

Convolutional Neural Networks

Hao Dong

2019, Peking University

Convolutional Neural Networks

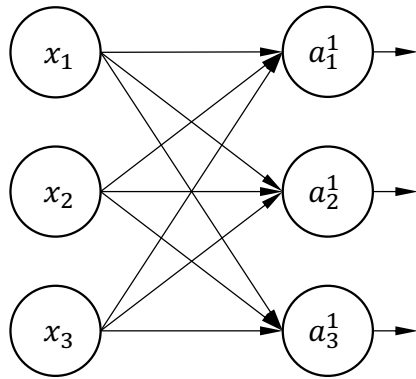
- Motivation
- Convolutional Algorithms
- Pooling Algorithms
- Hierarchical Representation Learning
- Convolutional Architectures
- Transposed Convolutional Algorithms
- Computer Vision Applications

Motivation

Motivation

- Curse of Dimensionality

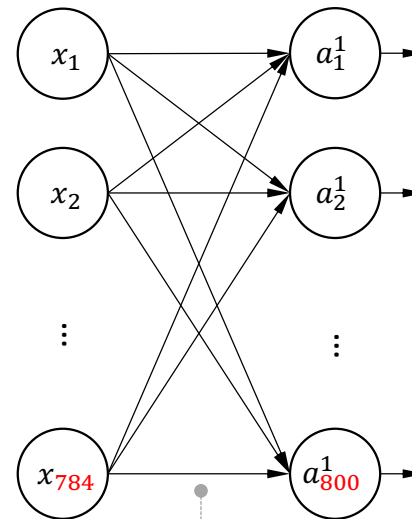
When we only have three inputs:



Fully connected layer

When we have a small image:

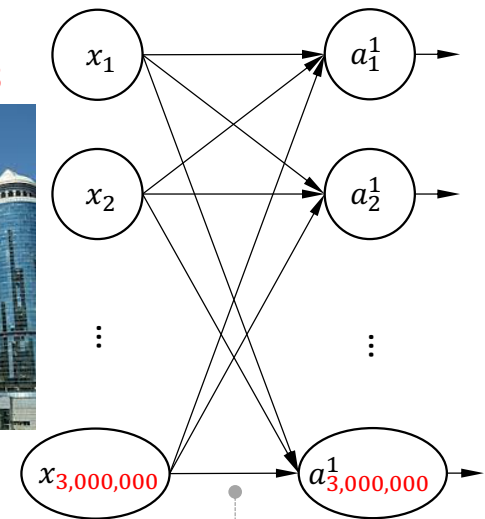
$28 \times 28 \times 1$



W is a 784×800 matrix

When we have a HD image:

$1K \times 1K \times 3$



W cannot fit into the memory ...

Motivation

- Eye != Fully Connected Layer

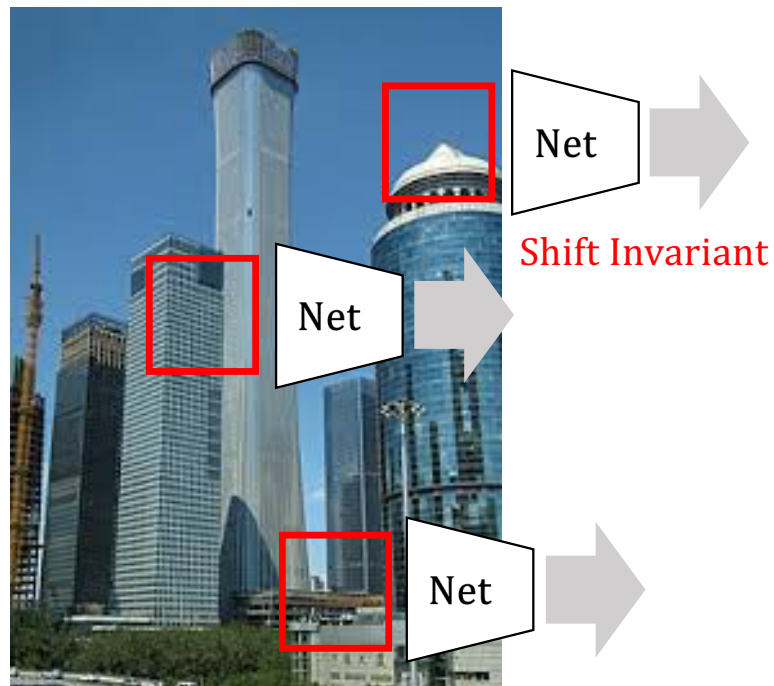
- Do human analysis each pixel using different ways? **NO**



Shift Invariant

Motivation

- Eye != Fully Connected Layer



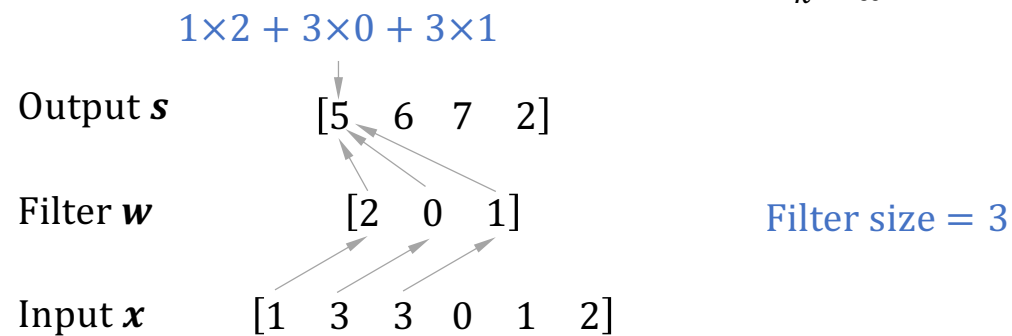
- Parameters shared over space
- Sparse connectivity
- Equivariant representation

Convolutional Algorithms

Convolutional Algorithms

- Convolution on 1D vector

1D discrete convolution: $\mathbf{s}_t = (\mathbf{x} * \mathbf{w})_t = \sum_{k=-\infty}^{\infty} \mathbf{x}_k \mathbf{w}_{t-k}$



- The more similar the filter (kernel) and local path are, the higher value on the output

Convolutional Algorithms

- Edge detection on greyscale 2D images

2D discrete convolution: $s_{r,c} = ((x * W))_{r,c} = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} x_{r,c} w_{r-i,c-j}$

$h \times w \times 1$ image



$h \times w \times 1$ image



$$W = \begin{bmatrix} -1 & -1 & -1 \\ 2 & 2 & 2 \\ -1 & -1 & -1 \end{bmatrix}$$

Detect vertical edges

Convolutional Algorithms

- Edge detection on RGB Images

$$\mathbf{R} \quad \mathbf{W}_{:,1} = \begin{bmatrix} -1 & -1 & -1 \\ 2 & 2 & 2 \\ -1 & -1 & -1 \end{bmatrix}$$

$$\mathbf{G} \quad \mathbf{W}_{:,2} = \begin{bmatrix} -1 & -1 & -1 \\ 2 & 2 & 2 \\ -1 & -1 & -1 \end{bmatrix}$$

$$\mathbf{B} \quad \mathbf{W}_{:,3} = \begin{bmatrix} -1 & -1 & -1 \\ 2 & 2 & 2 \\ -1 & -1 & -1 \end{bmatrix}$$

