

Convolutional Neural Networks

Vishnu Lokhande

Department of Computer Science and Engineering
University at Buffalo, SUNY
`vishnulo@buffalo.edu`

March 5, 2025

How is convolution defined?

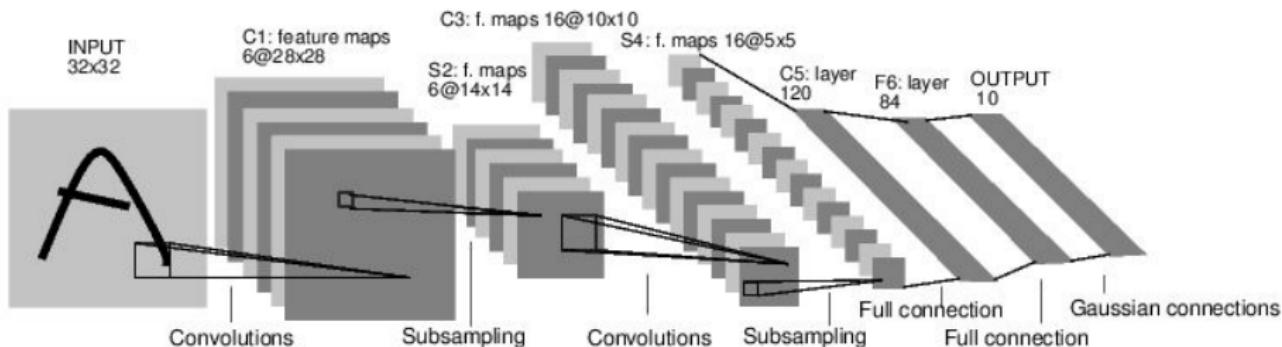
What are some useful properties of convolutional neural networks??

Introduction and History

- ➊ Convolution as a specialized type of linear operation:
 - ▶ Sparse interaction
 - ▶ Parameter sharing
- ➋ Convolutional neural networks (CNN) are simply neural networks that use convolution, instead of general matrix multiplication, in at least one of its layers.

Introduction and History

- ① CNN has been introduced decades ago (1980's) by Lecun:
 - ▶ First applied on classifying handwritten digits, obtained state-of-the-art performance.
 - ▶ Fast development due to the success of computational power acceleration and neural network training algorithms.



Conv filters were 5x5, applied at stride 1

Subsampling (Pooling) layers were 2x2 applied at stride 2
i.e. architecture is [CONV-POOL-CONV-POOL-CONV-FC]

Convolutional Layers¹

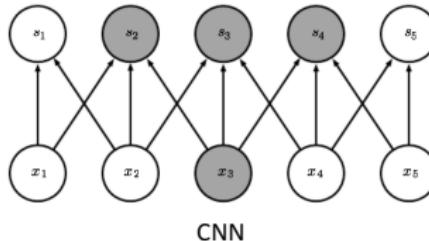
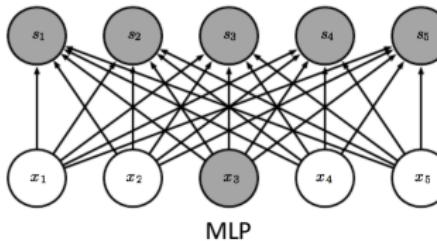
¹ Partially adapted from http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture5.pdf

Convolutional Neural Networks

- ➊ Scale up neural networks to process very large images/video sequences.
 - ▶ Sparse connections.
 - ▶ Parameter sharing.
- ➋ Automatically generalize across spatial translations of inputs.
- ➌ Applicable to any input that is laid out on a grid:
 - ▶ 1-D: signals
 - ▶ 2-D: gray images
 - ▶ 3-D: gray videos
 - ▶ ...

Key Idea

- 1 Replace matrix multiplication in neural nets with convolution:
 - ▶ MLP uses matrix multiplication: with m inputs and n outputs, $m \times n$ parameters and $O(m \times n)$ runtime per example.
 - ▶ CNN uses convolutional with sparse interactions and parameter sharing: with k connections for each input, k parameters and $O(k \times n)$ runtime per example.



- 2 Everything else stays the same:
 - ▶ Maximum likelihood.
 - ▶ Back-propagation.
 - ▶ ...

Convolution on Continuous Domains

1-D convolution:

For two functions $f : \mathbb{R} \mapsto \mathbb{R}$ and $g : \mathbb{R} \mapsto \mathbb{R}$

$$(f * g)(t) \triangleq \int f(\tau)g(t - \tau)d\tau$$

Extension to 2-D convolution:

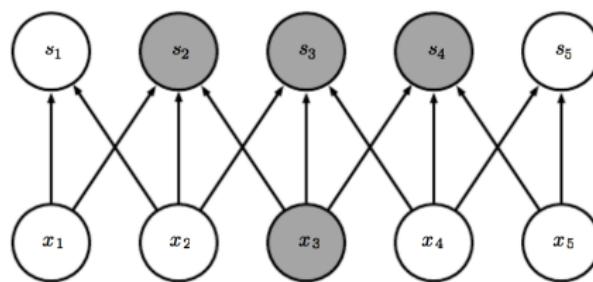
For two functions $f : \mathbb{R} \times \mathbb{R} \mapsto \mathbb{R} \times \mathbb{R}$ and $g : \mathbb{R} \times \mathbb{R} \mapsto \mathbb{R} \times \mathbb{R}$

$$(f * g)(t_1, t_2) \triangleq \int \int f(\tau_1, \tau_2)g(t_1 - \tau_1, t_2 - \tau_2)d\tau_1 d\tau_2$$

Convolution with Discrete Variables

- ① Real data is usually presented in discrete domain.
- ② Let \mathbf{x} be a 1-dimensional input, \mathbf{w} be a 1-dimensional kernel (filter), the output \mathbf{s} is defined via the 1-D convolution:

$$\mathbf{s}(t) \triangleq (\mathbf{x} * \mathbf{w})(t) = \sum_{i=-\infty}^{\infty} \mathbf{x}(i) \mathbf{w}(t - i)$$



- ③ In ML applications, input is a multidimensional array of data, and the kernel is a multidimensional array of parameters to be learned.

Two-dimensional Convolution

- If we use a 2-D image \mathbf{I} as input and use a 2-D kernel \mathbf{K} , we have²

$$\mathbf{S}(i, j) \triangleq (\mathbf{I} * \mathbf{K})(i, j) = \sum_m \sum_n \mathbf{I}(m, n) \mathbf{K}(i - m, j - n)$$

0	1	3
1	2	2
0	0	0

 $*$

0	1	2
2	2	0
0	1	2

 $=$?

²Take the out-of-bound elements to be zero.

Two-dimensional Convolution

- If we use a 2-D image \mathbf{I} as input and use a 2-D kernel \mathbf{K} , we have²

$$\mathbf{S}(i, j) \triangleq (\mathbf{I} * \mathbf{K})(i, j) = \sum_m \sum_n \mathbf{I}(m, n) \mathbf{K}(i - m, j - n)$$

$$\begin{array}{|c|c|c|} \hline 0 & 1 & 3 \\ \hline 1 & 2 & 2 \\ \hline 0 & 0 & 0 \\ \hline \end{array} \quad * \quad \begin{array}{|c|c|c|} \hline 0 & 1 & 2 \\ \hline 2 & 2 & 0 \\ \hline 0 & 1 & 2 \\ \hline \end{array} \quad = \quad \begin{array}{|c|c|c|} \hline 0 & 0 & 1 \\ \hline 0 & 3 & 10 \\ \hline 2 & 6 & 9 \\ \hline \end{array}$$

²Take the out-of-bound elements to be zero.

Commutativity of Convolution

- 1 Convolution is commutative:

$$\mathbf{S}(i, j) = \sum_m \sum_n \mathbf{I}(m, n) \mathbf{K}(i - m, j - n) = \sum_m \sum_n \mathbf{I}(i - m, j - n) \mathbf{K}(m, n)$$

- 2 Commutativity arises because we have flipped the kernel relative to the input:

- ▶ As m increases, index to the input increases, but index to the kernel decreases.

Cross-Correlation

- ① Cross-Correlation: same as convolution, but without flipping the kernel

$$\mathbf{S}(i, j) = \sum_m \sum_n \mathbf{I}(i + m, i + n) \mathbf{K}(m, n)$$

- ② Both referred to as convolution, whether kernel is flipped or not.
- ③ In ML, cross-correlation is preferable because of computational convenience.
- ④ Generally, the two representations are equivalent if the kernel is learned.

Cross-Correlation Convolution Examples

$$\mathbf{S}(i, j) = \sum_m \sum_n \mathbf{I}(i + m, i + n) \mathbf{K}(m, n)$$

3 ₀	3 ₁	2 ₂	1	0
0 ₂	0 ₁	1 ₀	3	1
3 ₀	1 ₁	2 ₂	2	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0
0 ₂	0 ₁	1 ₀	3	1
3 ₀	1 ₁	2 ₂	2	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0 ₃
0 ₀	1 ₂	3 ₀	1 ₀	2 ₁
3 ₁	2 ₀	2 ₂	3 ₂	1 ₀
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0 ₀	0 ₁	1 ₂	3	1
3 ₂	1 ₁	2 ₀	2	3
2 ₀	0 ₁	2 ₂	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0
0 ₀	1 ₁	3 ₀	1 ₀	2 ₁
3 ₁	2 ₂	2 ₀	3 ₀	3
2	0 ₀	1 ₂	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0 ₃
0 ₀	1 ₂	3 ₀	1 ₀	2 ₁
3 ₁	2 ₂	2 ₀	3 ₀	3 ₂
2	0 ₀	1 ₂	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0 ₀	1	3	1	1
3 ₀	1 ₁	2 ₂	2	3
2 ₂	0 ₂	0 ₀	2	2
2 ₀	0 ₁	2 ₂	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0
0 ₀	1 ₁	3 ₀	1 ₀	2 ₁
3 ₁	2 ₂	2 ₀	3 ₀	3 ₂
2	0 ₂	0 ₀	2 ₀	2
2	0 ₀	1 ₂	0 ₂	1

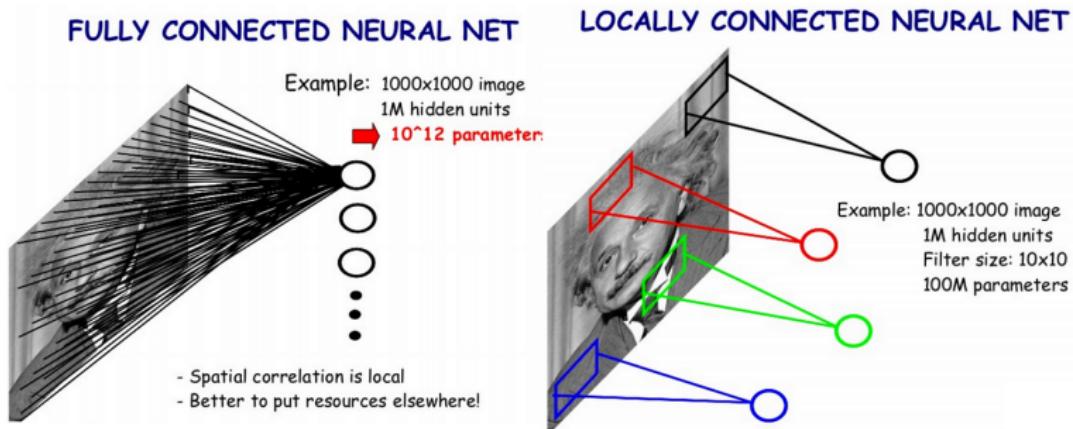
12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0 ₃
0 ₀	1	3	1	1
3 ₁	2 ₀	2 ₁	3 ₂	3 ₂
2	0 ₂	2 ₂	2 ₀	2 ₀
2	0 ₀	0 ₁	0 ₂	1 ₃

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

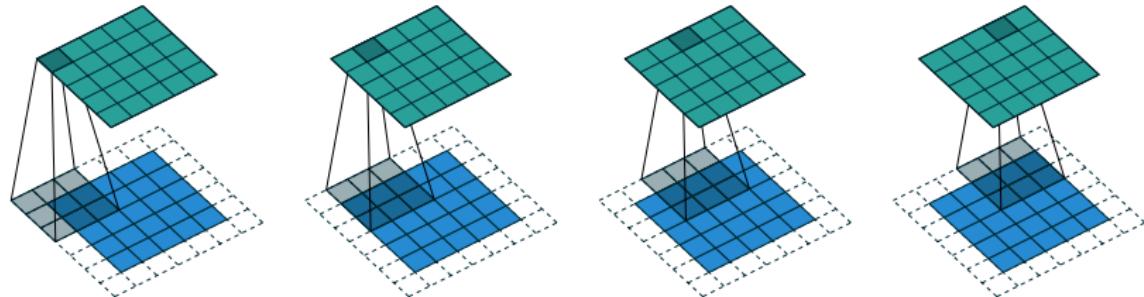
Why do We Want Convolution?

- Parameter sharing.
- Best utilization of data local information.



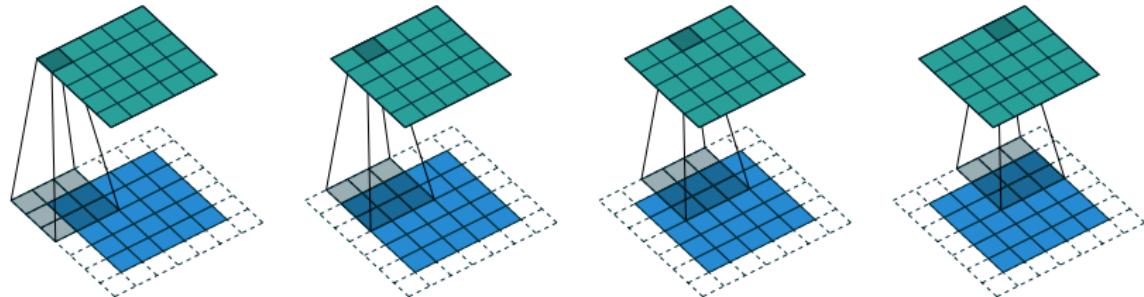
Padding

- Augment the input with zeros on the boundary.
- Padding can make the size of output equal to the size of input.



