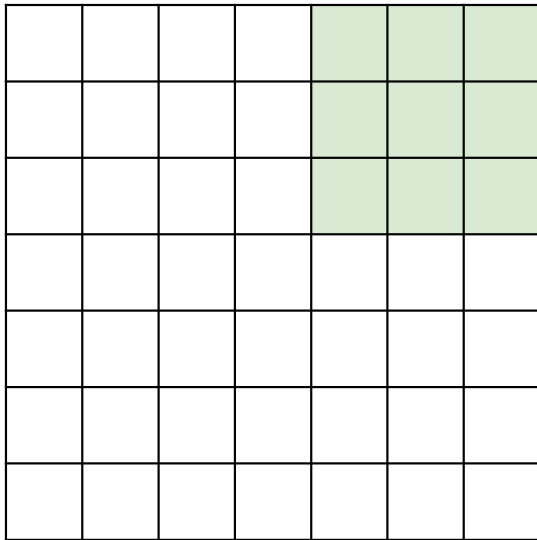Input: 7x7
Filter: 3x3
Stride: 2

# 1. Strided Convolution

Input: 7x7
Filter: 3x3　　Output: 3x3
Stride: 2

# 1. Strided Convolution

Input: 7x7
Filter: 3x3        Output: 3x3
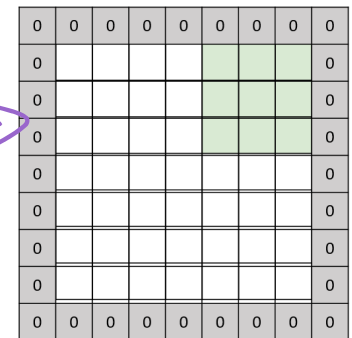Stride: 2

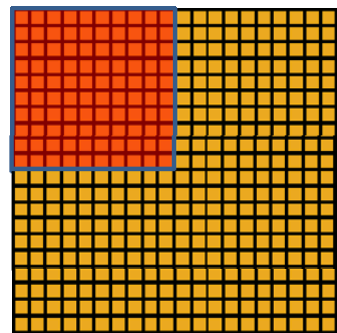In general:
Input: W
Filter: K
Padding: P
Stride: S

Output dimension: (W − K + 2P) / S + 1
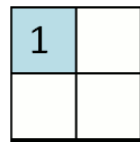(one dimension of the output square)
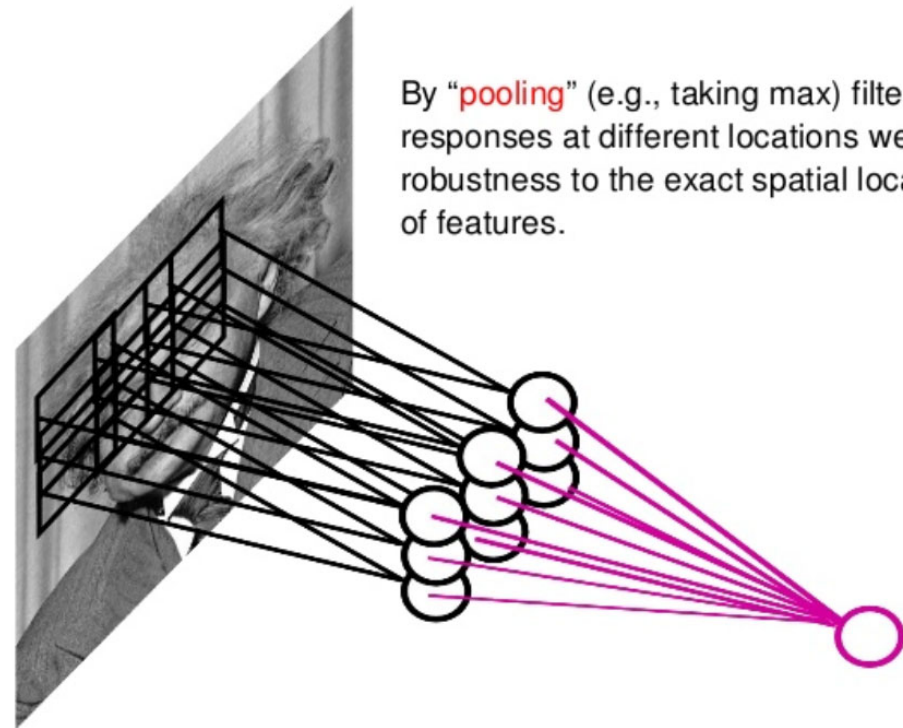
# 2. Pooling layers downsample its inputs