\mathbf{W} is a $3 \times 3 \times 3$ tensor

Height of filter size

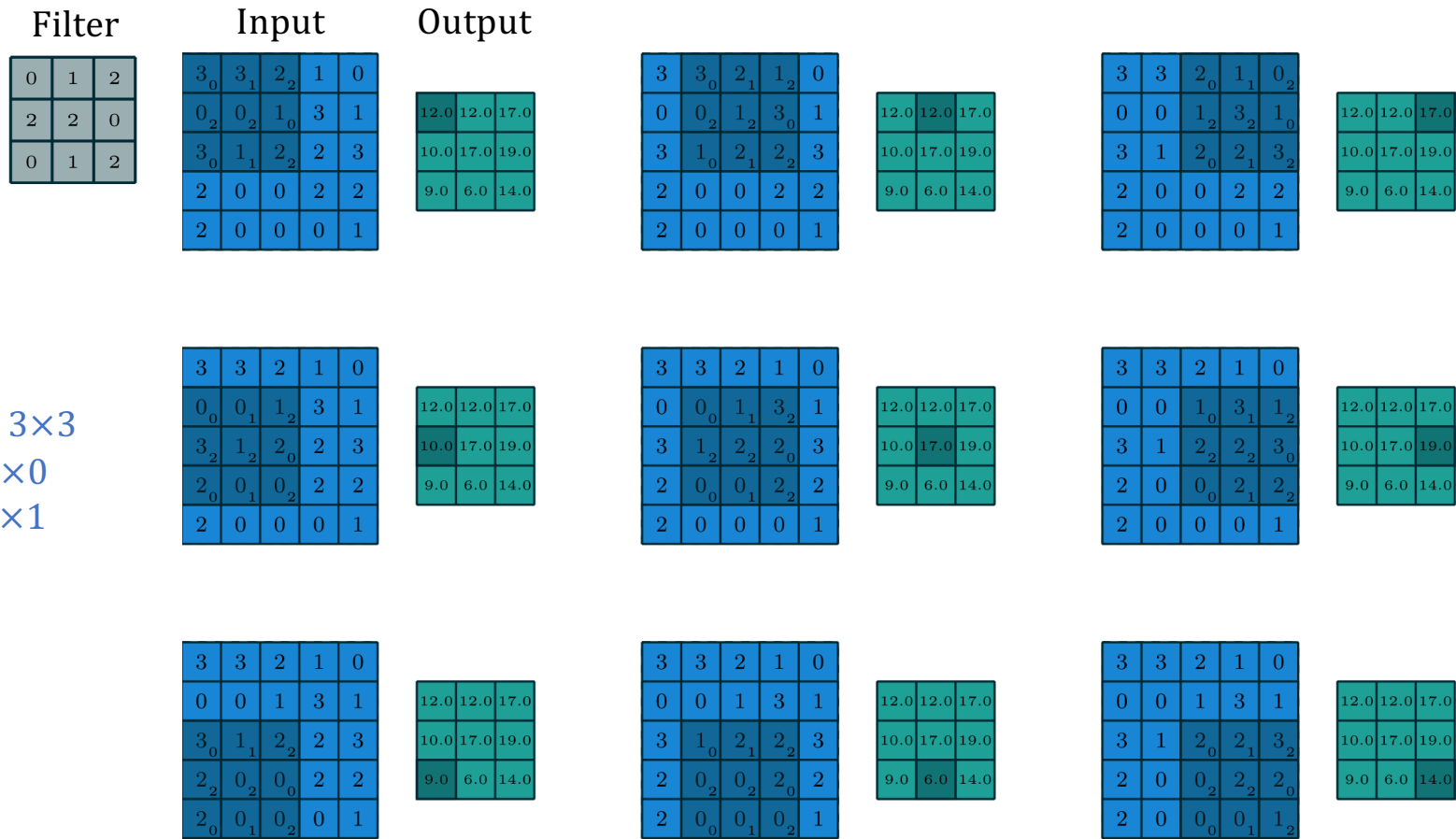
Width of filter size

Depth of filter size == Number of input channels

Each input channel needs one filter

Convolutional Algorithms

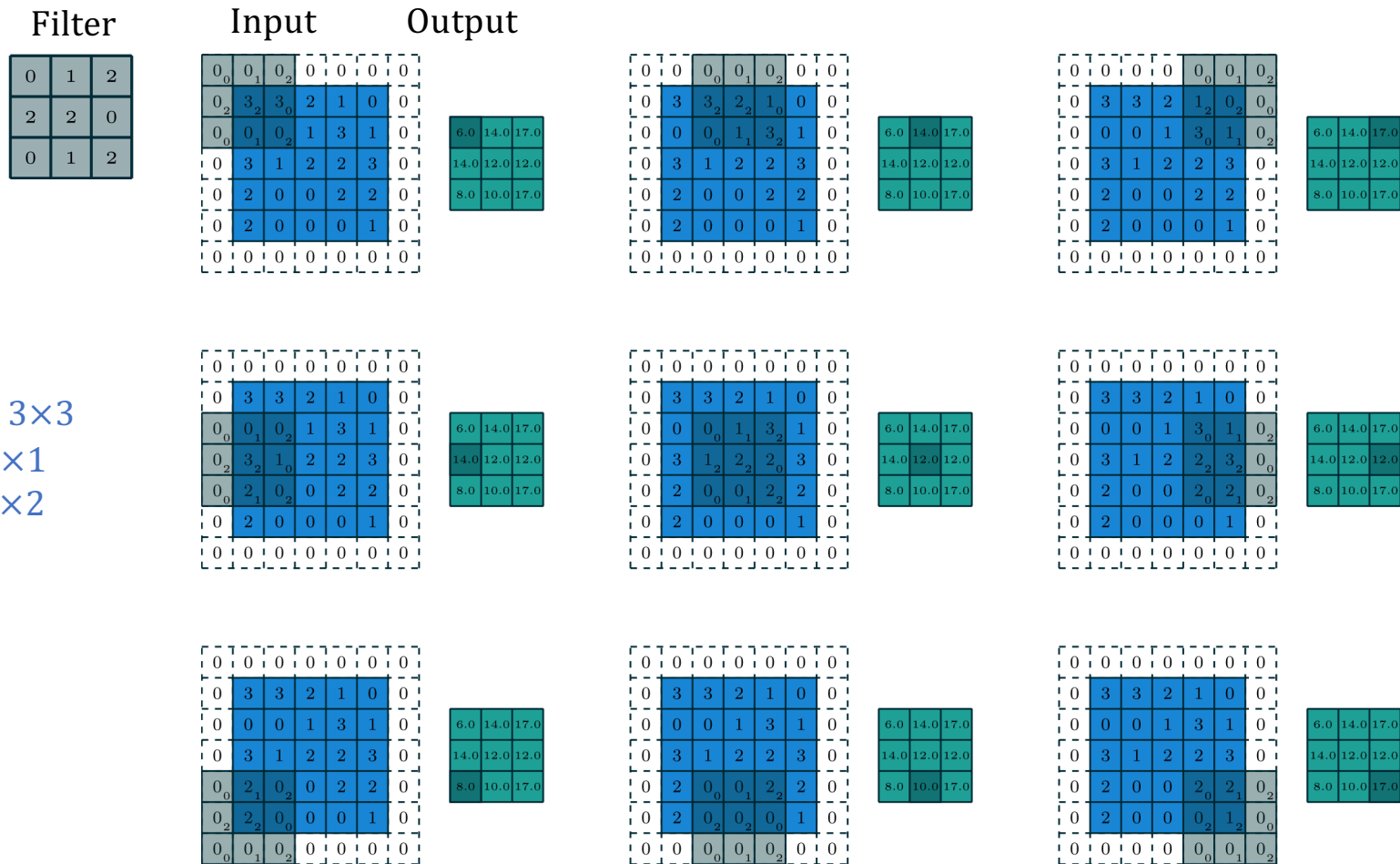
- No padding, unit strides



Filter size = 3×3
padding = 0×0
Strides = 1×1

Convolutional Algorithms

- Zero padding, non-unit strides



Filter size = 3×3
padding = 1×1
Strides = 2×2

Convolutional Algorithms

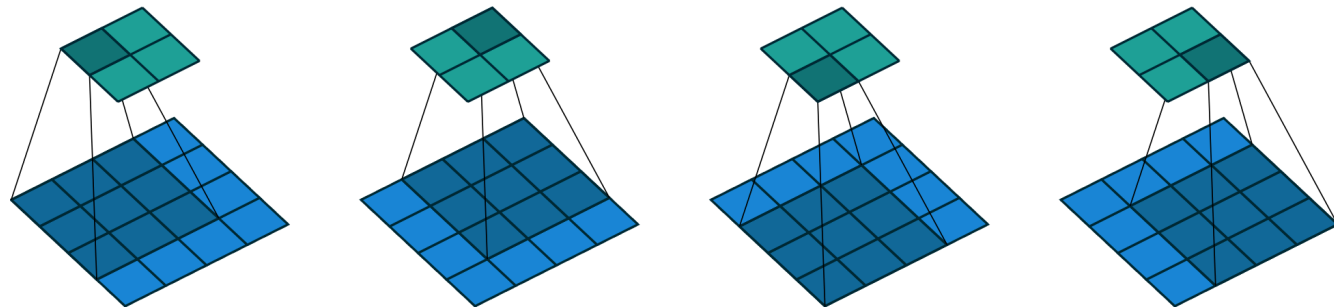
- Side Views

No padding and no strides

Filter size = 3×3

padding = 0×0

Strides = 1×1

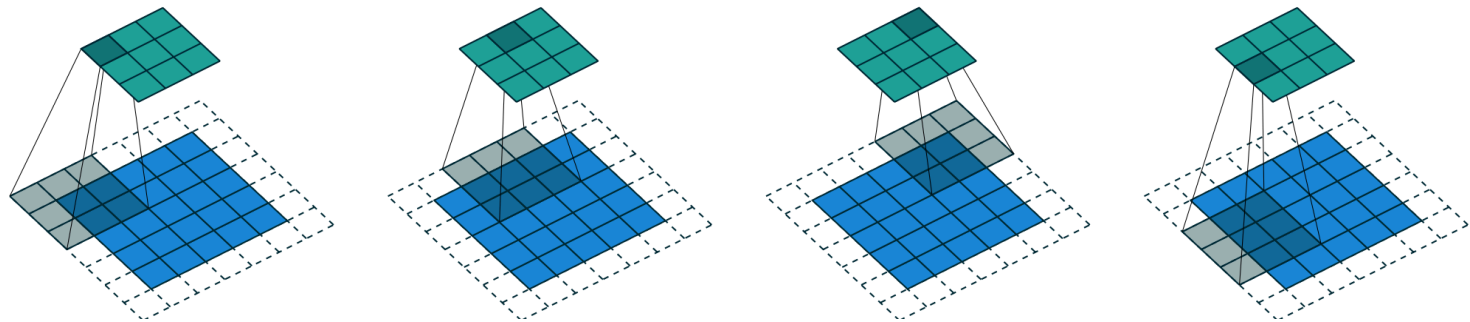


Padding and strides

Filter size = 3×3

padding = 1×1

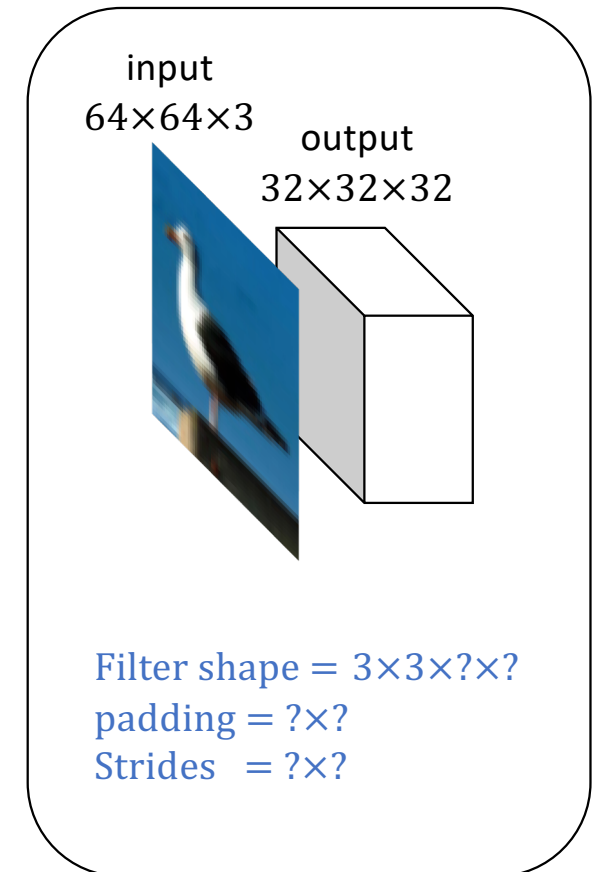
Strides = 2×2



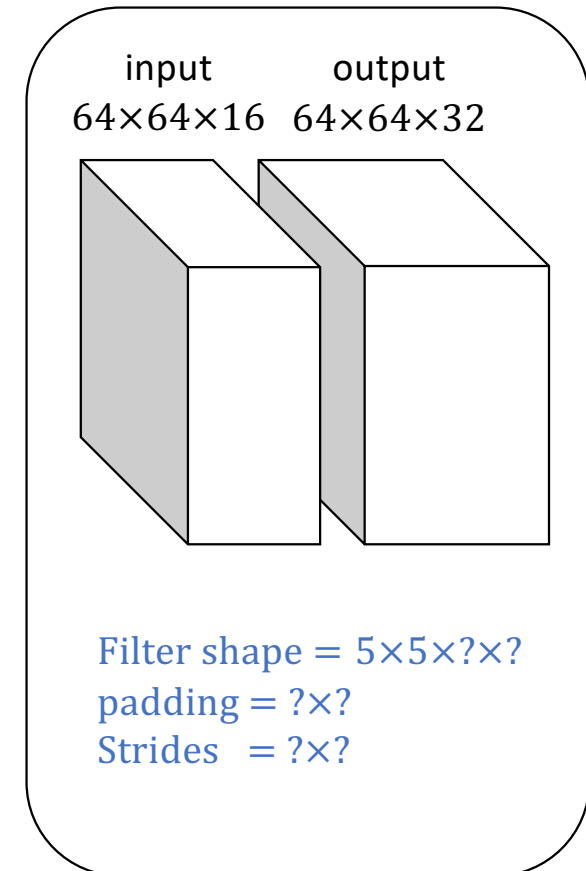
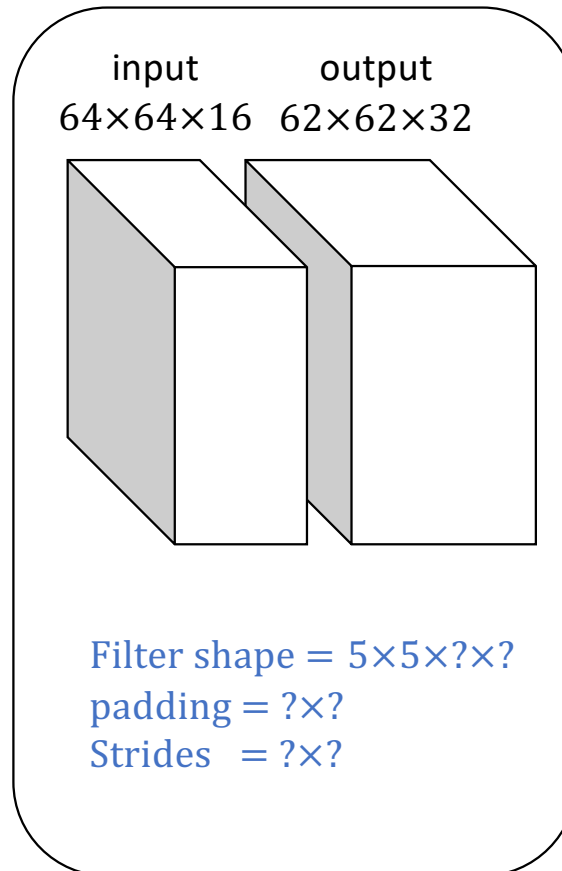
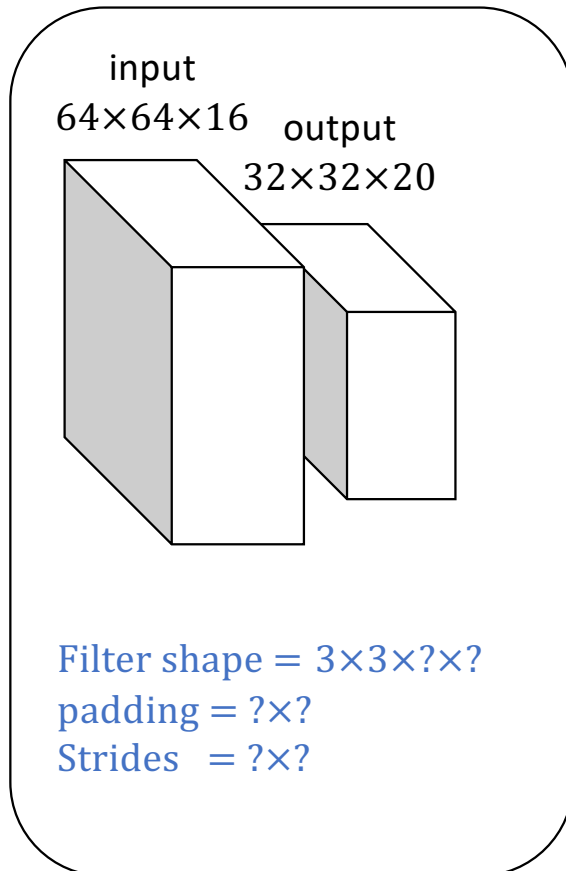
Convolutional Algorithms

- Filter (kernel) shapes
 - To convolute RGB images, a filter have 3 channels
 - The filter shape is
(filter_height \times filter_width \times input_channels \times n_filters)
- Number of filters (channels)
 - One filter can have one corresponding output feature map
 - The number of output channels == The number of filters
 - The output shape is (feature_height \times feature_width \times n_filters)

Rely on input size, padding, strides, and filter size



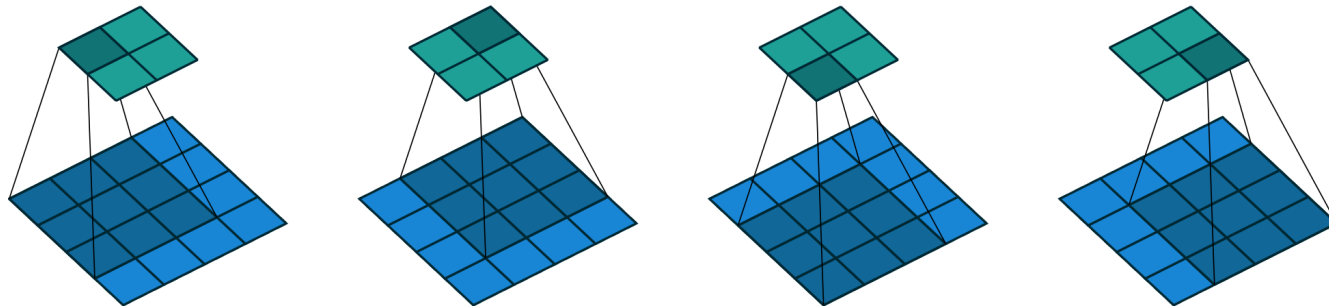