Padding

- Augment the input with zeros on the boundary.
- Padding can make the size of output equal to the size of input.



General rules for the output sizes and padding? \Rightarrow later

Convolution with Strides

- Convolution with skipped indexes.
- Convolution with stride can be viewed as a way of doing dimensional reduction.

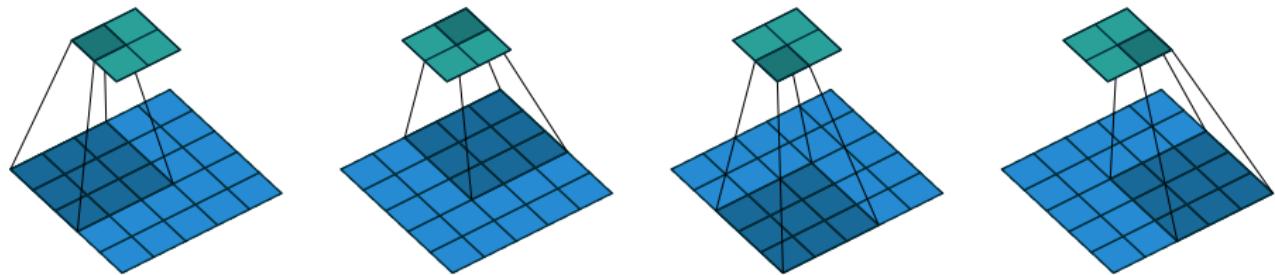


Figure: Convolution with stride of size 2×2 (sometimes simply say "stride 2"). 5×5 input, 2×2 output.

Convolution with Strides

- Convolution with skipped indexes.
- Convolution with stride can be viewed as a way of doing dimensional reduction.

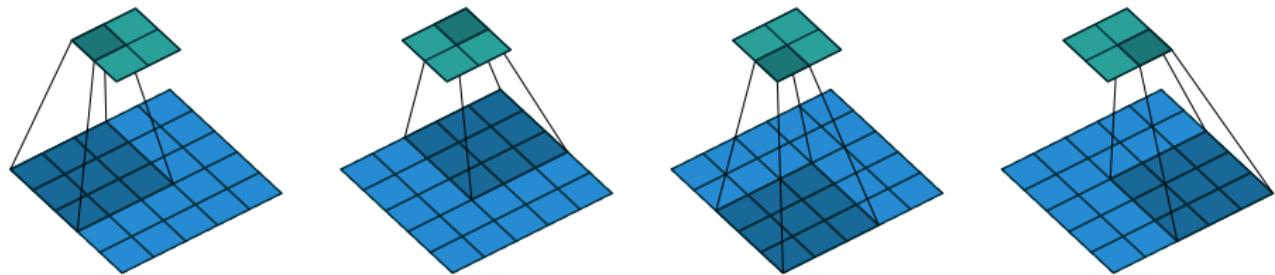


Figure: Convolution with stride of size 2×2 (sometimes simply say "stride 2"). 5×5 input, 2×2 output.

Relation of output size w.r.t. stride?

Convolution with Strides

- Convolution with skipped indexes.
- Convolution with stride can be viewed as a way of doing dimensional reduction.

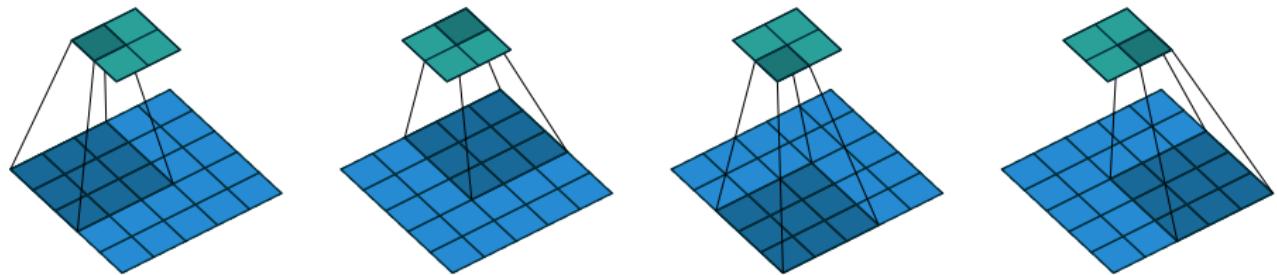


Figure: Convolution with stride of size 2×2 (sometimes simply say "stride 2"). 5×5 input, 2×2 output.

Relation of output size w.r.t. stride

$$\text{output size} = (\text{input size} - \text{filter size}) / \text{stride} + 1$$

Convolution with Strides

- Convolution with skipped indexes.
- Convolution with stride can be viewed as a way of doing dimensional reduction.

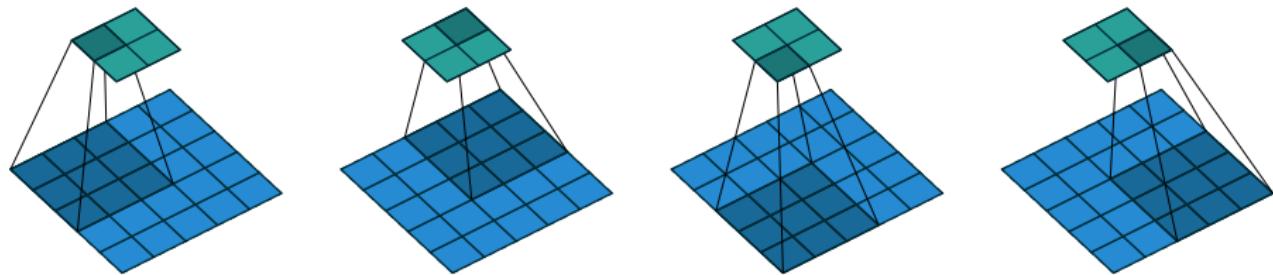


Figure: Convolution with stride of size 2×2 (sometimes simply say "stride 2"). 5×5 input, 2×2 output.

Relation of output size w.r.t. stride

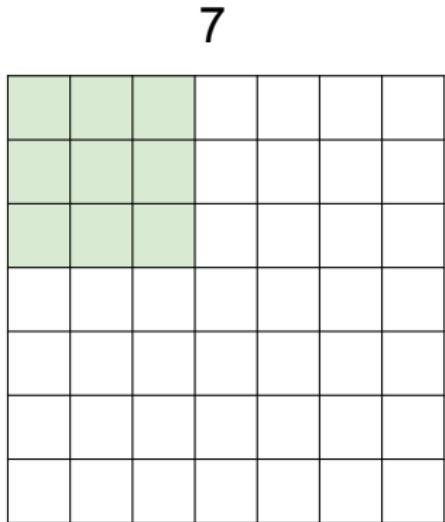
$$\text{output size} = (\text{input size} - \text{filter size}) / \text{stride} + 1$$

When considering padding of size P :

$$\text{output size} = (\text{input size} + 2P - \text{filter size}) / \text{stride} + 1$$

A Closer Look at Strides

- Cannot use arbitrary stride sizes.



7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

doesn't fit!
cannot apply 3x3 filter on
7x7 input with stride 3.

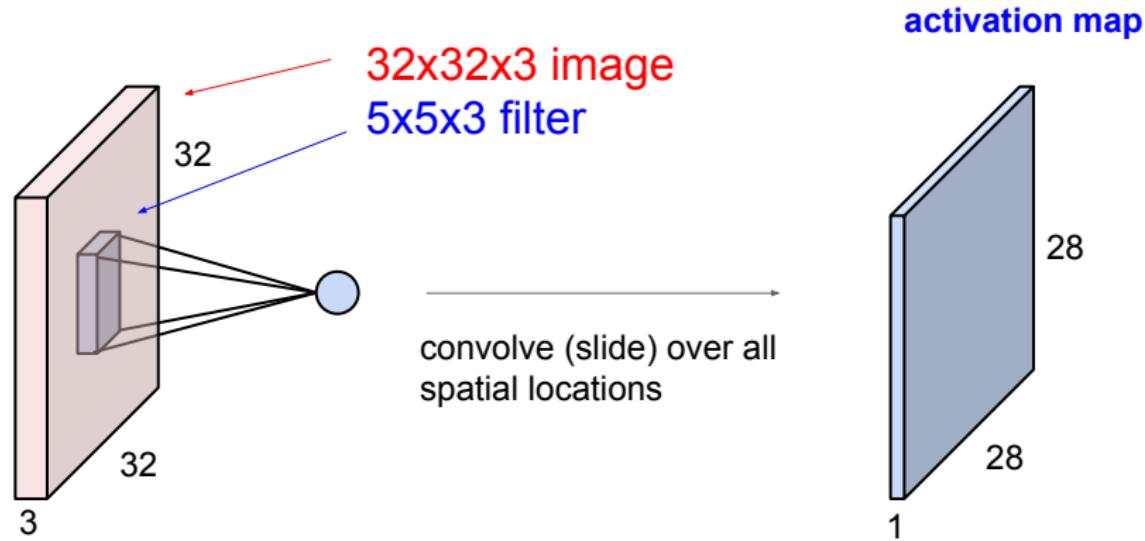
Practical Setting

- In practice, common to set convolutional layers with stride 1, filters of size $F \times F$, and zero-padding with $(F - 1)/2$:
 - ▶ will preserve input size

0	0	0	0	0	0			
0								
0								
0								
0								

Extension to 3-dimensional tensors

Convolution Layer

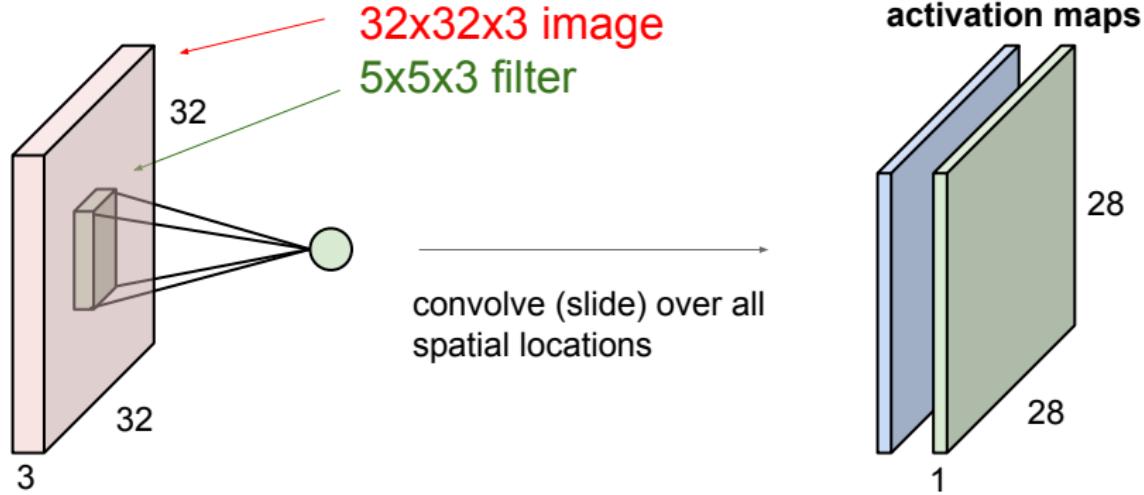


- Convolution and sum over the third dimension.

Extension to 3-dimensional tensors

Convolution Layer

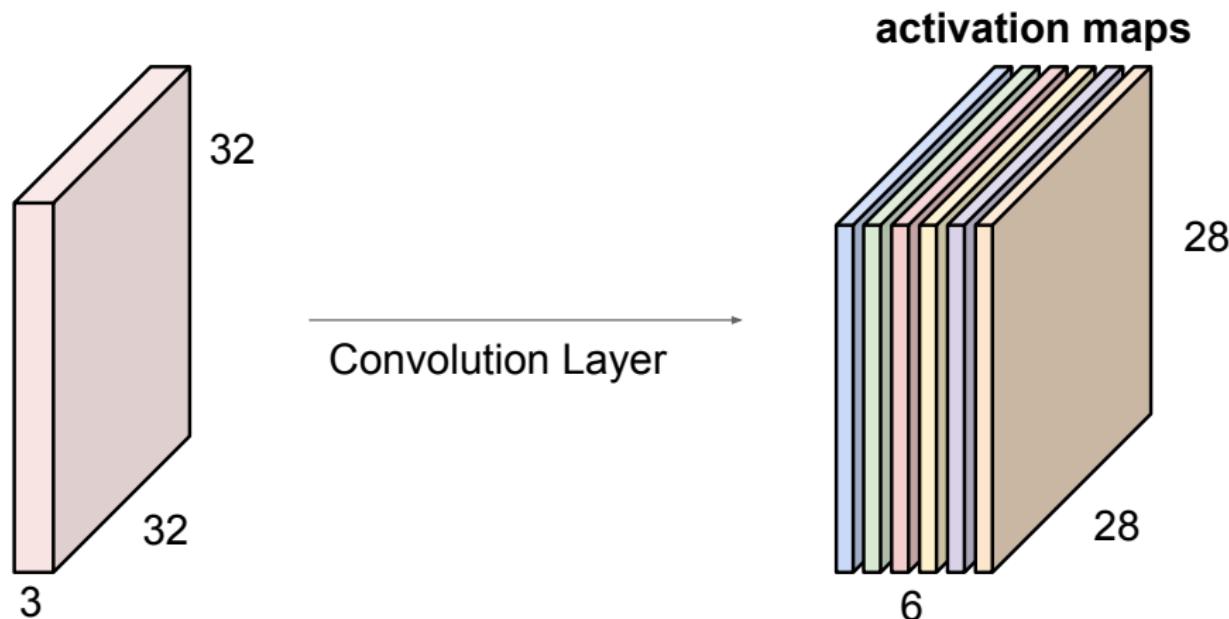
consider a second, green filter



- Convolution and sum over the third dimension.

Extension to 3-dimensional tensors

If we had six $5 \times 5 \times 3$ filters^a, we'll get 6 separate activation maps.



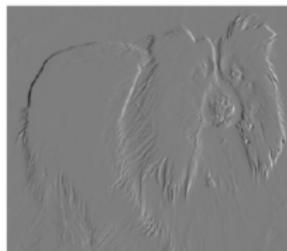
We stack these up to get a “new image” of size 28x28x6!