Also adds some local translational *invariance* (by summing/averaging):
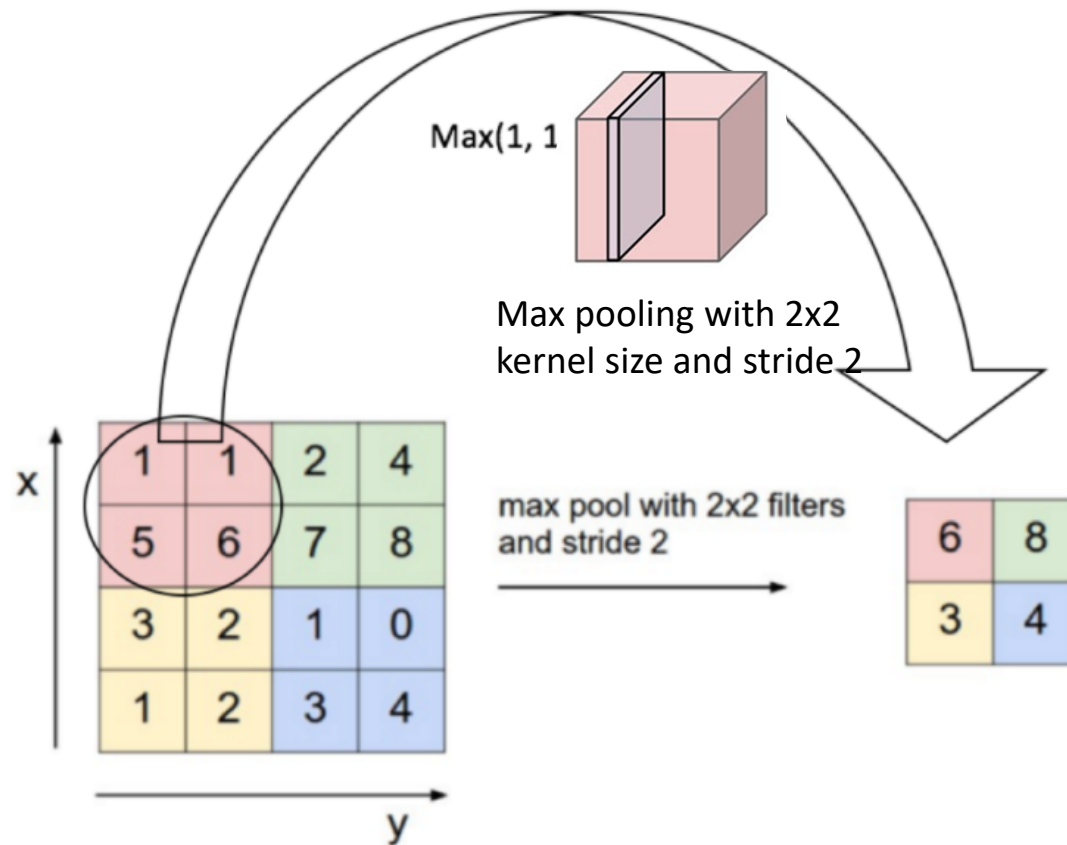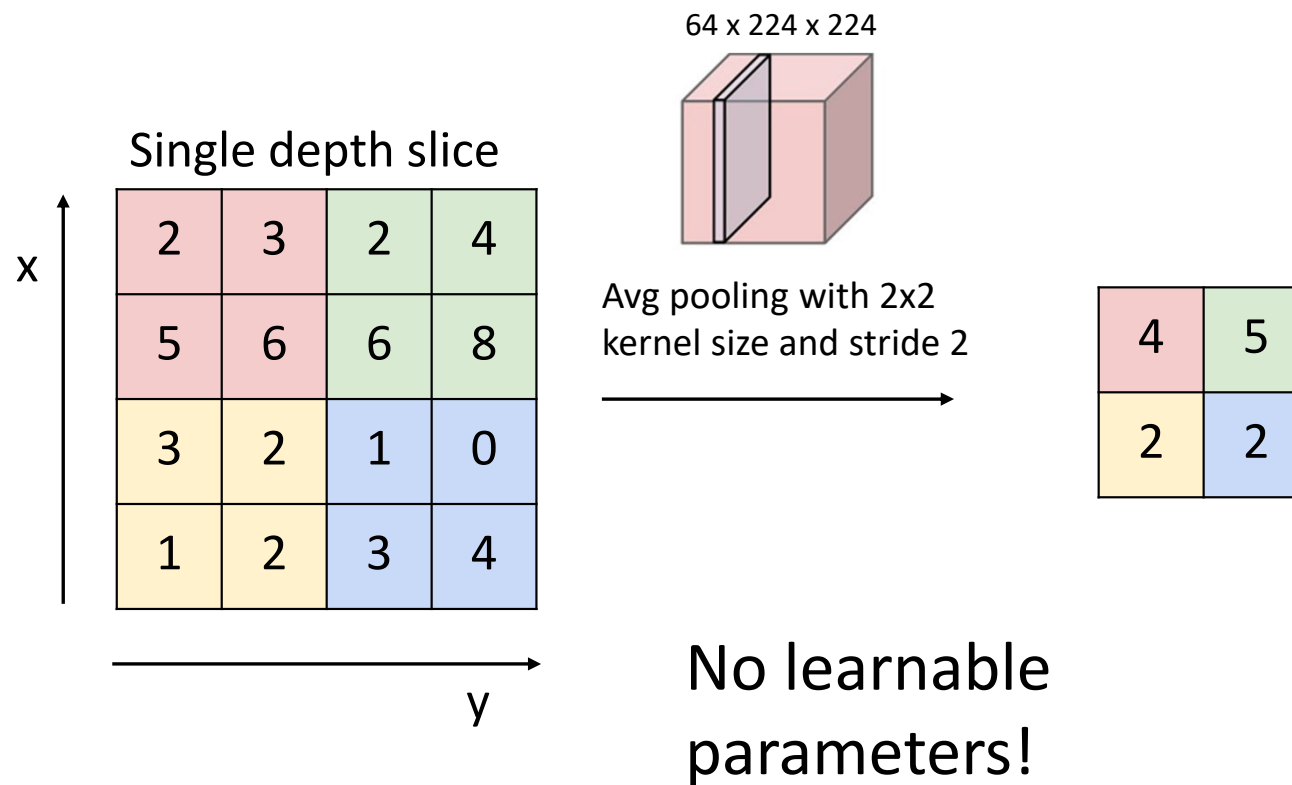