Convolutional Algorithms



Convolutional Algorithms

- Receptive field

Receptive field = 3×3



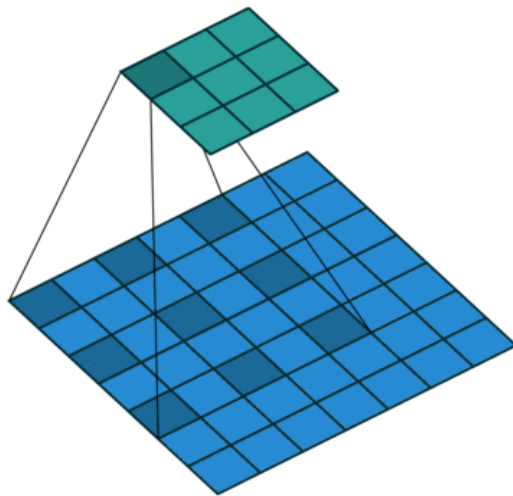
Convolutional Algorithms

- Receptive field

Layer #	Kernel Size	Stride	Dilation	Padding	Input Size	Output Size	Receptive Field
1	3	1	1	2	256	256	3
2	3	2	1	1	256	128	5
3	3	1	1	2	128	128	9
4	3	1	1	2	128	128	13
5	3	2	1	1	128	64	17
6	3	1	1	2	64	64	25
7	3	1	1	2	64	64	33
8	3	2	1	1	64	32	41

Convolutional Algorithms

- Dilated convolution: Larger Receptive Field



Filter

0	1	2
2	2	0
0	1	2



New Filter

0	0	1	0	2
0	0	0	0	0
2	0	2	0	0
0	0	0	0	0
0	0	1	0	2

Dilation rate = 2×2

Receptive field : $3 \times 3 \rightarrow 5 \times 5$

Convolutional Algorithms

- Convolution on 3D volume
 - Filter (kernel) shapes
 - To convolute MRI images, a filter have 1 channels
 - The filter shape is ($\text{filter_depth} \times \text{filter_height} \times \text{filter_width} \times \text{input_channels} \times \text{n_filters}$)
 - Number of filters (channels)
 - One filter can have one corresponding output feature **volume**
 - The number of output channels == The number of filters
 - The output shape is ($\text{feature_depth} \times \text{feature_height} \times \text{feature_width} \times \text{n_filters}$)



Pooling Algorithms

Pooling Algorithms

- Motivation: Shift Invariant

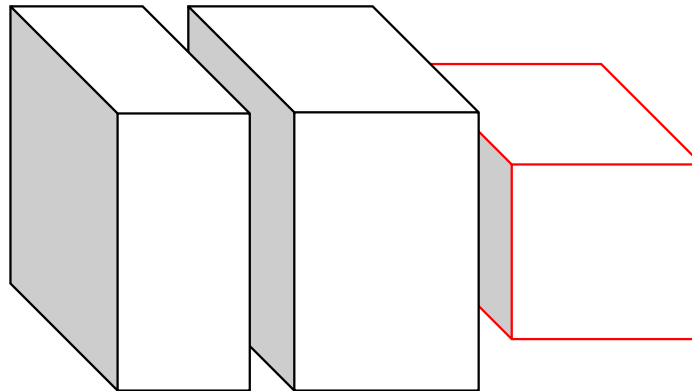
- Small shifts should not effect the result