^aWe sometimes ignore the last dimension for simplicity, written as filter size of 5×5 .

Example: Efficiency of Convolution for Edge Detection

- ➊ Image on right formed by taking each pixel of input image and subtracting the value of its neighboring pixel on the left:
 - ▶ this is a measure of all the vertically oriented edges in input image, useful for object detection.
 - ▶ can be implemented as a convolution with a 1×2 kernel:
 $K = (1, -1)$.

Input
image



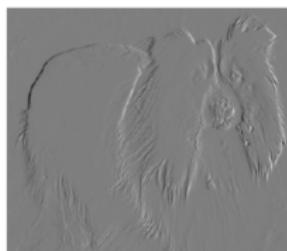
Both images are 280 pixels tall
Input image is 320 pixels wide
Output image is 319 pixels wide

- ➋ Number of operations:

Example: Efficiency of Convolution for Edge Detection

- ➊ Image on right formed by taking each pixel of input image and subtracting the value of its neighboring pixel on the left:
 - ▶ this is a measure of all the vertically oriented edges in input image, useful for object detection.
 - ▶ can be implemented as a convolution with a 1×2 kernel:
 $K = (1, -1)$.

Input
image



Both images are 280 pixels tall
Input image is 320 pixels wide
Output image is 319 pixels wide

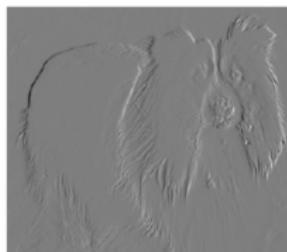
- ➋ Number of operations:

- ▶ Convolution with one kernel of size 1×2 :
- ▶ MLP:

Example: Efficiency of Convolution for Edge Detection

- ➊ Image on right formed by taking each pixel of input image and subtracting the value of its neighboring pixel on the left:
 - ▶ this is a measure of all the vertically oriented edges in input image, useful for object detection.
 - ▶ can be implemented as a convolution with a 1×2 kernel:
 $K = (1, -1)$.

Input
image



Both images are 280 pixels tall
Input image is 320 pixels wide
Output image is 319 pixels wide

- ➋ Number of operations:

- ▶ Convolution with one kernel of size 1×2 : $319 \times 280 \times 1 \times 2$.
- ▶ MLP: $320 \times 280 \times 319 \times 280$.

Output Volume Size w.r.t. Padding and Stride

- Input volume: 1×1 , 5×5 filter with stride 1, pad 2.
- Output volume size?
 - ▶ 1×1

Output Volume Size w.r.t. Padding and Stride

- Input volume: 1×1 , 5×5 filter with stride 1, pad 2.
- Output volume size?
 - ▶ 1×1

Output Volume Size w.r.t. Padding and Stride

- Input volume: 2×2 , 5×5 filter with stride 1, pad 2.
- Output volume size?
 - ▶ 2×2

Output Volume Size w.r.t. Padding and Stride

- Input volume: 2×2 , 5×5 filter with stride 1, pad 2.
- Output volume size?
 - ▶ 2×2

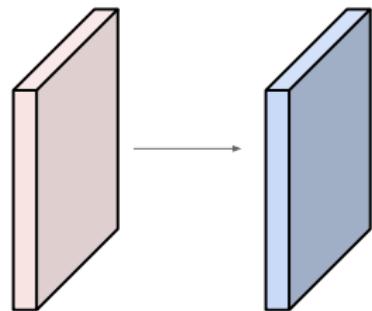
Output Volume Size w.r.t. Padding and Stride

Examples

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

Output volume size: ?

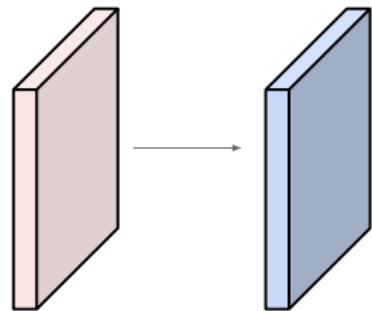


Output Volume Size w.r.t. Padding and Stride

Examples

Input volume: **32x32x3**

10 **5x5** filters with stride **1**, pad **2**



Output volume size:

$(32+2*2-5)/1+1 = 32$ spatially, so

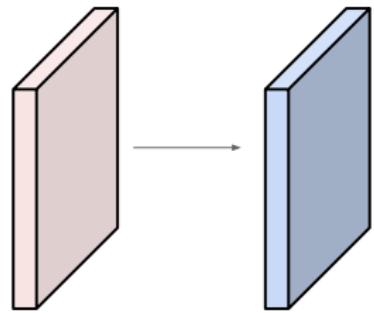
32x32x10

Output Volume Size w.r.t. Padding and Stride

Examples

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2



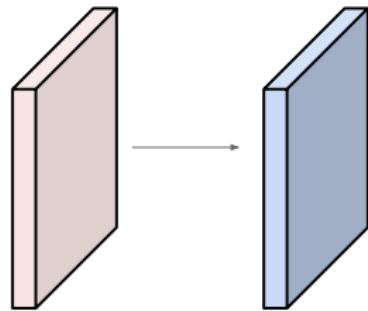
Number of parameters in this layer?

Output Volume Size w.r.t. Padding and Stride

Examples

Input volume: **32x32x3**

10 **5x5** filters with stride 1, pad 2



Number of parameters in this layer?

each filter has $5*5*3 + 1 = 76$ params (+1 for bias)

$$\Rightarrow 76 * 10 = 760$$

Summary

To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$.
- Requires four hyperparameters:
 - Number of filters K .
 - Filter size F .
 - Stride S .
 - Amount of zero-padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$:

$$W_2 = (W_1 - F + 2P)/S + 1$$

$$H_2 = (H_1 - F + 2P)/S + 1$$

$$D_2 = K$$

- With parameter sharing, it introduces $F^2 D_1$ weights per filter, for a total of $F^2 D_1 K$ weights and K biases.
- In the output volume, the d -th depth slices (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input with a stride of S , and then offset by d -th bias.

Summary

To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$.
- Requires four hyperparameters:
 - Number of filters K .
 - Filter size F .
 - Stride S .
 - Amount of zero-padding P .

Common settings:

$K = (\text{powers of 2, e.g. } 32, 64, 128, 512)$

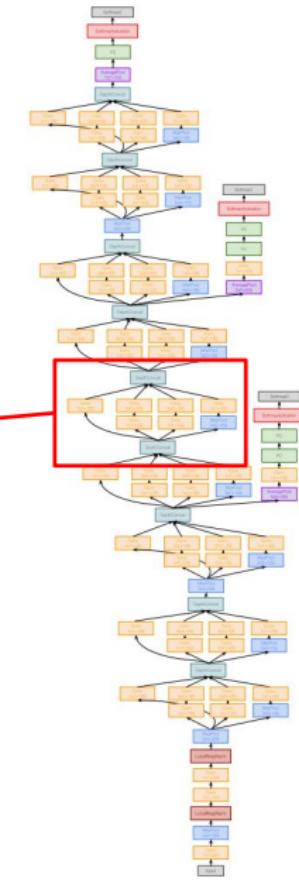
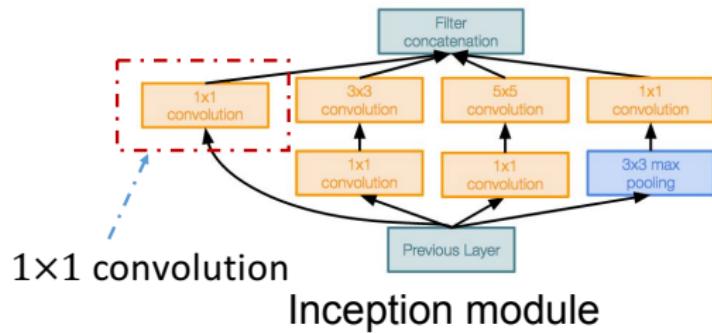
- $F = 3, S = 1, P = 1$
- $F = 5, S = 1, P = 2$
- $F = 5, S = 2, P = ? \text{ (whatever fits)}$
- $F = 1, S = 1, P = 0$

The Power of 1×1 Convolution³

³Partially adapted from http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture9.pdf

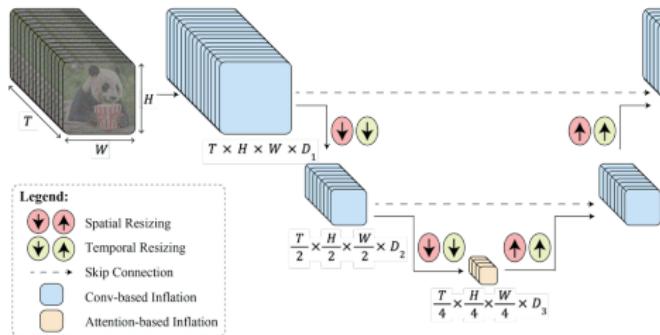
GoogleNet

Stack the module on top of each other!

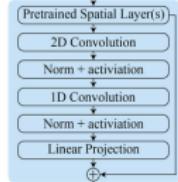


Text-to-Video Generation

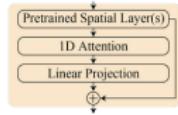
(a) Space-Time UNet (STUNet)



(b) Convolution-based Inflation Block



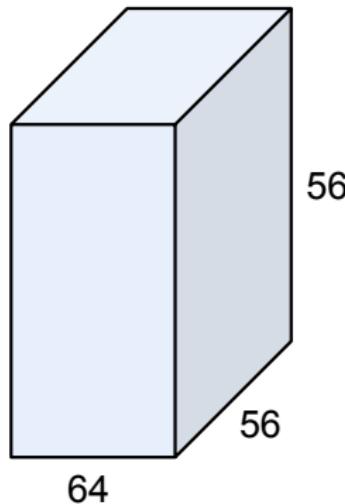
(c) Attention-based Inflation Block



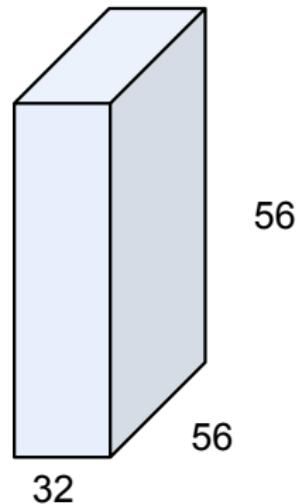
1x1 Convolution Layers Make Perfect Sense

Used in GoogleNet (inception model)

Perform like dimension reduction/extension on the third dimension.



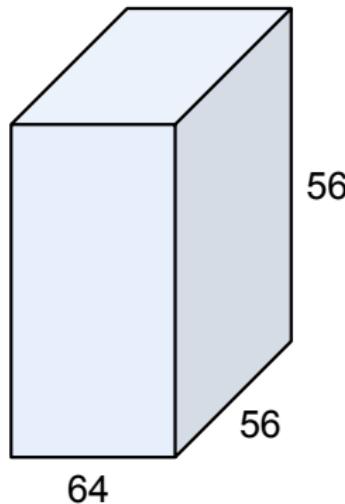
1x1 CONV
with 32 filters
→
(each filter has size
 $1 \times 1 \times 64$, and performs a
64-dimensional dot product)



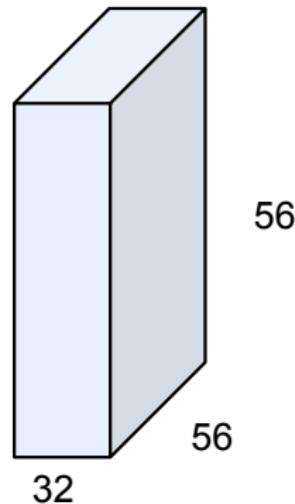
1x1 Convolution Layers Make Perfect Sense

Used in GoogleNet (inception model)

Perform like dimension reduction/extension on the third dimension.



1x1 CONV
with 32 filters
→
(each filter has size
 $1 \times 1 \times 64$, and performs a
64-dimensional dot product)



What happens with the 1×1 convolution here?

The Power of Small Filters

Why do we need small filters?

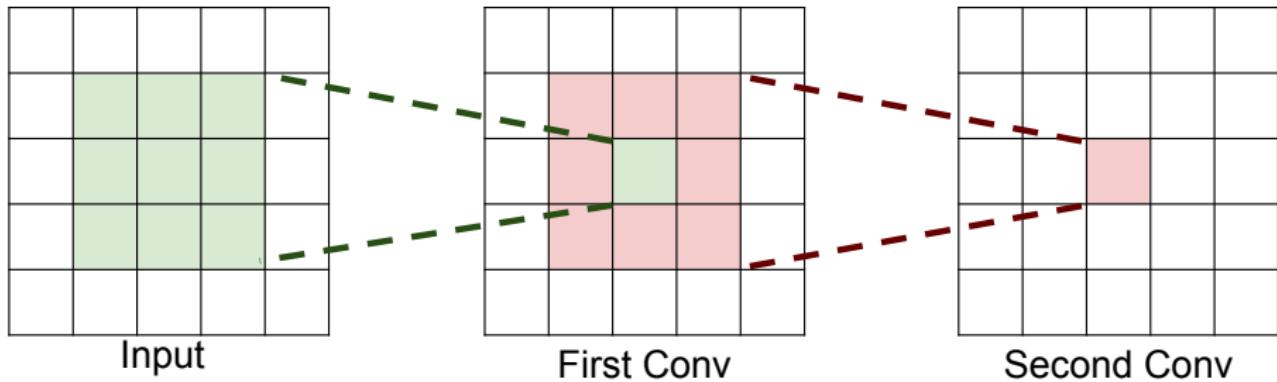
Proof Idea

- Define two kinds of CNNs which have different filter sizes such that they have the same representation power.
- Show the CNN with smaller filter size has less parameters.