Convolved feature

Pooled feature

By "pooling" (e.g., taking max) filter responses at different locations we gain robustness to the exact spatial location of features.
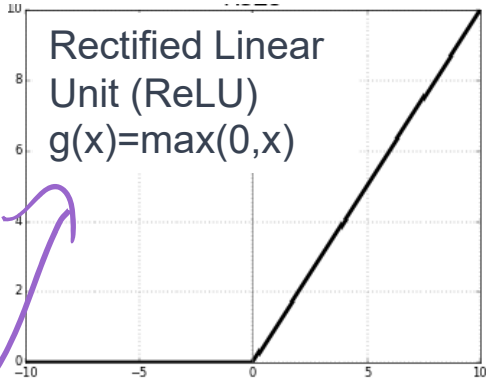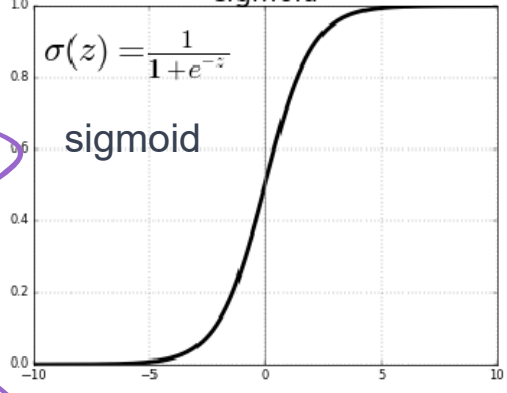
73

**Ranzato** f

# 2. Max pooling

No learnable parameters!

# 2. Average pooling

Single depth slice



64 x 224 x 224

Avg pooling with 2x2
kernel size and stride 2

No learnable
parameters!

# Side note: sigmoid vs ReLU non-linearity in NNs



$$O_i = g\left(\sum_j W_{ij}\, g\left(\sum_K W_{jK}\, X_K\right)\right)$$

$\sigma(z) = \frac{1}{1+e^{-z}}$

sigmoid

Rectified Linear Unit (ReLU)
$g(x)=\max(0,x)$

ReLU:
1. Gradient doesn't die in one direction.
2. More efficient to compute.
3. Easier to get exactly zero activations: sparsity.