Pooling Algorithms

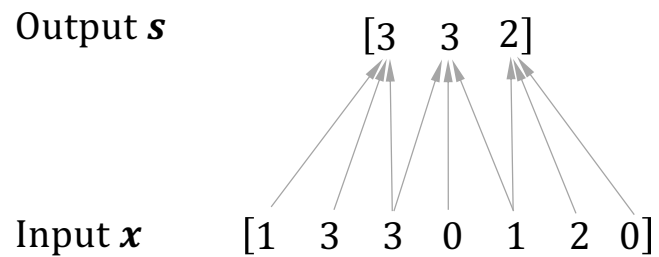
- Motivation: Features aggregation

input conv output after pooling
 $64 \times 64 \times 16$ $64 \times 64 \times 32$ $32 \times 32 \times 32$



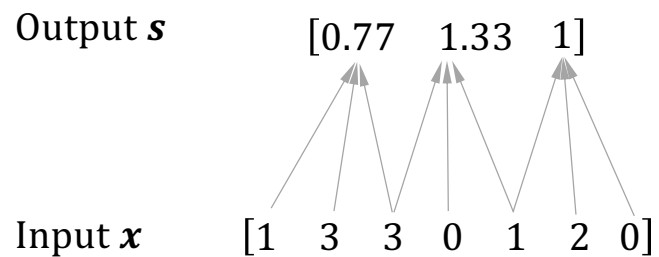
Pooling Algorithms

- Pooling on 1D vector
 - Max Pooling



Filter size = 3
 Padding = 0
 Stride = 2

- Mean (Average) Pooling

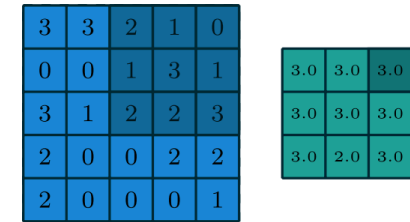
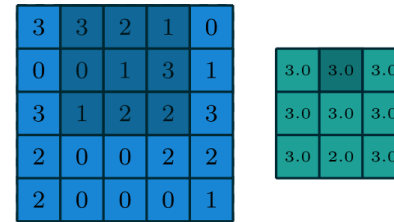
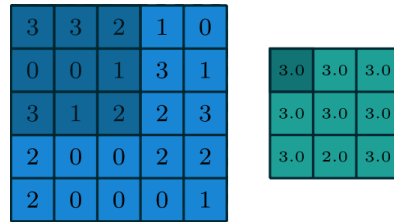


Filter size = 3
 Padding = 0
 Stride = 2

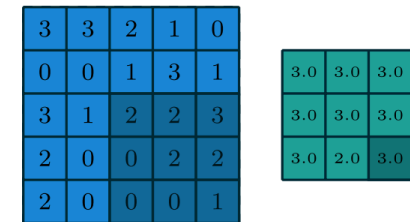
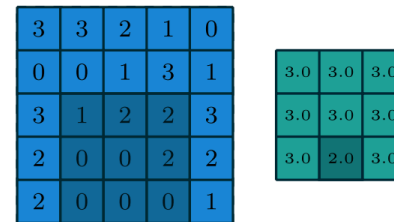
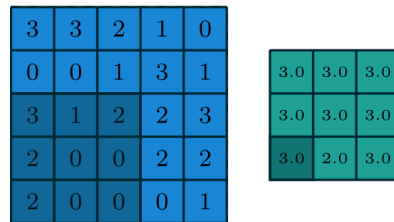
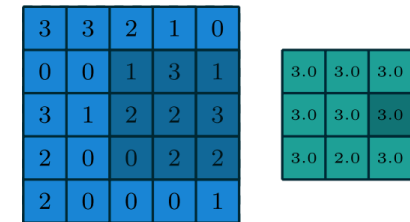
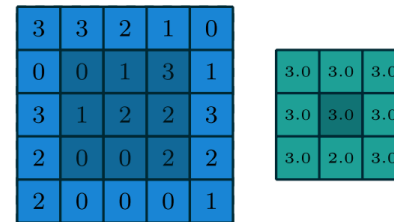
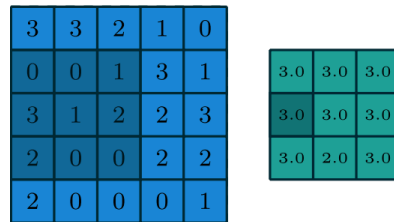
Receptive field is increased

Pooling Algorithms

- Max Pooling on 2D matrix



Filter size = 3×3
padding = 0×0
Strides = 1×1



Pooling Algorithms

- Mean Pooling on 2D matrix

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

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

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

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

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

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

Filter size = 3×3
padding = 0×0
Strides = 1×1

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

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

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

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

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

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

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

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

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

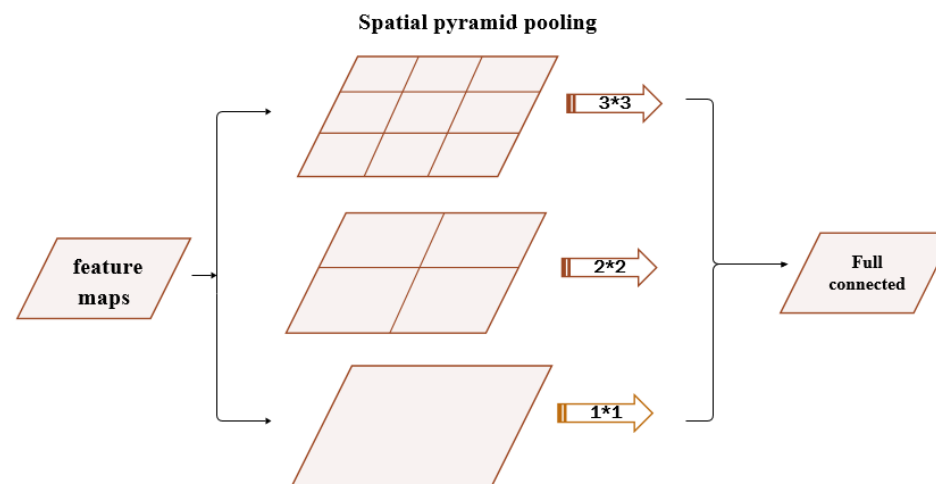
1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

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

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

Pooling Algorithms

- Spatial Pyramid Pooling



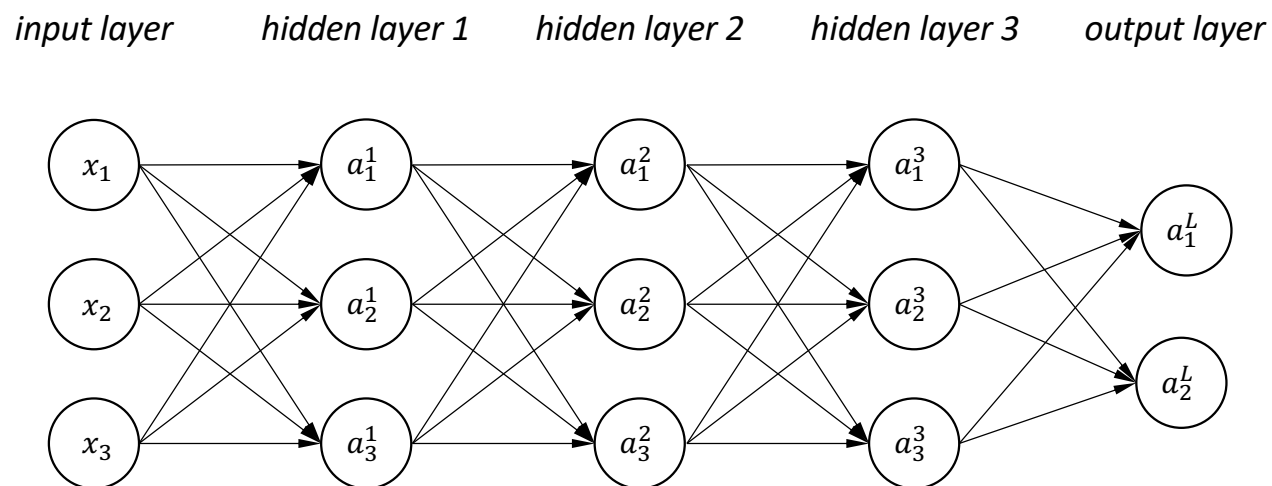
Pooling Algorithms

- Pooling on 3D volume
 -

Hierarchical Representation Learning

Hierarchical Representation Learning

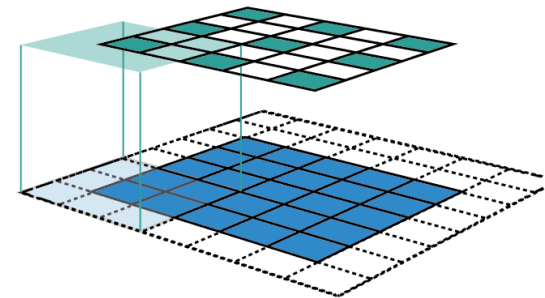
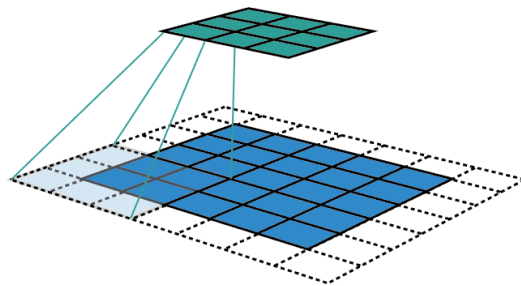
- Motivation