The Power of Small Filters

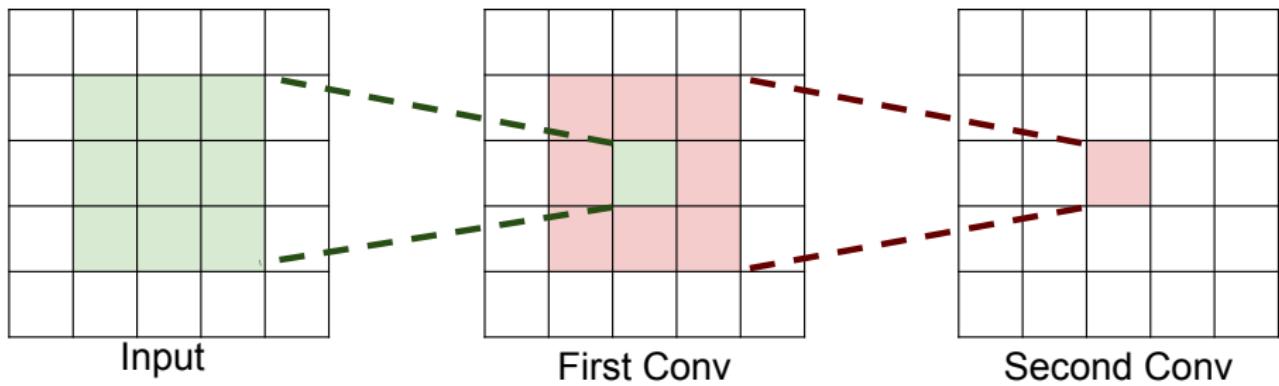
Suppose we stack two 3×3 (filter size) conv layers (stride 1):

- each neuron sees 3×3 region of previous activation map.



The Power of Small Filters

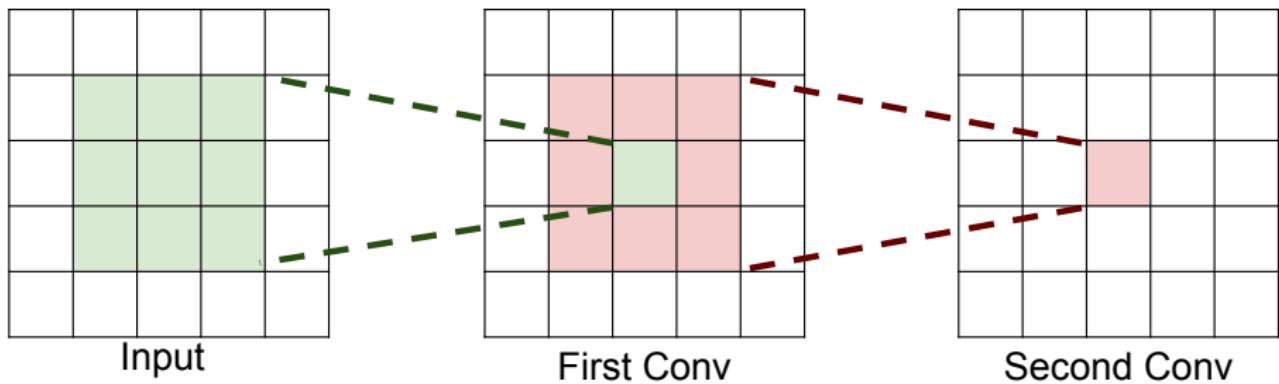
Question: How big of a region in the input does a neuron on the second conv layer see?



The Power of Small Filters

Question: How big of a region in the input does a neuron on the second conv layer see?

Answer: 5×5

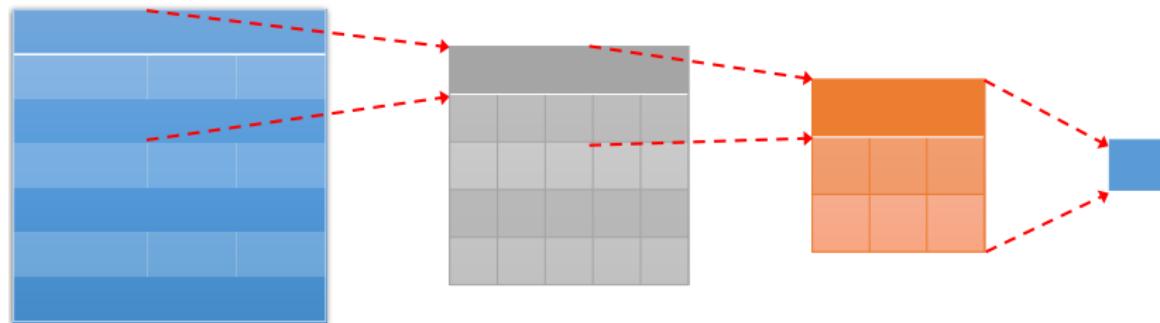


The Power of Small Filters

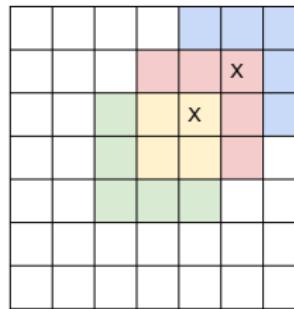
Question: If we stack **three** 3x3 conv layers, how big of an input region does a neuron in the third layer see?

The Power of Small Filters

Question: If we stack **three** 3x3 conv layers, how big of an input region does a neuron in the third layer see?



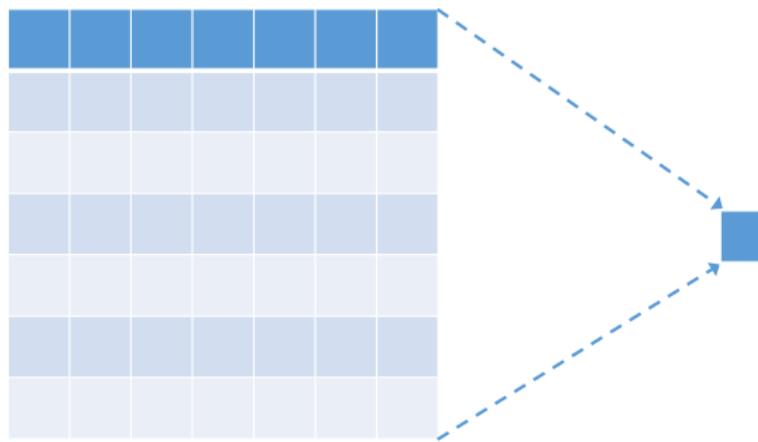
Answer: 7×7



The Power of Small Filters

Question: If we stack **three** 3x3 conv layers, how big of an input region does a neuron in the third layer see?

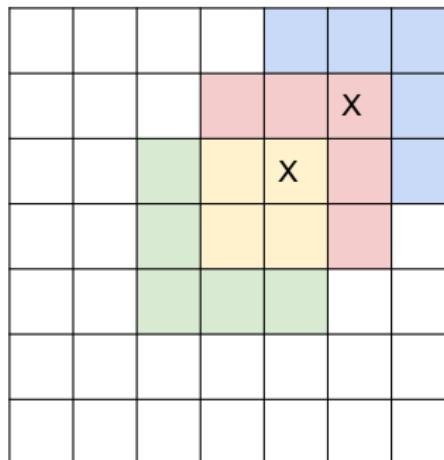
- If we use 1 convolution layer, we need to use a kernel of 7×7 in order to let a neuron in the output layer see a 7×7 region in the input.



The Power of Small Filters

Question: If we stack **three** 3x3 conv layers, how big of an input region does a neuron in the third layer see?

Answer: 7 x 7



Three 3×3 conv
gives similar
representational
power as a single
 7×7 convolution

While it has less parameters!

The Power of Small Filters

Suppose input is $H \times W \times C$ and we use convolutions with C filters to preserve depth (stride 1, padding to preserve H, W).

- One CONV with 7×7 filters.
- Number of weights:

$$= C \times (7 \times 7 \times C) = 49C^2$$

- Number of multiply-adds:

$$= (H \times W \times C) \times (7 \times 7 \times C) = 49HWC^2$$

- Three CONV with 3×3 filters.
- Number of weights:

$$= 3 \times C \times (3 \times 3 \times C) = 27C^2$$

fewer parameters = GOOD

- Number of multiply-adds:

$$= 3 \times (H \times W \times C) \times (3 \times 3 \times C) = 27HWC^2$$

less compute = GOOD

The Power of Small Filters

Suppose input is $H \times W \times C$ and we use convolutions with C filters to preserve depth (stride 1, padding to preserve H, W).

- One CONV with 7×7 filters.
- Number of weights:

$$= C \times (7 \times 7 \times C) = 49C^2$$

- Number of multiply-adds:

$$= (H \times W \times C) \times (7 \times 7 \times C) = 49HWC^2$$

- Three CONV with 3×3 filters.
- Number of weights:

$$= 3 \times C \times (3 \times 3 \times C) = 27C^2$$

fewer parameters = GOOD

- Number of multiply-adds:

$$= 3 \times (H \times W \times C) \times (3 \times 3 \times C) = 27HWC^2$$

less compute = GOOD

The Power of Small Filters

Suppose input is $H \times W \times C$ and we use convolutions with C filters to preserve depth (stride 1, padding to preserve H, W).

- One CONV with 7×7 filters.
- Number of weights:

$$= C \times (7 \times 7 \times C) = 49C^2$$

- Number of multiply-adds:

$$= (H \times W \times C) \times (7 \times 7 \times C) = 49HWC^2$$

- Three CONV with 3×3 filters.
- Number of weights:

$$= 3 \times C \times (3 \times 3 \times C) = 27C^2$$

fewer parameters = GOOD

- Number of multiply-adds:

$$= 3 \times (H \times W \times C) \times (3 \times 3 \times C) = 27HWC^2$$

less compute = GOOD

The Power of Small Filters

Suppose input is $H \times W \times C$ and we use convolutions with C filters to preserve depth (stride 1, padding to preserve H, W).

- One CONV with 7×7 filters.
- Number of weights:

$$= C \times (7 \times 7 \times C) = 49C^2$$

- Number of multiply-adds:

$$= (H \times W \times C) \times (7 \times 7 \times C) = 49HWC^2$$

- Three CONV with 3×3 filters.
- Number of weights:

$$= 3 \times C \times (3 \times 3 \times C) = 27C^2$$

fewer parameters = GOOD

- Number of multiply-adds:

$$= 3 \times (H \times W \times C) \times (3 \times 3 \times C) = 27HWC^2$$

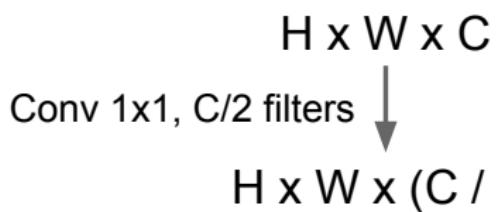
less compute = GOOD

The Power of Small Filters

Why stop at 3×3 filters? Why not try 1×1 ?

The Power of Small Filters

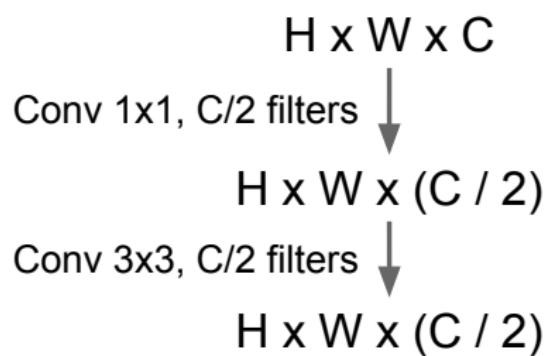
Why stop at 3×3 filters? Why not try 1×1 ?



1. “bottleneck” 1×1 conv to reduce dimension

The Power of Small Filters

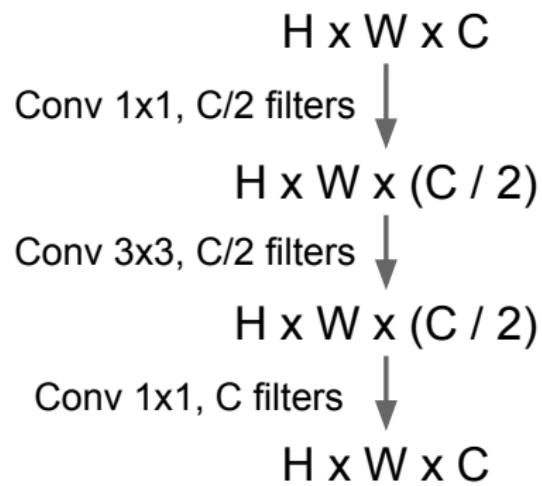
Why stop at 3×3 filters? Why not try 1×1 ?



1. “bottleneck” 1×1 conv to reduce dimension
2. 3×3 conv at reduced dimension

The Power of Small Filters

Why stop at 3×3 filters? Why not try 1×1 ?