# Putting it altogether! ConvNets: conv + ReLU + pooling



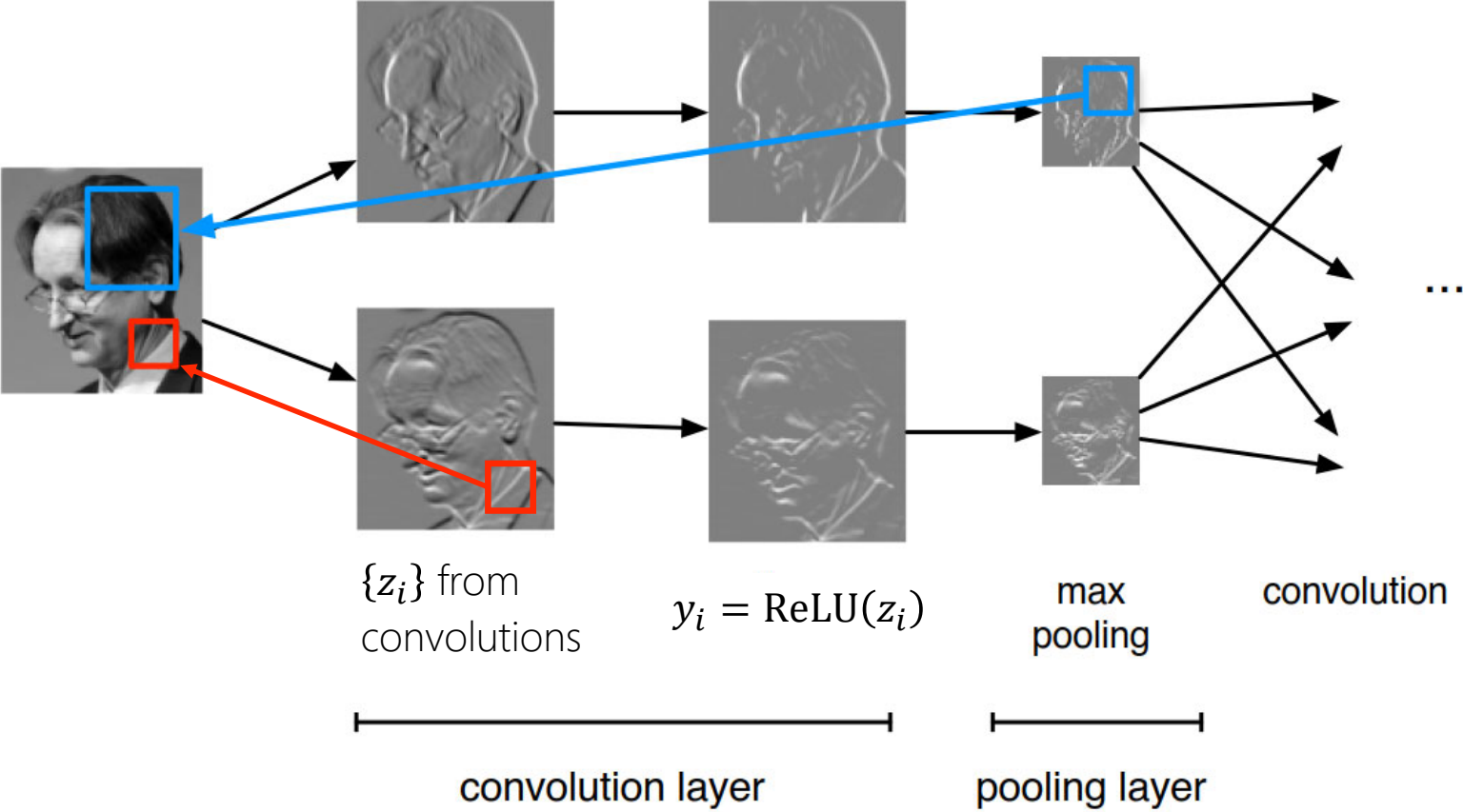$\{z_i\}$ from convolutions

$y_i = \mathrm{ReLU}(z_i)$

max pooling

convolution

convolution layer

pooling layer

# Putting it altogether! ConvNets: conv + ReLU + pooling

Receptive field increases



$\{z_i\}$ from convolutions

$y_i = \text{ReLU}(z_i)$

max pooling

convolution

convolution layer

pooling layer

# Example CNN architecture



| Inputs 3@32x32 | Feature maps 32@18x18 | Feature maps 32@10x10 | Feature maps 48@6x6 | Feature maps 48@4x4 | Hidden units 768 | Hidden units 500 | Outputs 2 |

Convolution 5x5 kernel | Max-pooling 2x2 kernel | Convolution 5x5 kernel | Max-pooling 2x2 kernel | Flatten | Fully connected | Fully connected

# Training CNNs

Gradient descent with back-propagation algorithm.

1. Goal is still MLE/ maximize cross-entropy.
2. Shared weights (via one convolution filter) → sum over gradient for each use of one filter.
3. Max-pooling → gradient only gets back-propagated through the neuron that "won" the max pool—technically this is a "sub-gradient".