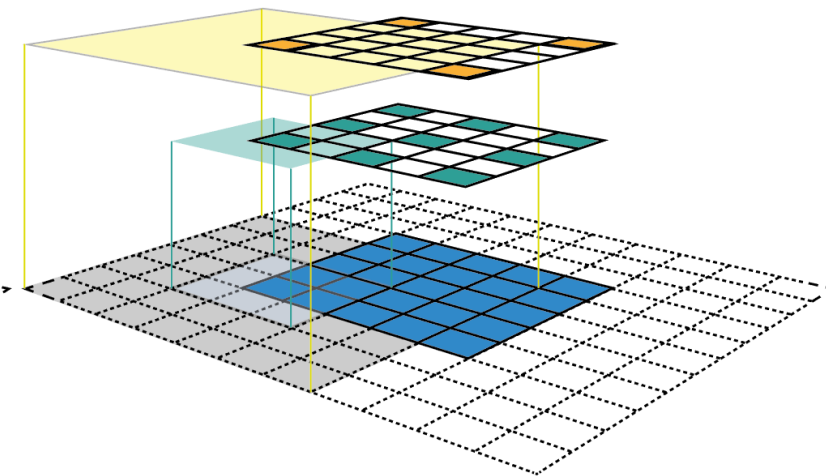
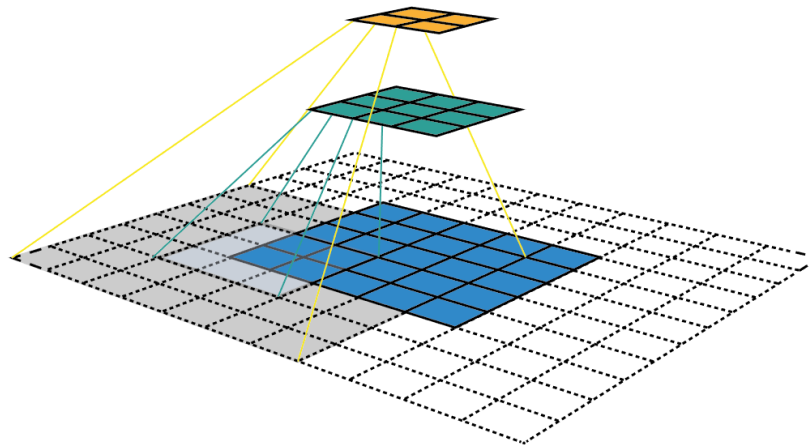
Hierarchical Representation Learning

- Large receptive field

1 layer



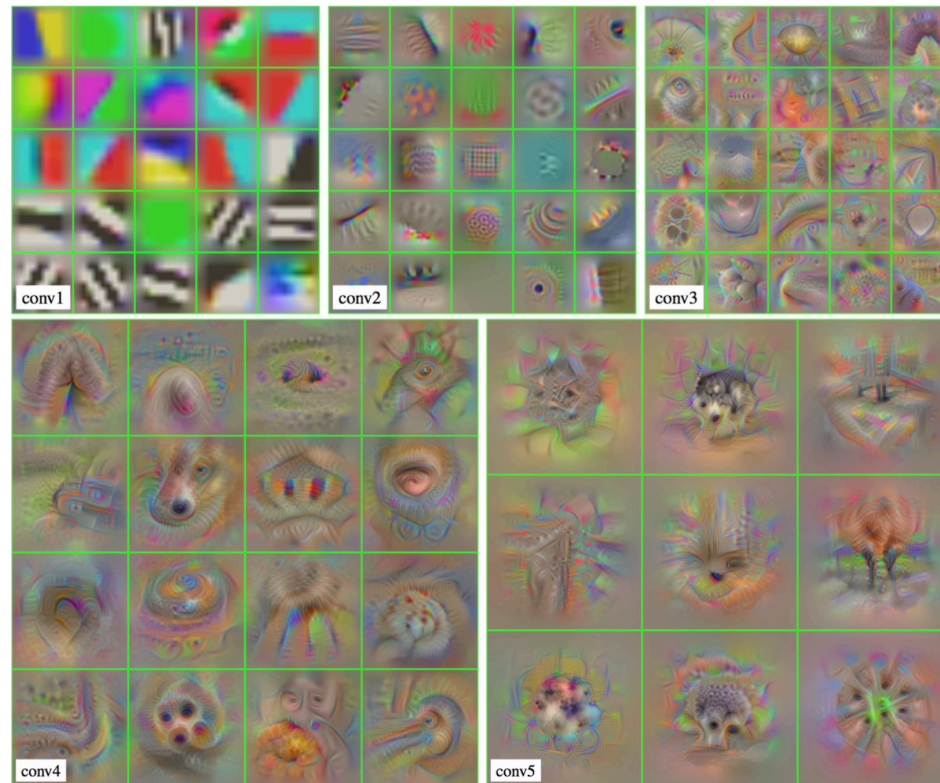
2 layers



Hierarchical Representation Learning

- Visualizing the activation outputs

Activation maximization of the first filters of each convolutional layer in a well-known CNN architecture called VGG. The notations of “conv1” through “conv5” distinct hidden layers, where a larger number represents a deeper hidden layer.

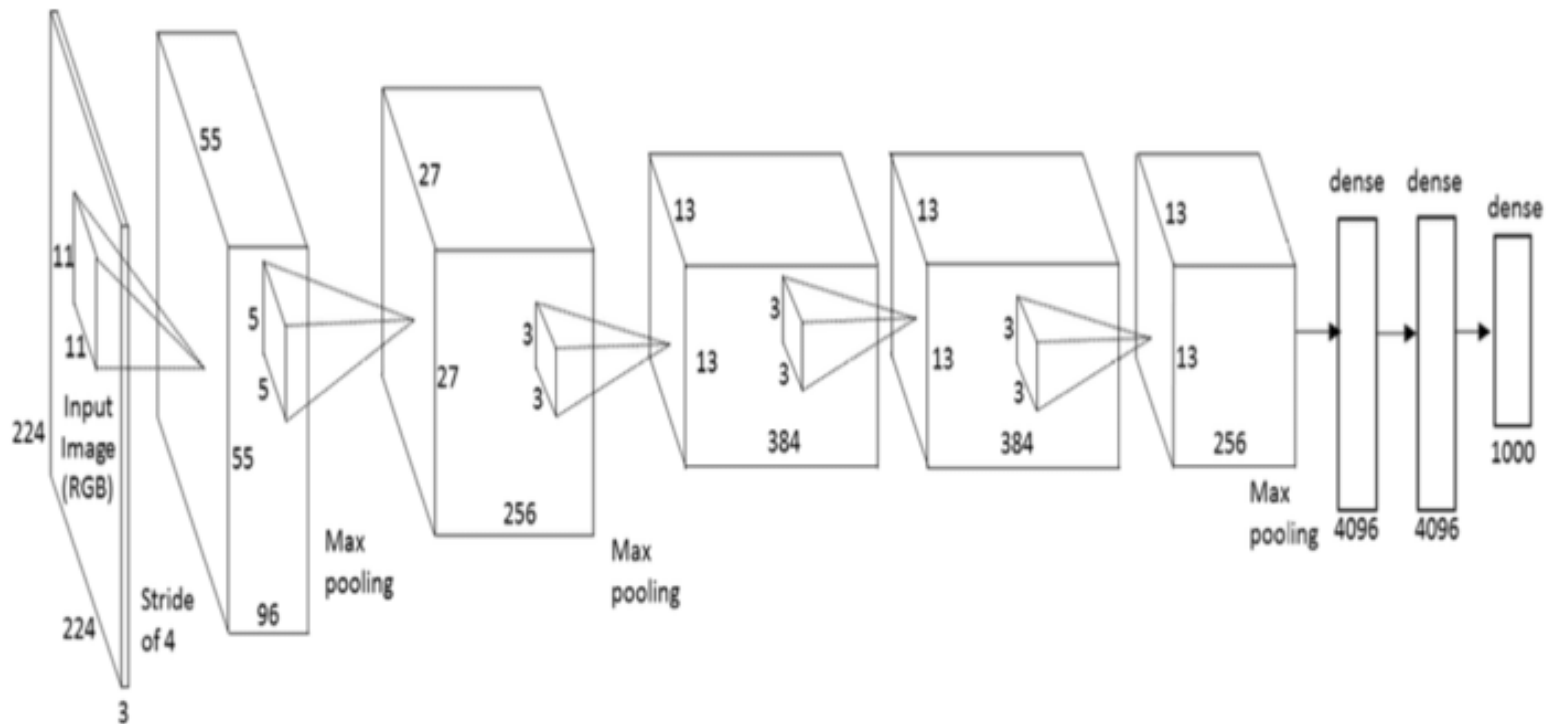


A. Mahendran and A. Vedaldi, “Visualizing deep convolutional neural networks using natural pre-images,”
International Journal of Computer Vision (IJCV), vol. 120, no. 3, pp. 233–255, 2016.

Convolutional Architectures

Convolutional Encoding Architectures

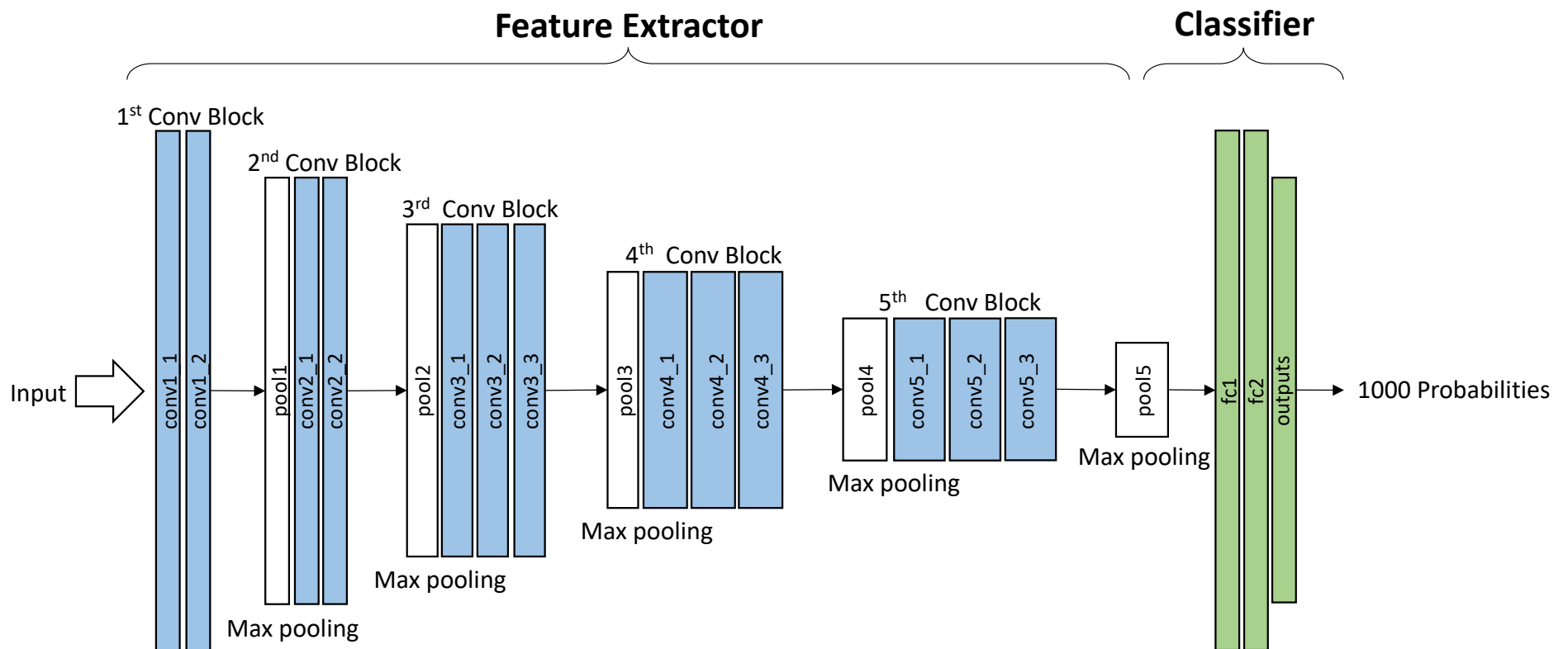
- AlexNet: make the history



In 2012, AlexNet achieves 10% performance gain on ImageNet compared to its previous methods