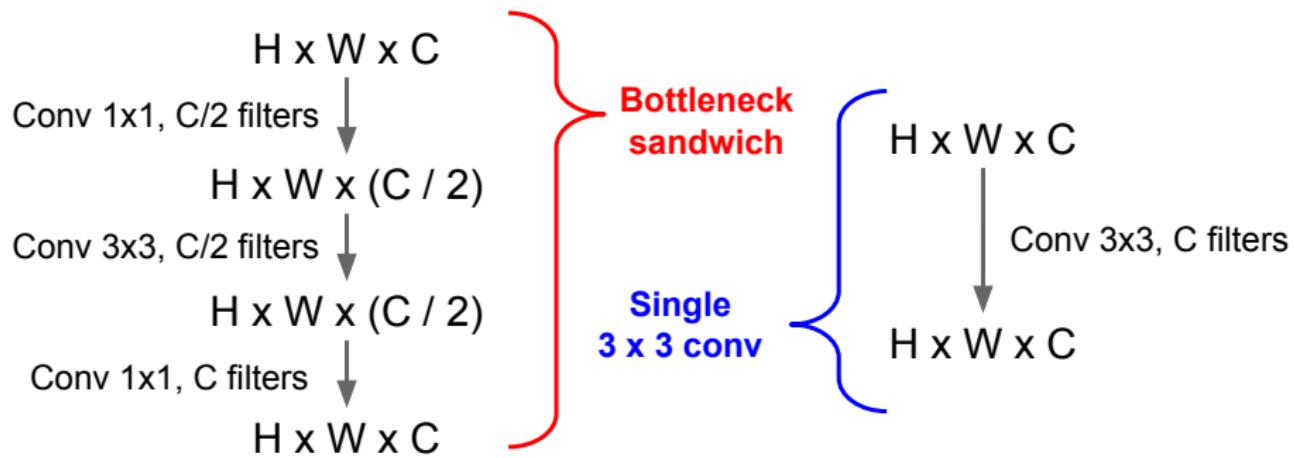


1. “bottleneck” 1×1 conv to reduce dimension
2. 3×3 conv at reduced dimension
3. Restore dimension with another 1×1 conv

[Seen in Lin et al, “Network in Network”, GoogLeNet, ResNet]

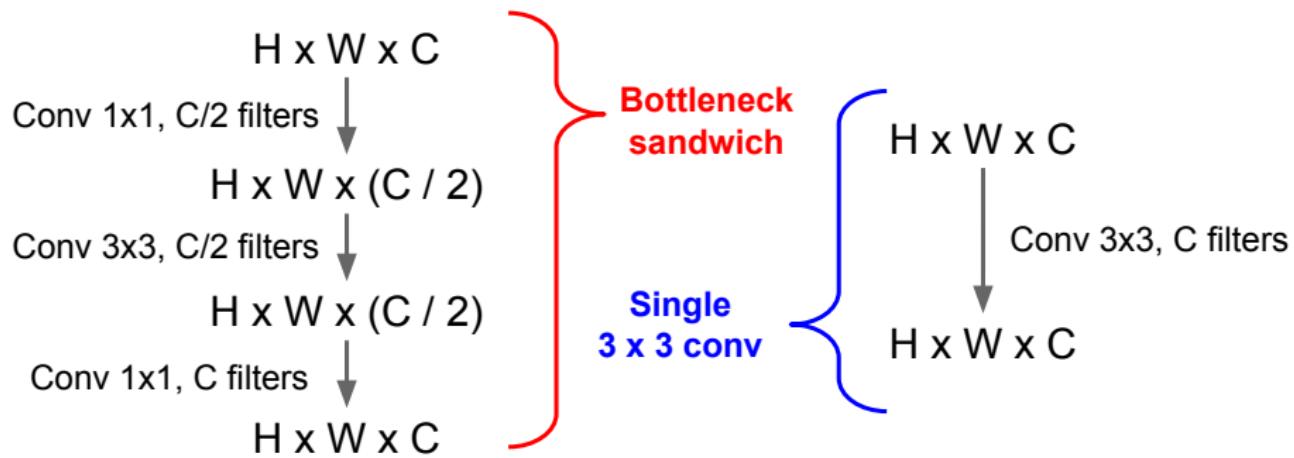
The Power of Small Filters

Why stop at 3×3 filters? Why not try 1×1 ?



The Power of Small Filters

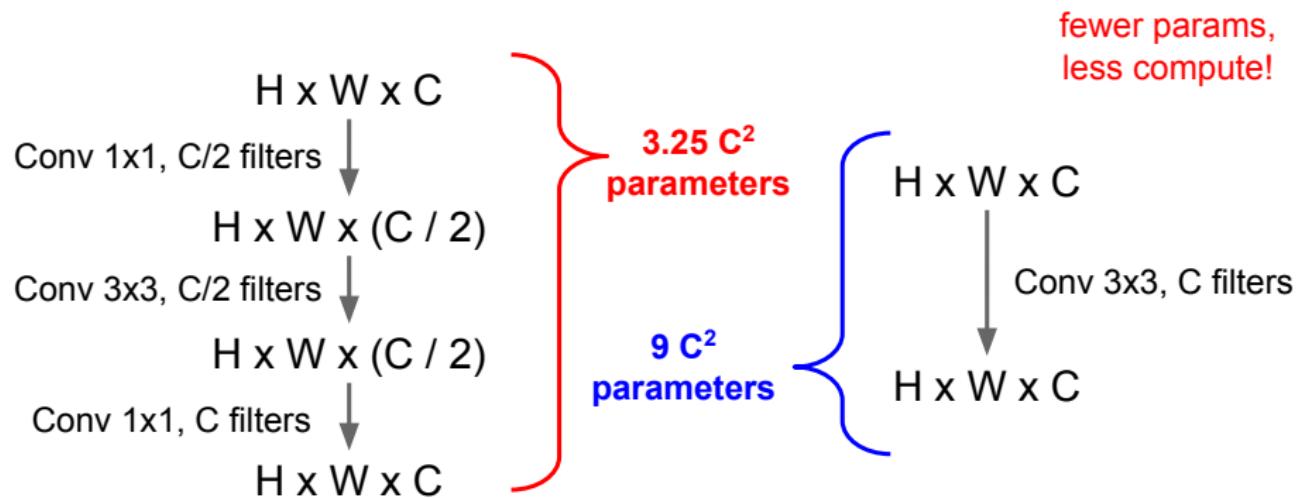
Why stop at 3×3 filters? Why not try 1×1 ?



How many parameters?

The Power of Small Filters

Why stop at 3×3 filters? Why not try 1×1 ?



The Power of Small Filters

Still using 3×3 filters . . . can we break it up?

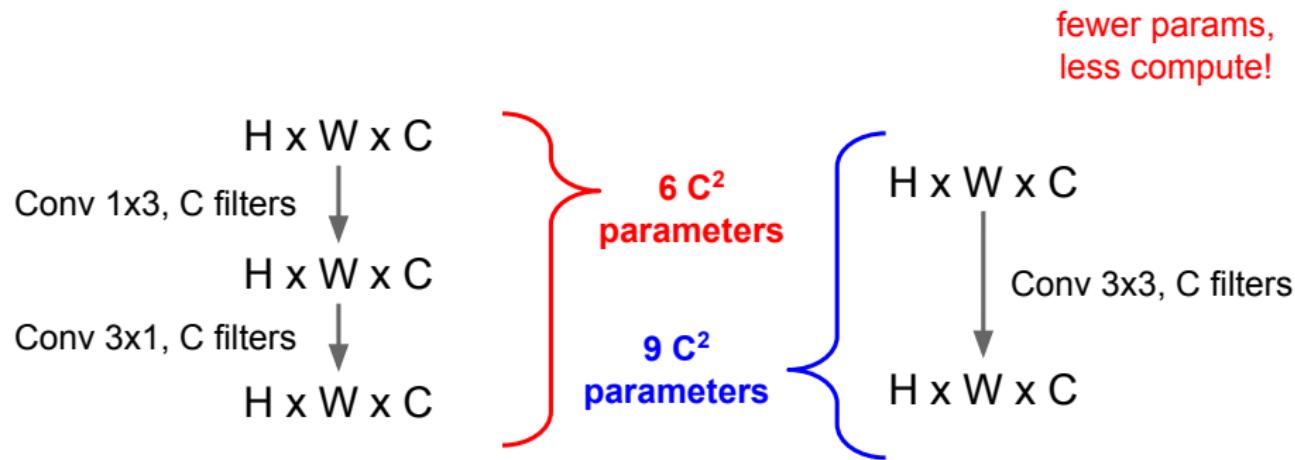
The Power of Small Filters

Still using 3×3 filters ... can we break it up?



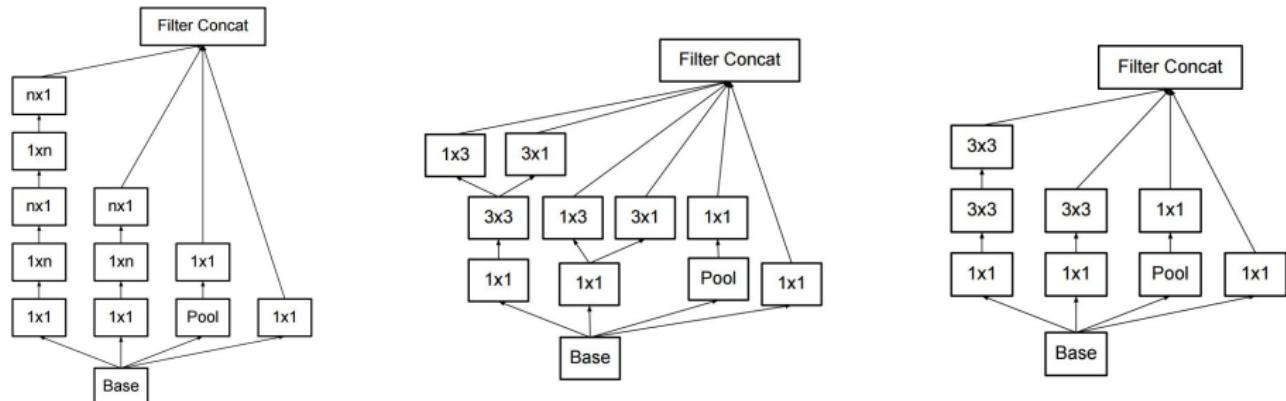
The Power of Small Filters

Still using 3×3 filters ... can we break it up?



The Power of Small Filters

Latest version of GoogLeNet incorporates all these ideas



Szegedy et al, "Rethinking the Inception Architecture for Computer Vision"

How to Stack Convolutions: Recap

- ① Replace large convolutions (5×5 , 7×7) with stacks of 3×3 convolutions.
- ② 1×1 “bottleneck” convolutions are very efficient.
- ③ Can factor $N \times N$ convolutions into $1 \times N$ and $N \times 1$.
- ④ All of the above give fewer parameters and less compute.

Some Properties of Convolution

Equivalence of Convolution to Translation

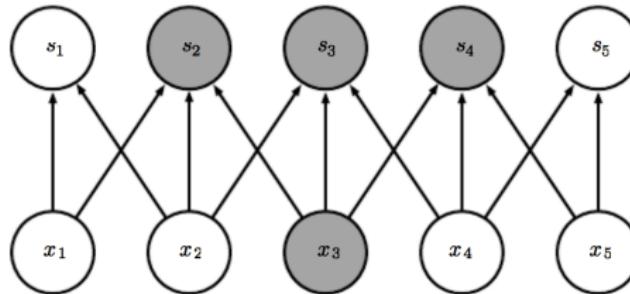
- ➊ The particular form of parameter sharing leads to equivalence to translation:
 - ▶ meaning that if the input changes, the output changes in the same way.
- ➋ If g is a function that translates the input, *i.e.*, that shifts it, then the convolution function is equivalent to g :
 - ▶ $\mathbf{I}(x, y)$ is image brightness at point (x, y) .
 - ▶ $\mathbf{I}'(x, y) = g(\mathbf{I}) = \mathbf{I}(x - 1, y)$, *i.e.*, shifts every pixel of \mathbf{I} one unit to the right.
 - ▶ If we apply g to \mathbf{I} and then apply convolution, the output will be the same as if we applied convolution to \mathbf{I}' , then applied transformation g to the output.
- ➌ Convolution is equivalent to translation.

Question

Is convolutional operator linear or nonlinear?

Convolution as a Structured Linear Operator

- Convolution can be viewed as multiplication by a matrix, with some constraints:
 - Construct a matrix \mathbf{A} , such that the following 1-D convolution satisfies: $\mathbf{s} = \mathbf{x} * \mathbf{w} = \mathbf{A}\mathbf{x}$.

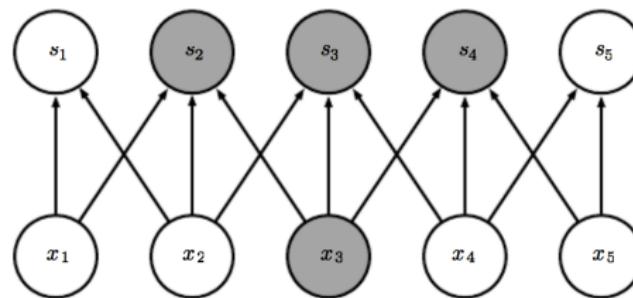


Convolution as a Structured Linear Operator

- Convolution can be viewed as multiplication by a matrix, with some constraints:
 - Construct a matrix \mathbf{A} , such that the following 1-D convolution satisfies: $\mathbf{s} = \mathbf{x} * \mathbf{w} = \mathbf{A} \mathbf{x}$.

$$\mathbf{s} = \underbrace{\begin{bmatrix} w_2 & w_3 & 0 & 0 & 0 \\ w_1 & w_2 & w_3 & 0 & 0 \\ 0 & w_1 & w_2 & w_3 & 0 \\ 0 & 0 & w_1 & w_2 & w_3 \\ 0 & 0 & 0 & w_1 & w_2 \end{bmatrix}}_{\mathbf{A}}$$

- \mathbf{A} is called the univariate Toeplitz matrix.



Convolution as a Structured Linear Operator

- For a 2-D convolution $\mathbf{S} = \mathbf{X} * \mathbf{W}$ mat(A vec(X))
 - “vec” means vectorizing a matrix into a vector row by row; “mat” means reshaping a vector into a matrix of the original size.

$$\underbrace{\mathbf{S}}_{2 \times 2} = \underbrace{\begin{bmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{bmatrix}}_{\mathbf{X}:3 \times 3} * \underbrace{\begin{bmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \end{bmatrix}}_{\mathbf{W}:2 \times 2}$$

Convolution as a Structured Linear Operator

- For 2-D convolution $\mathbf{S} = \mathbf{X} * \mathbf{W} = \text{mat}(\mathbf{A} \text{vec}(\mathbf{X}))$ (\mathbf{A} is a doubly block circulant matrix):

$$\mathbf{A} = \begin{bmatrix} W_{11} & W_{12} & 0 & W_{21} & W_{22} & 0 & 0 & 0 & 0 \\ 0 & W_{11} & W_{12} & 0 & W_{21} & W_{22} & 0 & 0 & 0 \\ 0 & 0 & 0 & W_{11} & W_{12} & 0 & W_{21} & W_{22} & 0 \\ 0 & 0 & 0 & 0 & W_{11} & W_{12} & 0 & W_{21} & W_{22} \end{bmatrix}$$

$$\underbrace{\mathbf{S}_{2 \times 2}}_{= \underbrace{\begin{bmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{bmatrix}}_{\mathbf{X}:3 \times 3} * \underbrace{\begin{bmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \end{bmatrix}}_{\mathbf{W}:2 \times 2}}$$

Transposed Convolution

- ① Recalled $\mathbf{S} = \mathbf{X} * \mathbf{W} = \text{mat}(\mathbf{A} \text{vec}(\mathbf{X}))$.
- ② We can go from a smaller size matrix (e.g. \mathbf{S}) to a larger size matrix (e.g. \mathbf{X}) by using the transpose of \mathbf{A} :

$$\mathbf{X}' = \text{mat}(\mathbf{A}^T \text{vec}(\mathbf{S}))$$

- ③ Sometimes referred to “Deconvolution”⁴.

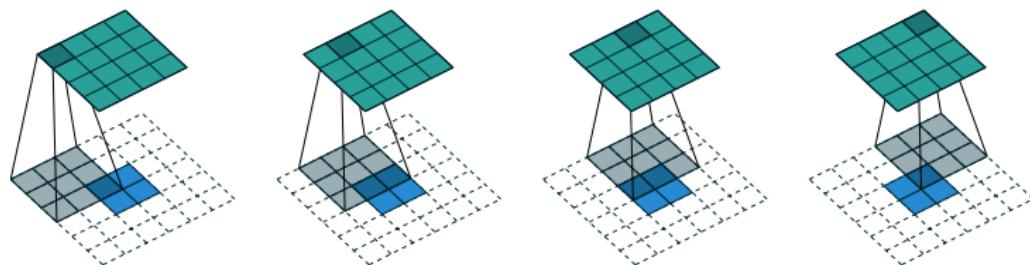
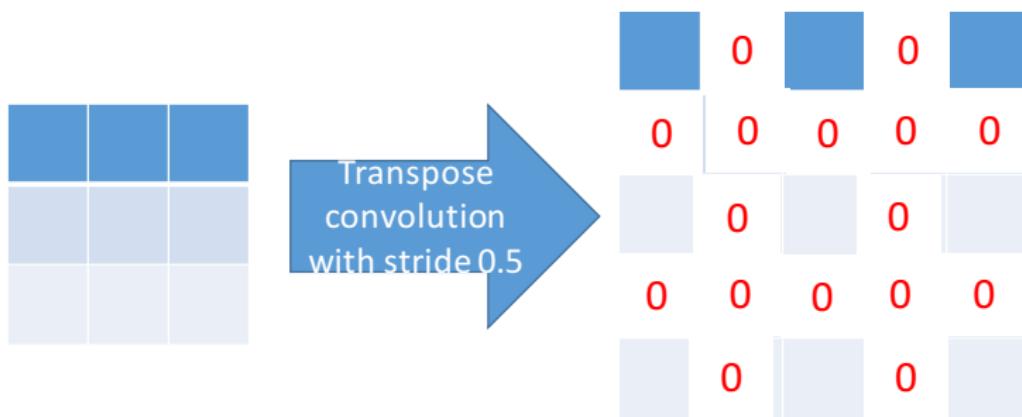


Figure: The transpose of convolving performs like an inversion of the standard convolution.

⁴There are several definitions of deconvolution. This is one of the mostly references.

Stride in Transposed Convolution

- ① In convolution, the stride is defined as an integer $s \geq 1$.
- ② In transposed convolution, the stride is a fractional number, i.e., $s \in (0, 1]$.
- ③ Equivalent to adding $(\frac{1}{s} - 1)$ zeros between each two adjacent inputs.



Transposed Convolution

Relationship between Convolution and Transpose Convolution

A convolution with kernel size $k \times k$, stride s and zero-padding p and whose input size $i \times i$ is such that $i + 2p - k$ is a multiple of s has an associated transposed convolution described by $k' = k$, s' and p' and o' . Then the output size of deconvolution can be expressed as

Transposed Convolution

Relationship between Convolution and Transpose Convolution

A convolution with kernel size $k \times k$, stride s and zero-padding p and whose input size $i \times i$ is such that $i + 2p - k$ is a multiple of s has an associated transposed convolution described by $k' = k$, s' and p' and o' . Then the output size of deconvolution can be expressed as

$$o' = \frac{i + 2p - k}{ss'} + \frac{2p' - k + 1}{s'} + 1.$$

Basic Gradient Computation

Notation

- \mathcal{L} : the loss function.
- $\mathbf{H}^{(1)}$: first hidden units. $H_{cij}^{(1)}$ indexes position c (typically a 2-dimensional index) within feature map i for example j .
- $\mathbf{H}^{(2)}$: second hidden units. $H_{cij}^{(2)}$ indexes position c (typically a 2-dimensional index) within feature map i for example j .
- \mathbf{V} : visible units with the same index-format as \mathbf{H} .
- $W^{(1)}$: weights defining the kernel for the first layer. $W_{cij}^{(2)}$ indexes the weight at position c within the kernel, connecting visible channel i to hidden channel j .
- $W^{(2)}$: weights defining the kernel for the second layer.
- Assume all strides to be one, and ignore the biases.

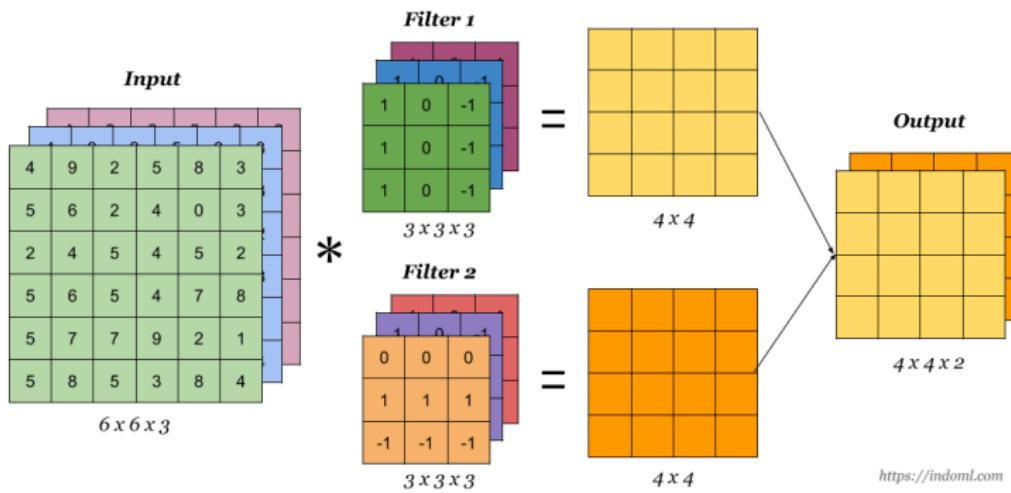
$$H_{cij}^{(1)} = \sum_{k,m} W_{kmi}^{(1)} V_{c+k, m, j}$$

$$H_{cij}^{(2)} = \sum_{k,m} W_{kmi}^{(2)} H_{c+k, m, j}^{(1)}$$

Notation

$$H_{cij}^{(1)} = \sum_{k,m} W_{kmi}^{(1)} V_{c+k,m,j}$$

$$H_{cij}^{(2)} = \sum_{k,m} W_{kmi}^{(2)} H_{c+k,m,j}^{(1)}$$



<https://indoml.com>

Basic Gradients

In order to apply BP:

What is the gradients of $\frac{\partial \mathcal{L}}{\partial H_{cij}^{(i)}}$ and $\frac{\partial \mathcal{L}}{\partial W_{cij}^{(i)}}$?

Basic Gradients

The term $\frac{\partial \mathcal{L}}{\partial H_{kjn}^{(1)}}$ is calculated as

$$\frac{\partial \mathcal{L}}{\partial H_{cij}^{(1)}} = \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(2)}} \frac{\partial H_{kmn}^{(2)}}{\partial H_{cij}^{(1)}}$$

Basic Gradients

The term $\frac{\partial \mathcal{L}}{\partial H_{kjn}^{(1)}}$ is calculated as

$$\frac{\partial \mathcal{L}}{\partial H_{cij}^{(1)}} = \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(2)}} \frac{\partial H_{kmn}^{(2)}}{\partial H_{cij}^{(1)}}$$

Basic Gradients

The term $\frac{\partial \mathcal{L}}{\partial H_{kij}^{(1)}}$ is calculated as

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial H_{cij}^{(1)}} &= \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(2)}} \frac{\partial H_{kmn}^{(2)}}{\partial H_{cij}^{(1)}} \\ &= \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(2)}} \frac{\partial \sum_{p,q} W_{pqm}^{(2)} H_{k+p,q,n}^{(1)}}{\partial H_{cij}^{(1)}}\end{aligned}$$

Basic Gradients

The term $\frac{\partial \mathcal{L}}{\partial H_{kjn}^{(1)}}$ is calculated as

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial H_{cij}^{(1)}} &= \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(2)}} \frac{\partial H_{kmn}^{(2)}}{\partial H_{cij}^{(1)}} \\ &= \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(2)}} \frac{\partial \sum_{p,q} W_{pqm}^{(2)} H_{k+p,q,n}^{(1)}}{\partial H_{cij}^{(1)}} \\ &= \sum_{k,m} \frac{\partial \mathcal{L}}{\partial H_{kmj}^{(2)}} \frac{\partial \sum_{p,q} W_{pqm}^{(2)} H_{k+p,q,j}^{(1)}}{\partial H_{cij}^{(1)}}\end{aligned}$$

Basic Gradients

The term $\frac{\partial \mathcal{L}}{\partial H_{kjn}^{(1)}}$ is calculated as

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial H_{cij}^{(1)}} &= \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(2)}} \frac{\partial H_{kmn}^{(2)}}{\partial H_{cij}^{(1)}} \\ &= \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(2)}} \frac{\partial \sum_{p,q} W_{pqm}^{(2)} H_{k+p,q,n}^{(1)}}{\partial H_{cij}^{(1)}} \\ &= \sum_{k,m} \frac{\partial \mathcal{L}}{\partial H_{kmj}^{(2)}} \frac{\partial \sum_{p,q} W_{pqm}^{(2)} H_{k+p,q,j}^{(1)}}{\partial H_{cij}^{(1)}} \\ &= \sum_{k,m} \frac{\partial \mathcal{L}}{\partial H_{kmj}^{(2)}} W_{c-k,i,j}^{(2)}\end{aligned}$$

- The term $\frac{\partial \mathcal{L}}{\partial H_{kmn}^{(2)}}$ is backpropagated from last layer.

Basic Gradients

The gradient of the loss function with respect to the weight $W_{cij}^{(1)}$ is given by

Basic Gradients

The gradient of the loss function with respect to the weight $W_{cij}^{(1)}$ is given by

$$\frac{\partial \mathcal{L}}{\partial W_{cij}^{(1)}} = \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(1)}} \frac{\partial H_{kmn}^{(1)}}{\partial W_{cij}^{(1)}}$$

Basic Gradients

The gradient of the loss function with respect to the weight $W_{cij}^{(1)}$ is given by

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial W_{cij}^{(1)}} &= \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(1)}} \frac{\partial H_{kmn}^{(1)}}{\partial W_{cij}^{(1)}} \\ &= \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(1)}} \frac{\partial \sum_{p,q} W_{pqm}^{(1)} V_{k+p,q,n}}{\partial W_{cij}^{(1)}}\end{aligned}$$

Basic Gradients

The gradient of the loss function with respect to the weight $W_{cij}^{(1)}$ is given by

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial W_{cij}^{(1)}} &= \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(1)}} \frac{\partial H_{kmn}^{(1)}}{\partial W_{cij}^{(1)}} \\ &= \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(1)}} \frac{\partial \sum_{p,q} W_{pqm}^{(1)} V_{k+p,q,n}}{\partial W_{cij}^{(1)}} \\ &= \sum_{k,n} \frac{\partial \mathcal{L}}{\partial H_{kjn}^{(1)}} \frac{\partial \sum_{p,q} W_{pqj}^{(1)} V_{k+p,q,n}}{\partial W_{cij}^{(1)}}\end{aligned}$$

Basic Gradients

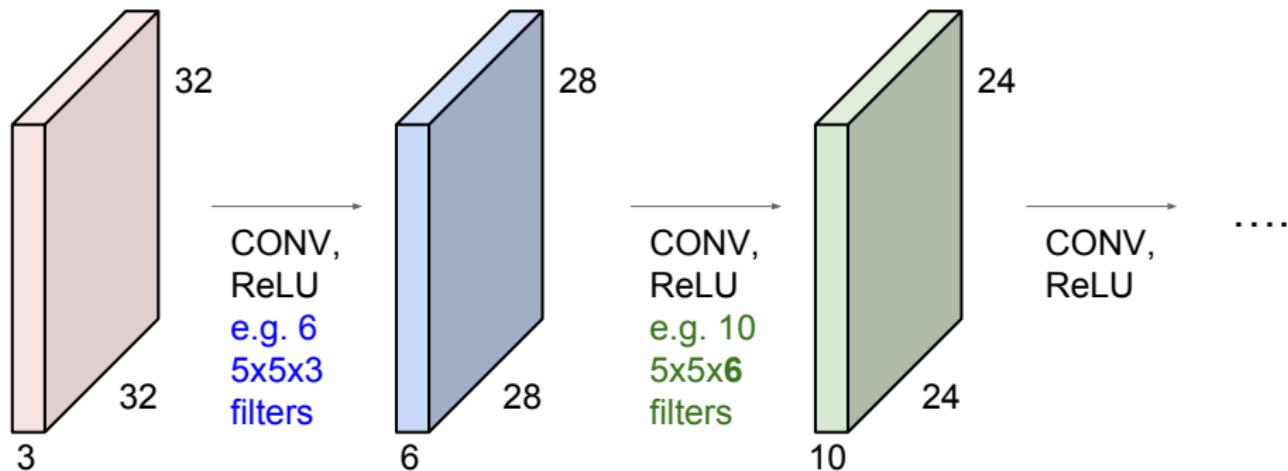
The gradient of the loss function with respect to the weight $W_{cij}^{(1)}$ is given by

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial W_{cij}^{(1)}} &= \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(1)}} \frac{\partial H_{kmn}^{(1)}}{\partial W_{cij}^{(1)}} \\ &= \sum_{k,m,n} \frac{\partial \mathcal{L}}{\partial H_{kmn}^{(1)}} \frac{\partial \sum_{p,q} W_{pqm}^{(1)} V_{k+p,q,n}}{\partial W_{cij}^{(1)}} \\ &= \sum_{k,n} \frac{\partial \mathcal{L}}{\partial H_{kjn}^{(1)}} \frac{\partial \sum_{p,q} W_{pqj}^{(1)} V_{k+p,q,n}}{\partial W_{cij}^{(1)}} \\ &= \sum_{k,n} \frac{\partial \mathcal{L}}{\partial H_{kjn}^{(1)}} V_{k+c,i,n}\end{aligned}$$