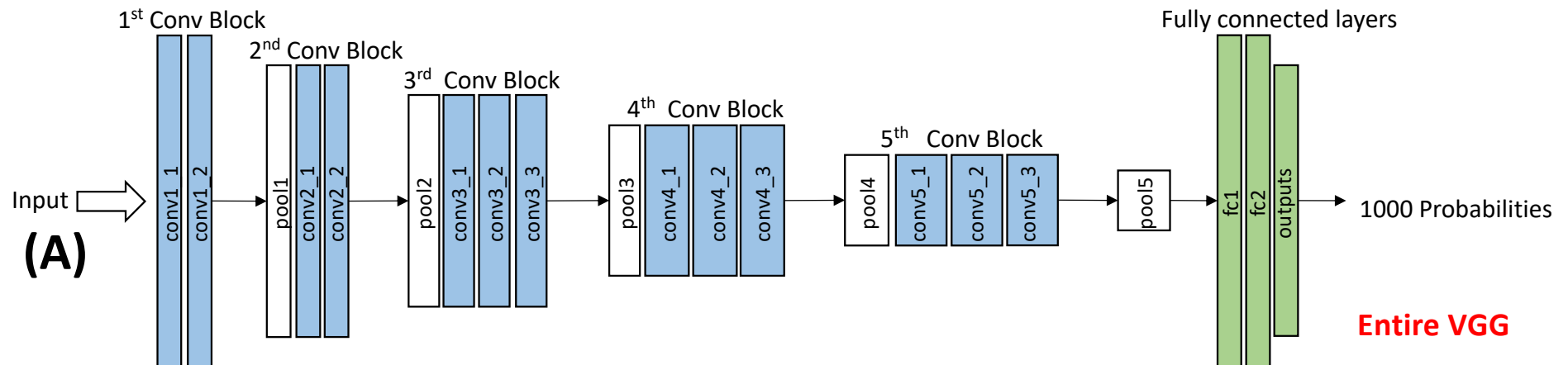
Convolutional Encoding Architectures

- VGG16



Convolutional Encoding Architectures

- VGG16

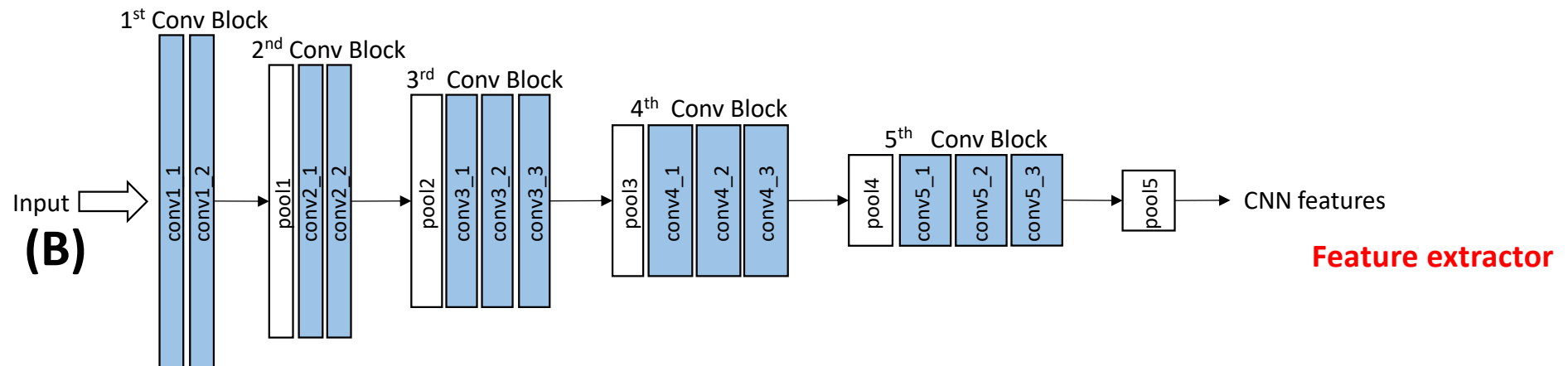


```

>>> # get the whole model, without pre-trained VGG parameters
>>> vgg = tl.models.vgg16()
>>> # get the whole model, restore pre-trained VGG parameters
>>> vgg = tl.models.vgg16(pretrained=True)
>>> # use for inferencing
>>> output = vgg(image, is_train=False)
>>> probs = tf.nn.softmax(output)[0].numpy()
    
```

Convolutional Encoding Architectures

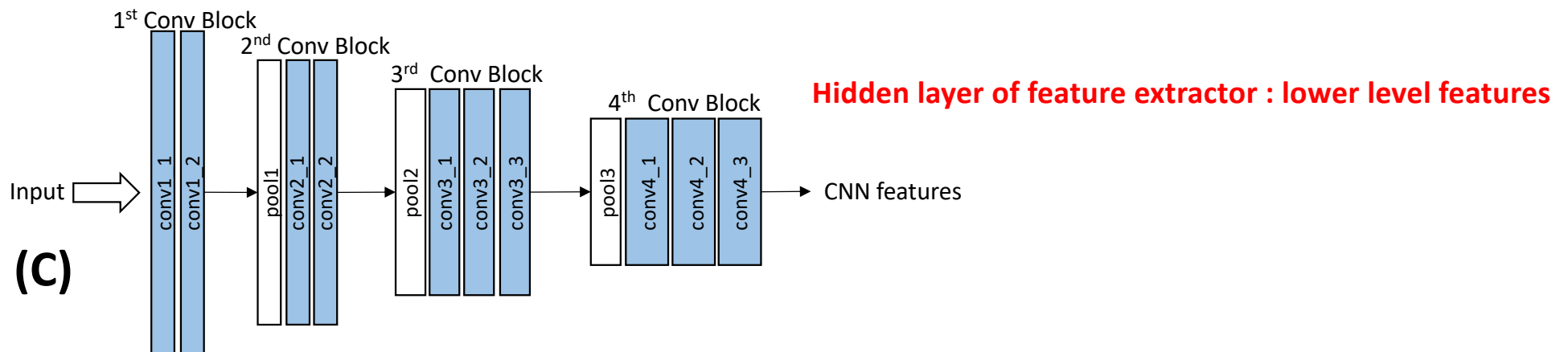
- VGG16



```
>>> vgg = tl.models.vgg16(pretrained=True, end_with='pool5', mode='static')
```

Convolutional Encoding Architectures

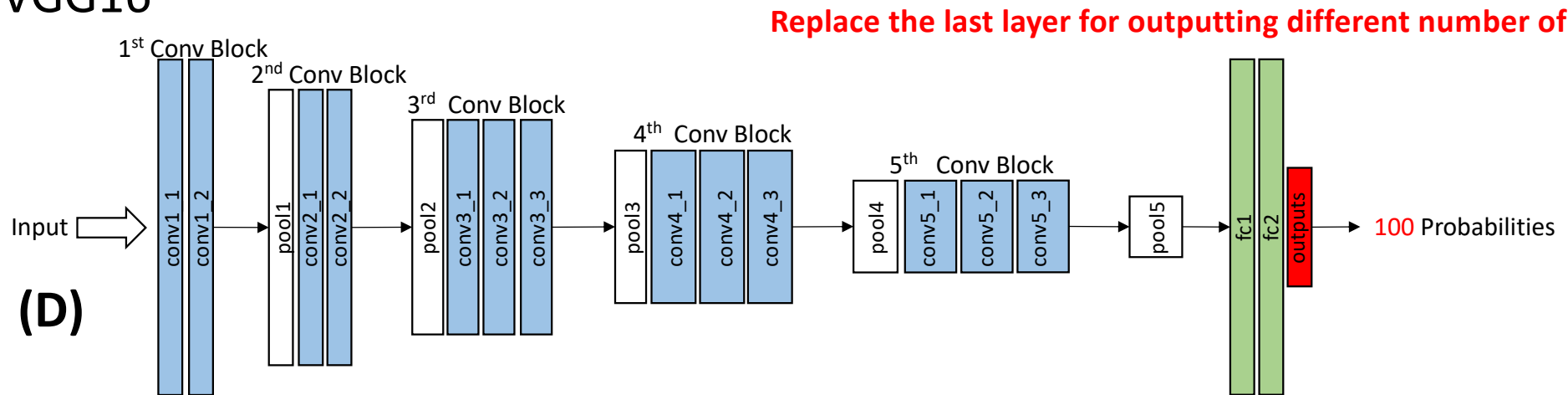
- VGG16



```
>>> vgg = tl.models.vgg16(pretrained=True, end_with='conv4_3', mode='static')
```