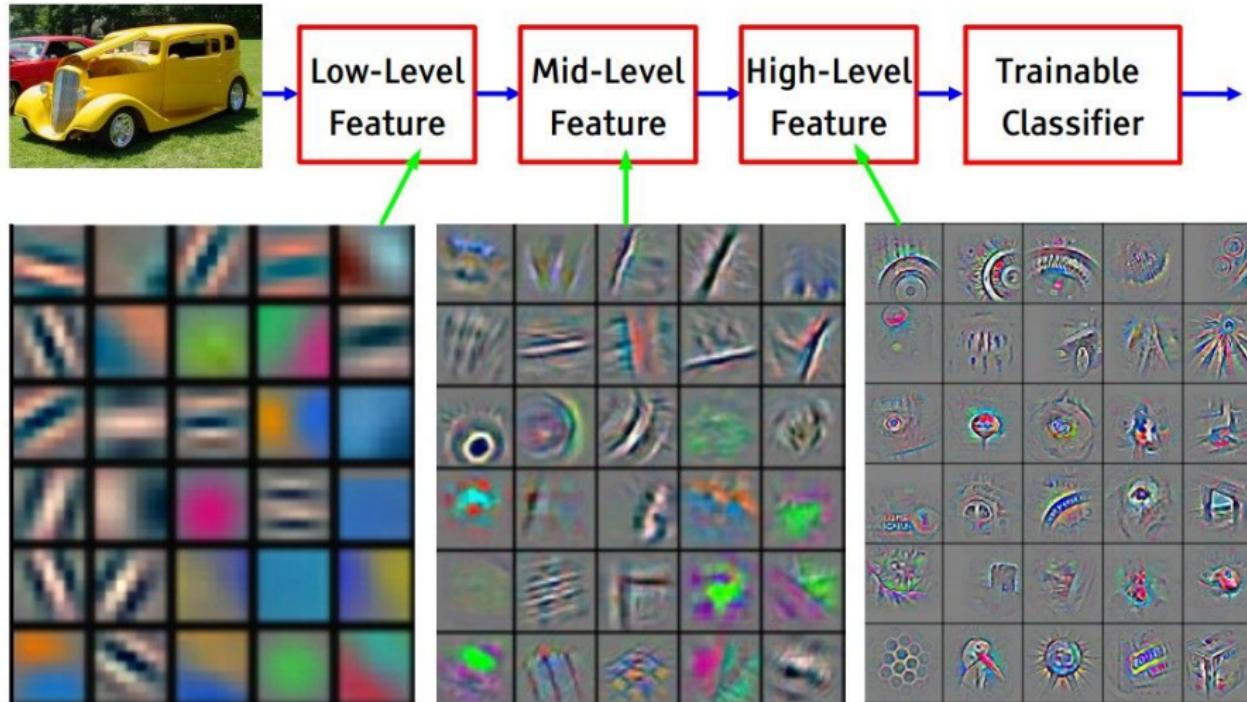
Convolutional Neural Networks (ConvNet)

Convolutional Neural Networks

ConvNet is a sequence of Convolutional Layers, interspersed with activation functions

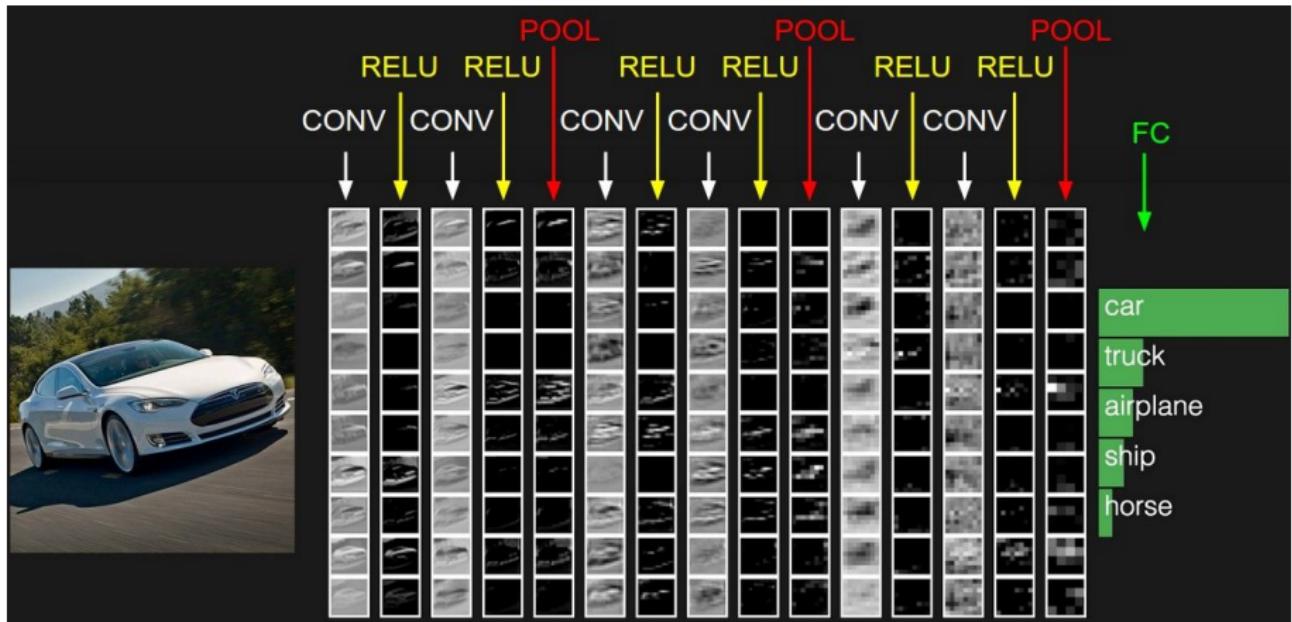


Features from ConvNet



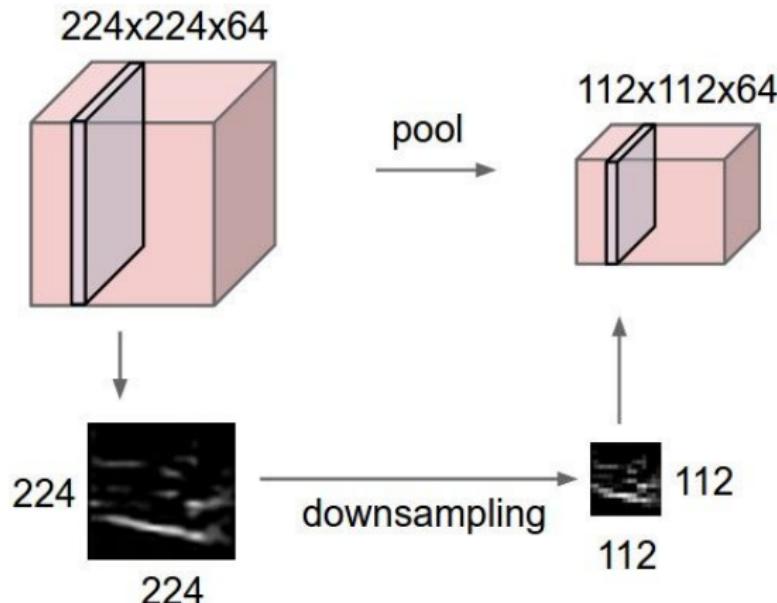
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

Two More Layers to Go: POOL/FC

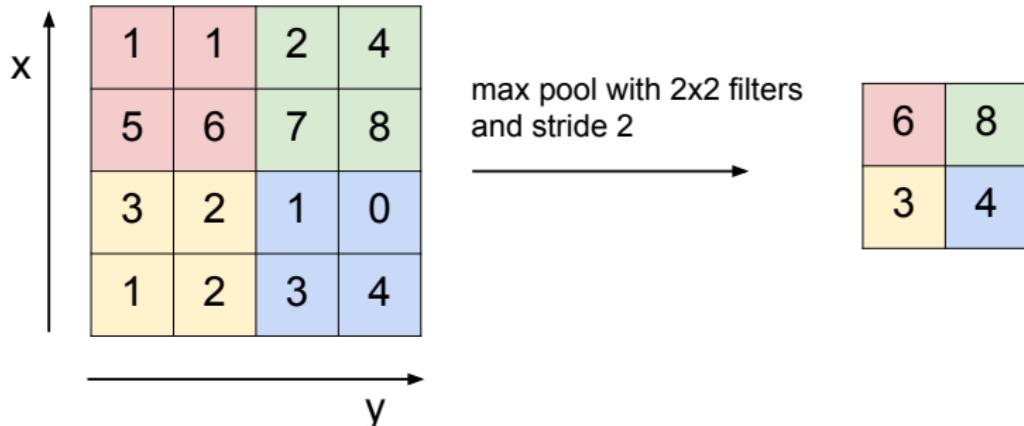


Pooling Layer

- ① Makes the representations smaller and more manageable.
- ② Operates over each activation map independently.



Max Pooling



Common settings:

- Filter size: 2×2 or 3×3
- Stride: 2

Max pooling is a non-linear operator.

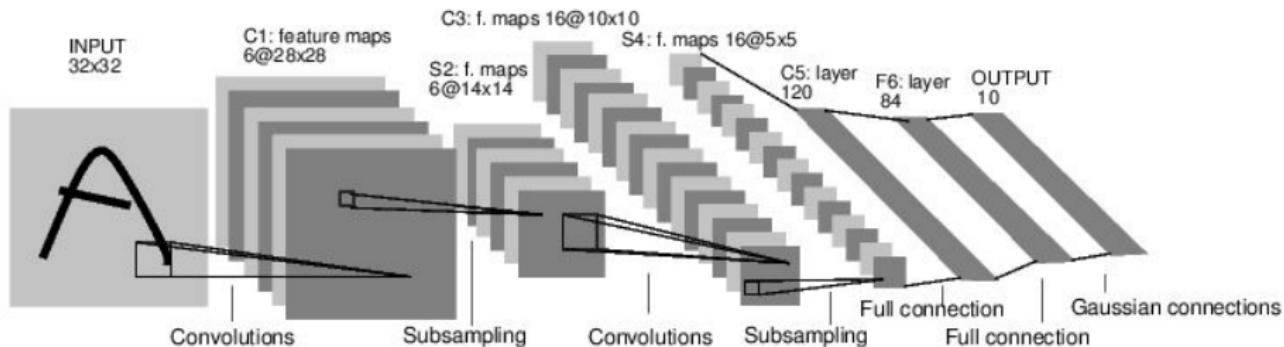
Max Pooling

Q: What is the gradient of a max pooling operation?

Case Study: CNN for MNIST



Case Study: LeNet-5 [LeCun *et al.*, 1998]



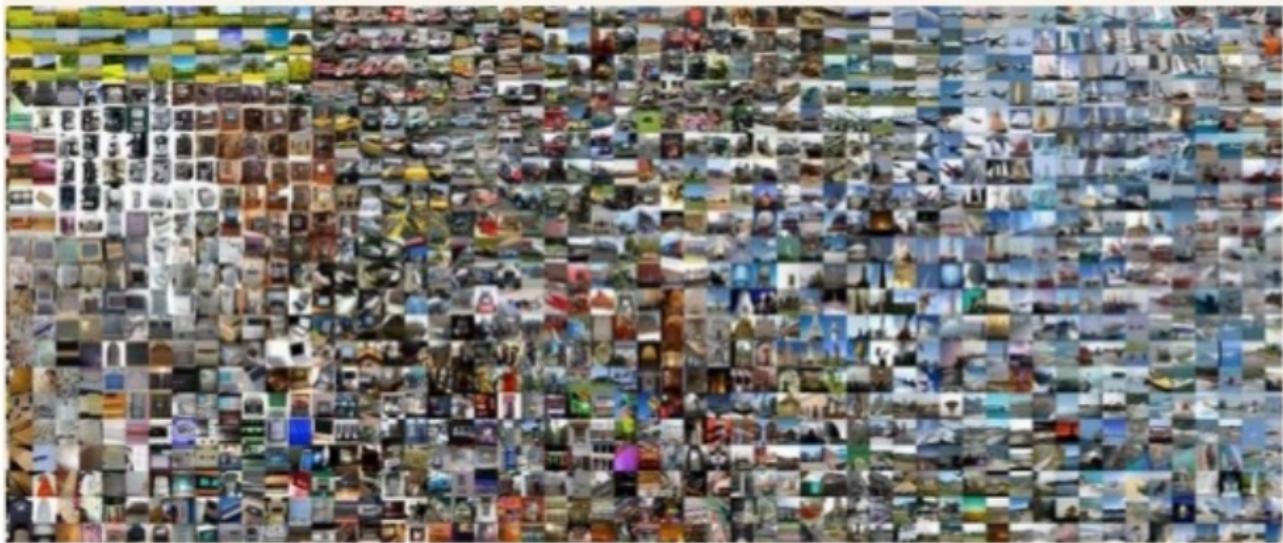
Conv filters were 5x5, applied at stride 1

Subsampling (Pooling) layers were 2x2 applied at stride 2
i.e. architecture is [CONV-POOL-CONV-POOL-CONV-FC]

Implementation

<https://github.com/sujaybabruwad/LeNet-in-Tensorflow/blob/master/LeNet-Lab.ipynb>

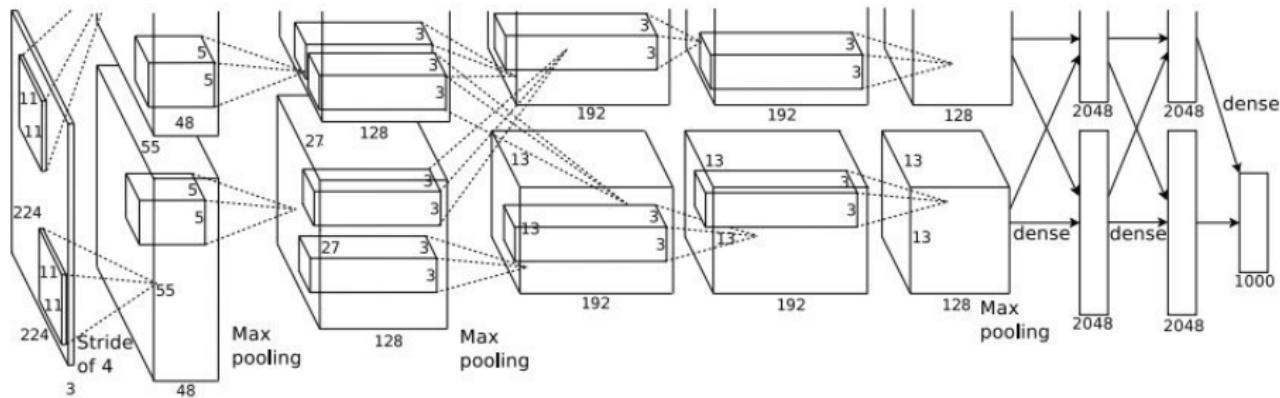
- About 14M images, 1000 classes.