Convolutional Encoding Architectures

- VGG16



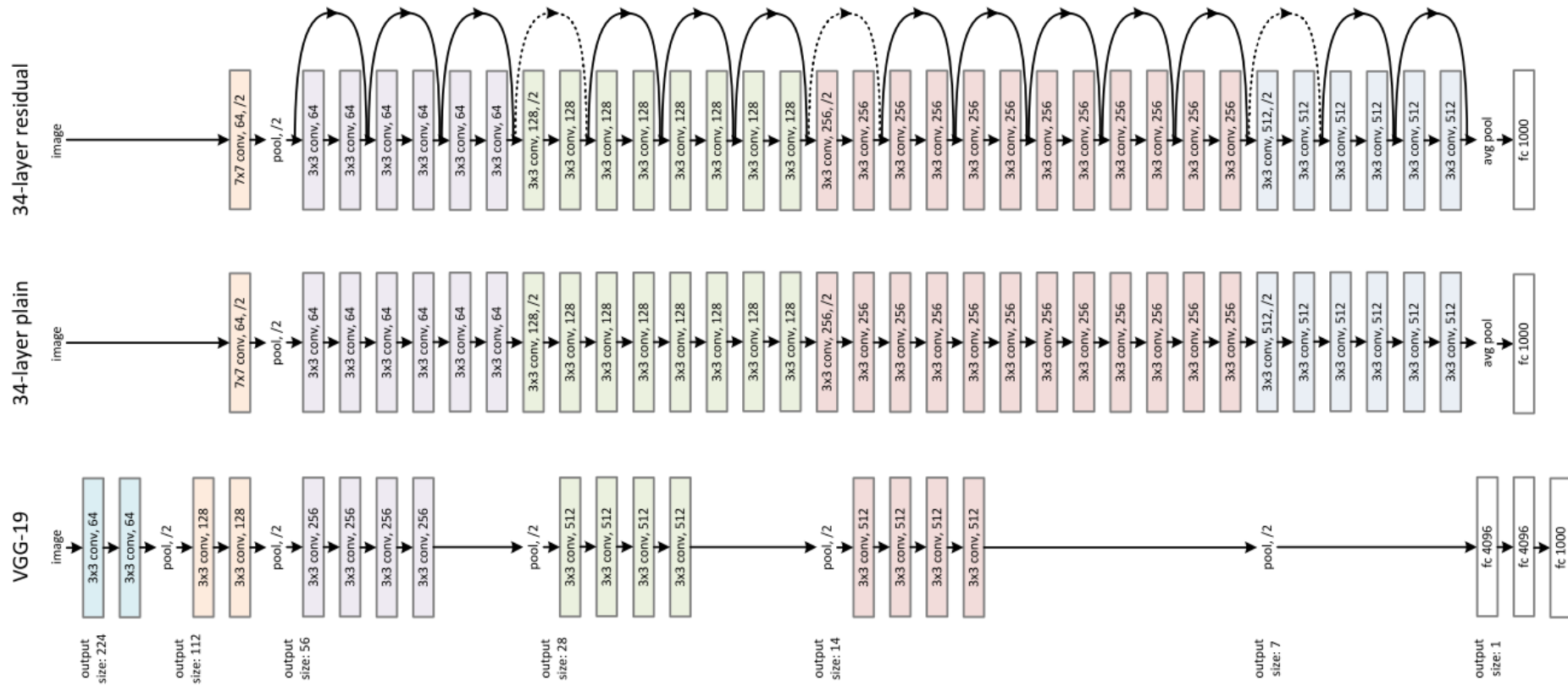
```

>>> # get VGG without the last layer
>>> cnn = tl.models.vgg16(end_with='fc2_relu', mode='static').as_layer()
>>> # add one more layer and build a new model
>>> ni = Input([None, 224, 224, 3], name="inputs")
>>> nn = cnn(ni)
>>> nn = tl.layers.Dense(n_units=100, name='out')(nn)
>>> model = tl.models.Model(inputs=ni, outputs=nn)
>>> # train your own classifier (only update the last layer)
>>> train_weights = model.get_layer('out').trainable_weights

```

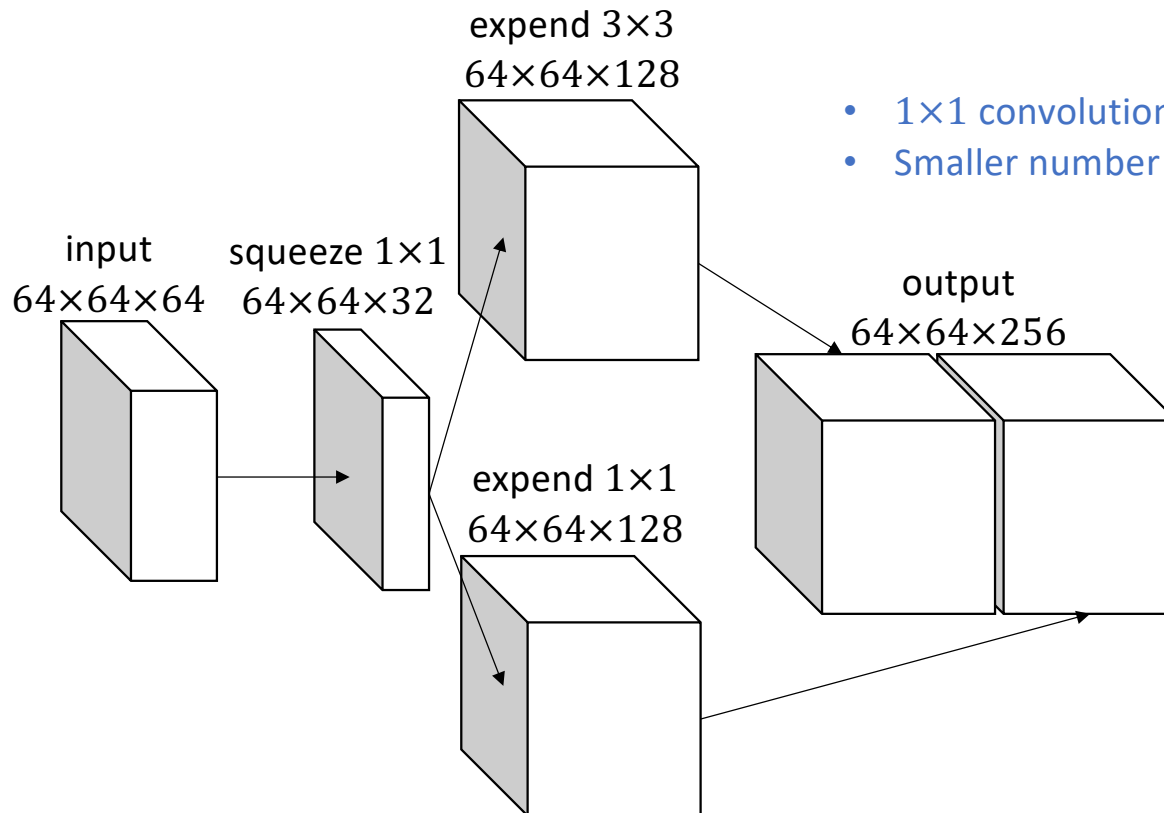
Convolutional Encoding Architectures

- ResNet



Convolutional Encoding Architectures

- SqueezeNet, MobileNet, ShuffleNet



- 1x1 convolution (pointwise conv) is 9 times faster than 3x3 convolution
- Smaller number of channels == Less computation

```
import time
import numpy as np
import tensorflow as tf
import tensorlayer as tl
from tensorlayer.models.imagenet_classes import class_names
```

```
squeezenet = tl.models.SqueezeNetV1(pretrained=True)
img1 = tl.vis.read_image('data/tiger.jpeg')
img1 = tl.prepro.imresize(img1, (224, 224)) / 255
img1 = img1.astype(np.float32)[np.newaxis, ...]
```

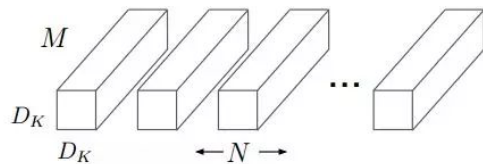
```
start_time = time.time()
output = squeezenet(img1, is_train=False)
prob = tf.nn.softmax(output)[0].numpy()
print(" End time : %.5s" % (time.time() - start_time))
preds = (np.argsort(prob)[::-1])[0:5]
for p in preds:
    print(class_names[p], prob[p])
```

Very Fast even on CPU

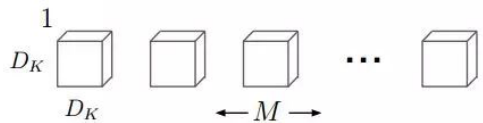
An example of "fire" block

Convolutional Encoding Architectures

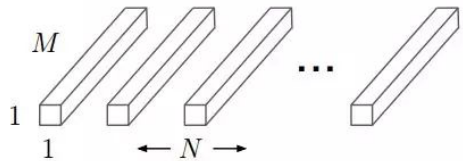
- SqueezeNet, MobileNet, ShuffleNet



(a) Standard Convolution Filters



(b) Depthwise Convolutional Filters



The number of multiplications of a **standard 2D CNN** layer without strides and padding:

- $filter_size_height \times filter_size_width \times height \times width \times in_channels \times output_channels$

The number of multiplications of a **depthwise 2D CNN** layer without strides and padding:

- $filter_size_height \times filter_size_width \times height \times width \times in_channels$

The number of multiplications of a **pointwise 2D CNN** layer without strides and padding:

- $height \times width \times in_channels \times output_channels$

(c) 1×1 Convolutional Filters called Pointwise Convolution in the con-
<http://blog.csdn.net/u011995719>