<http://www.image-net.org/>

Case Study: AlexNet [Krizhevsky et al., 2012]

- Input: $227 \times 227 \times 3$ images.
- First layer (CONV1): 96 11×11 filters applied at stride 4.



- Q: What are the output volume and number of parameters in each layer?

Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) 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

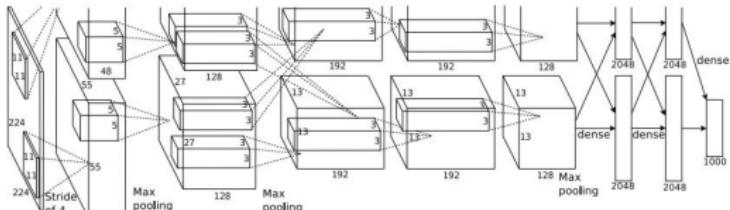
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



Details/Retrospectives:

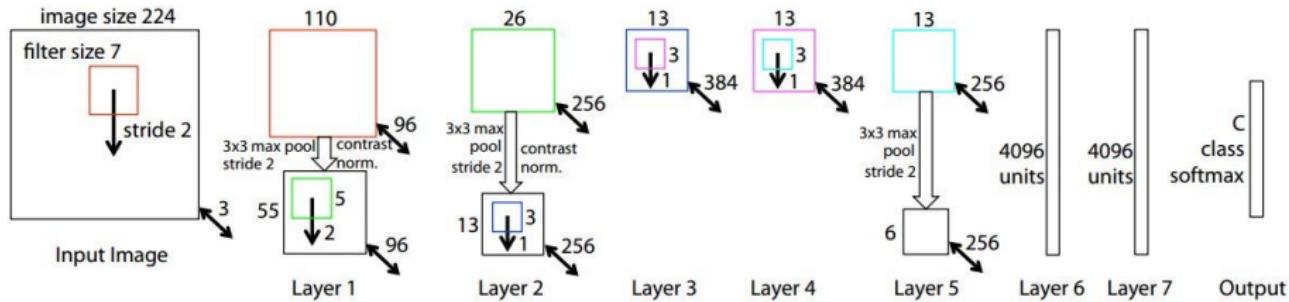
- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% -> 15.4%

Implementation

<https://github.com/vigneshthakkar/Deep-Nets/blob/master/AlexNet.py>

Case Study: ZFNet [Zeiler and Fergus, 2013]

- AlexNet but:
 - CONV1: change from $(11 \times 11 \text{ stride } 4)$ to $(7 \times 7 \text{ stride } 2)$
 - CONV3, 4, 5: instead of 384, 384, 256 filters use 512, 1024, 512
- ImageNet top 5 error: 15.4% \rightarrow 14.8%.



Implementation

<https://github.com/vigneshthakkar/Deep-Nets/blob/master/ZFNet.py>

Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Only 3x3 CONV stride 1, pad 1
and 2x2 MAX POOL stride 2

best model

11.2% top 5 error in ILSVRC 2013

->

7.3% top 5 error

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-128	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv1-256	conv3-256 conv3-256 conv1-512	conv3-256 conv3-256 conv3-256 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
FC-4096					
FC-4096					
FC-1000					
soft-max					

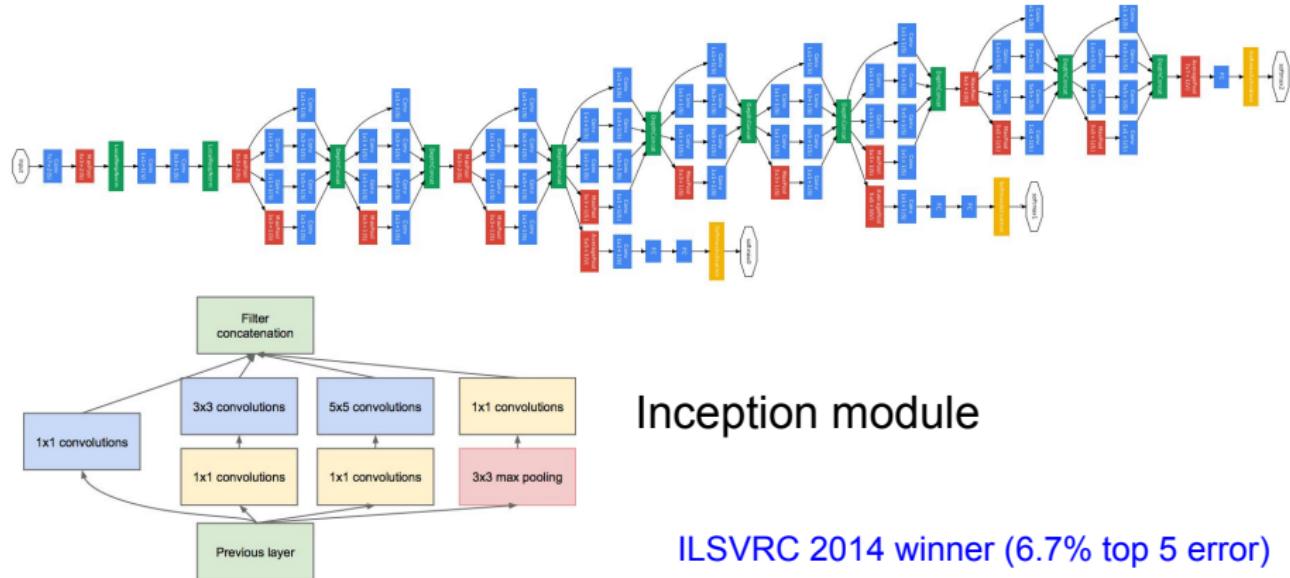
Table 2: Number of parameters (in millions).

Network	A,A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144

Implementation

<https://github.com/vigneshthakkar/Deep-Nets/blob/master/VGGNet.py>

Case Study: GoogLeNet [Szegedy *et al.*, 2014]



Case Study: GoogLeNet [Szegedy et al., 2014]

type	patch size/ stride	output size	depth	#1×1	#3×3 reduce	#3×3	#5×5 reduce	#5×5	pool proj	params	ops
convolution	7×7/2	112×112×64	1							2.7K	34M
max pool	3×3/2	56×56×64	0								
convolution	3×3/1	56×56×192	2		64	192				112K	360M
max pool	3×3/2	28×28×192	0								
inception (3a)		28×28×256	2	64	96	128	16	32	32	159K	128M
inception (3b)		28×28×480	2	128	128	192	32	96	64	380K	304M
max pool	3×3/2	14×14×480	0								
inception (4a)		14×14×512	2	192	96	208	16	48	64	364K	73M
inception (4b)		14×14×512	2	160	112	224	24	64	64	437K	88M
inception (4c)		14×14×512	2	128	128	256	24	64	64	463K	100M
inception (4d)		14×14×528	2	112	144	288	32	64	64	580K	119M
inception (4e)		14×14×832	2	256	160	320	32	128	128	840K	170M
max pool	3×3/2	7×7×832	0								
inception (5a)		7×7×832	2	256	160	320	32	128	128	1072K	54M
inception (5b)		7×7×1024	2	384	192	384	48	128	128	1388K	71M
avg pool	7×7/1	1×1×1024	0								
dropout (40%)		1×1×1024	0								
linear		1×1×1000	1							1000K	1M
softmax		1×1×1000	0								

Fun features:

- Only 5 million params!
(Removes FC layers completely)

Compared to AlexNet:

- 12X less params
- 2x more compute
- 6.67% (vs. 16.4%)

Implementation

<https://github.com/vigneshthakkar/Deep-Nets/blob/master/GoogLeNet.py>

Case Study: ResNet [He et al., 2015]

ILSVRC 2015 winner (3.6% top 5 error)

Microsoft
Research

MSRA @ ILSVRC & COCO 2015 Competitions

- **1st places** in all five main tracks

- ImageNet Classification: "*Ultra-deep*" (quote Yann) **152-layer** nets
- ImageNet Detection: **16%** better than 2nd
- ImageNet Localization: **27%** better than 2nd
- COCO Detection: **11%** better than 2nd
- COCO Segmentation: **12%** better than 2nd

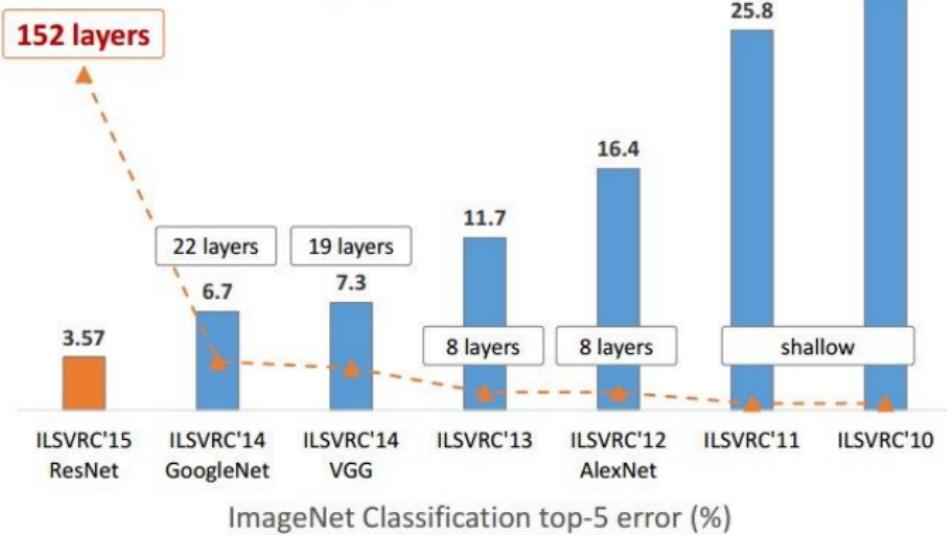
*improvements are relative numbers



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

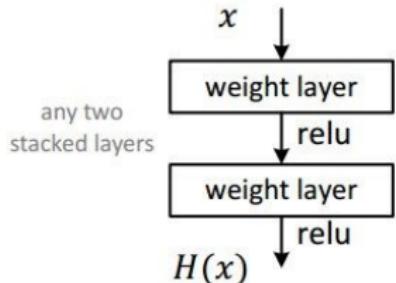
Slide from Kaiming He's recent presentation <https://www.youtube.com/watch?v=1PGLj-uKT1w>

Revolution of Depth

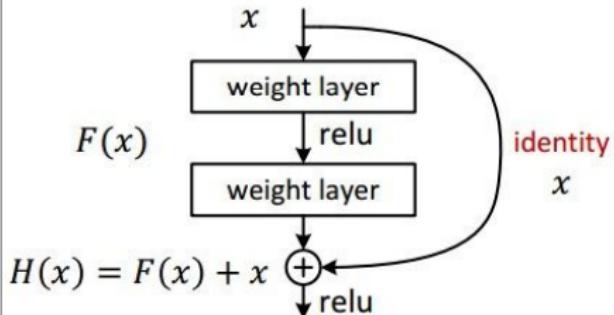


Case Study: ResNet [He et al., 2015]

- Plain net



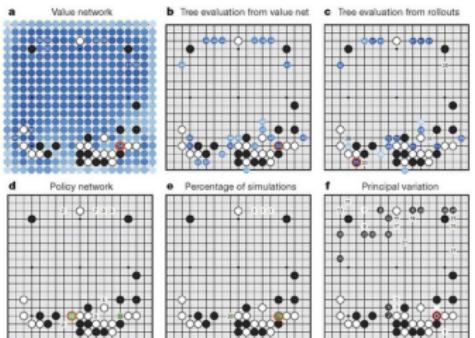
- Residual net



Implementation

<https://github.com/vigneshthakkar/Deep-Nets/blob/master/ResNet.py>

Case Study Bonus: DeepMind's AlphaGo



Case Study Bonus: DeepMind's AlphaGo

The input to the policy network is a $19 \times 19 \times 48$ image stack consisting of 48 feature planes. The first hidden layer zero pads the input into a 23×23 image, then convolves k filters of kernel size 5×5 with stride 1 with the input image and applies a rectifier nonlinearity. Each of the subsequent hidden layers 2 to 12 zero pads the respective previous hidden layer into a 21×21 image, then convolves k filters of kernel size 3×3 with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size 1×1 with stride 1, with a different bias for each position, and applies a softmax function. The match version of AlphaGo used $k = 192$ filters; [Fig. 2b](#) and [Extended Data Table 3](#) additionally show the results of training with $k = 128, 256$ and 384 filters.

policy network:

[$19 \times 19 \times 48$] Input

CONV1: 192 5×5 filters , stride 1, pad 2 => [$19 \times 19 \times 192$]

CONV2..12: 192 3×3 filters, stride 1, pad 1 => [$19 \times 19 \times 192$]

CONV: 1 1×1 filter, stride 1, pad 0 => [19×19] (*probability map of promising moves*)

Summary

- ➊ ConvNets stack CONV,POOL,FC layers.
- ➋ Trend towards smaller filters and deeper architectures.
- ➌ Trend towards getting rid of POOL/FC layers (just CONV and action layers).
- ➍ Typical architectures look like:

$[(\text{CONV-ReLU})^* N\text{-Pool?}]^* M\text{-(FC-ReLU)}^* K, \text{SOFTMAX}$

- N is usually up to ~ 5 , M is large, $0 \leq K \leq 2$.
- recent advances such as ResNet/GoogLeNet challenge this paradigm.