Convolutional Encoding Architectures

- SqueezeNet, MobileNet, ShuffleNet

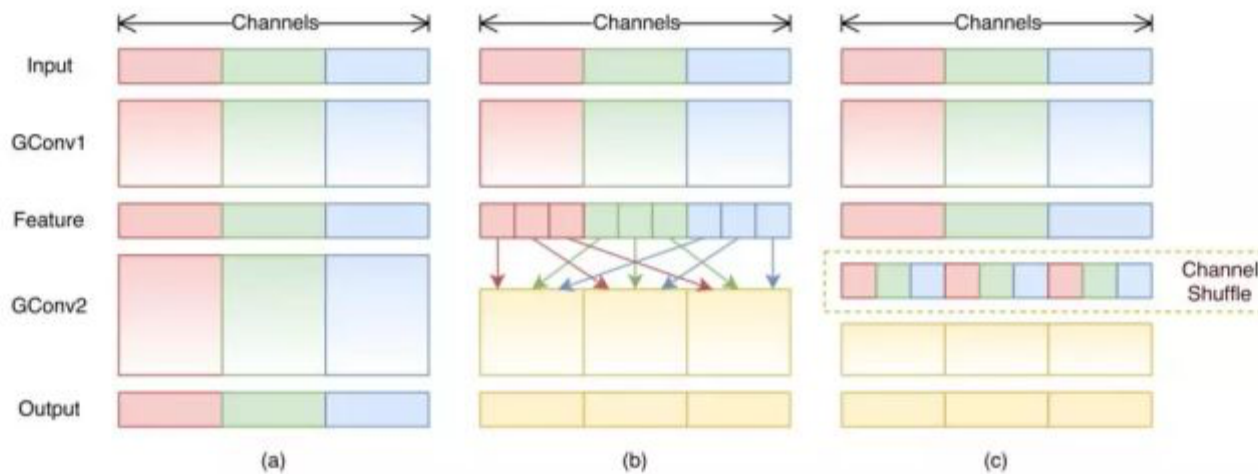


Figure 1. Channel shuffle with two stacked group convolutions. GConv stands for group convolution. a) two stacked convolution layers with the same number of groups. Each output channel only relates to the input channels within the group. No cross talk; b) input and output channels are fully related when GConv2 takes data from different groups after GConv1; c) an equivalent implementation to b) using channel shuffle.

Group Convolution: Split the channels into several “group” and perform standard convolution using different CNN layers.

Shuffle layer: Merge the information of different group by shuffling the channels.

Convolutional Encoding Architectures

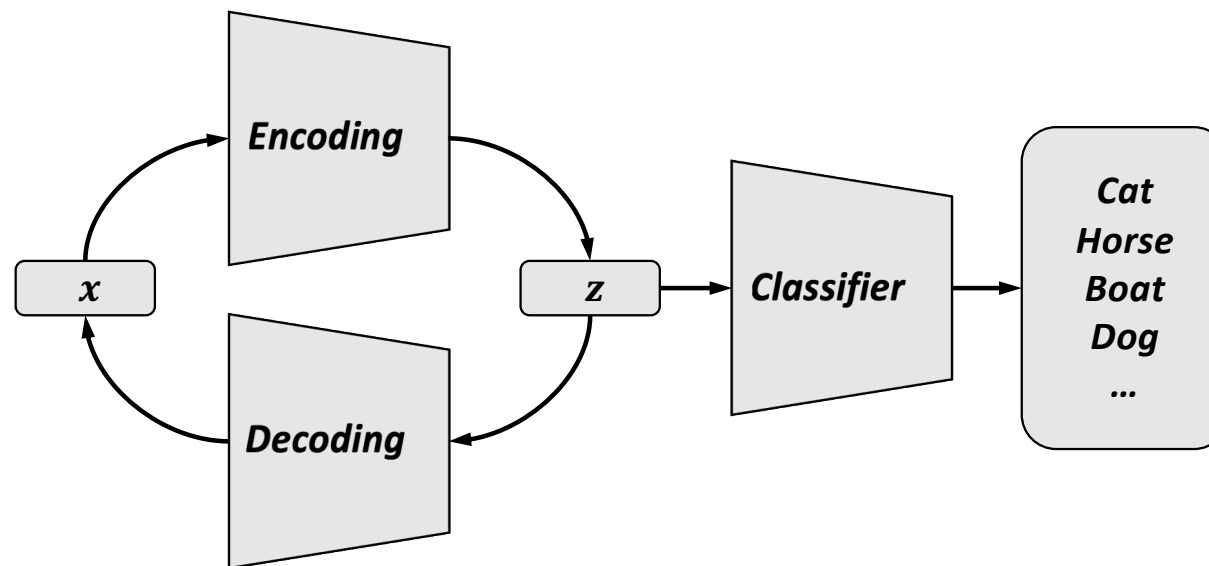
- SqueezeNet++
- MobileNetV2
- MobileNetV3
- ShuffleNetV2
- ...

Transposed Convolutional Algorithms

Transposed Convolutional Algorithms

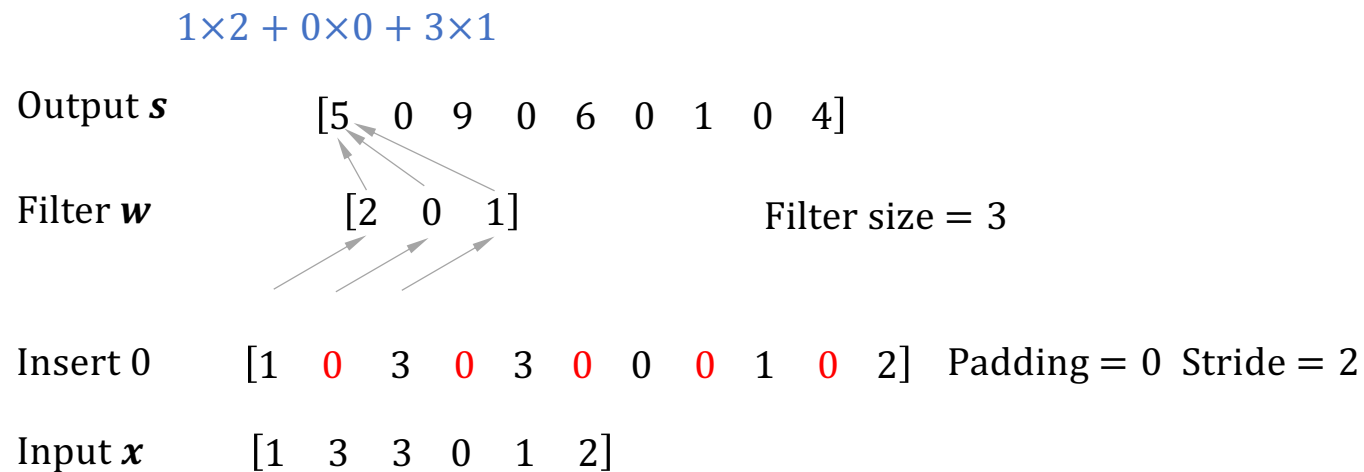
- Motivation

Convolution and pooling can only keep or reduce the size of the feature maps



Transposed Convolutional Algorithms

- Transposed convolution on 1D vector



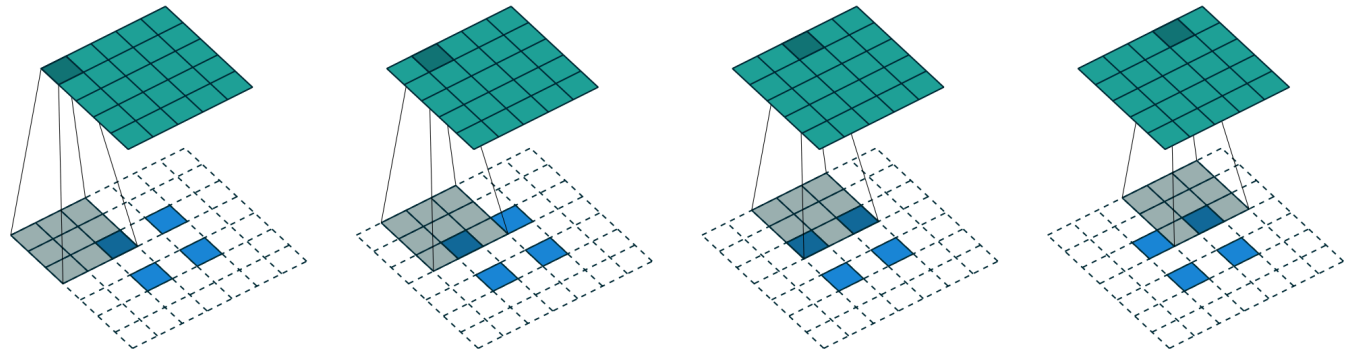
Transposed Convolutional Algorithms

- Inserting zeros between inputs

Padding, strides and transposed

Padding = 2×2

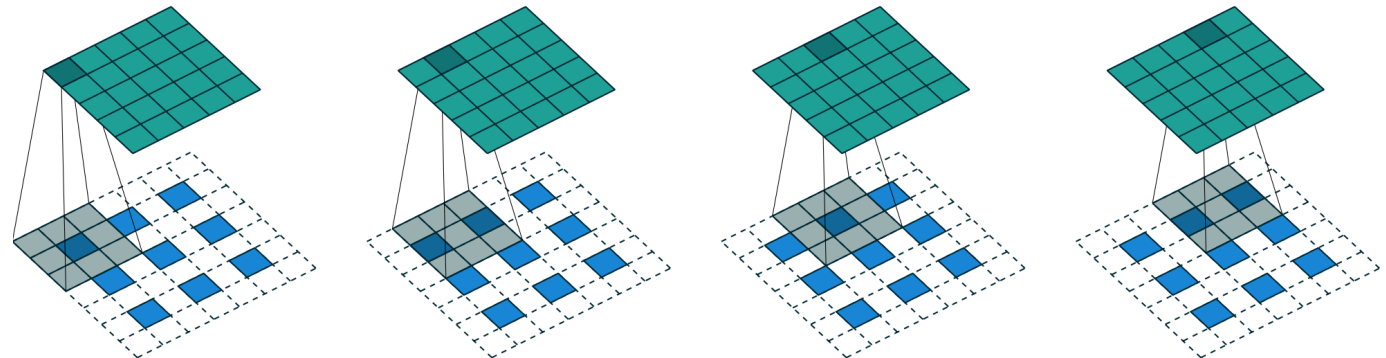
Strides = 2×2



Padding, strides and transposed

Padding = 1×1

Strides = 2×2



Transposed Convolutional Algorithms

- Resize convolution
- Subpixel convolution
- ...