# Recap

- Backpropagation
  - A modular way to compute the gradients of the loss w.r.t each parameter
  - Represent the computation using a graph
    - A node for each operation, and an edge for each variable
    - Each operation implements two functions: forward and backward
  - Apply chain rules against the graph
    - Forward the data through every node in topological order
    - Backward the gradient through every node in the reverse order

Slide credit: Wang Wei

# Recap

- Mini-batch stochastic gradient descent (SGD)
  - Reduces the chance of local optimal points and saddle points from GD
  - More stable than standard SGD
  - Extensions: Momentum, RMSProp, Adam
    - Exploiting historical updates
    - Adaptive learning rate per parameter
- Training tricks
  - Parameter initialization
  - Data normalization
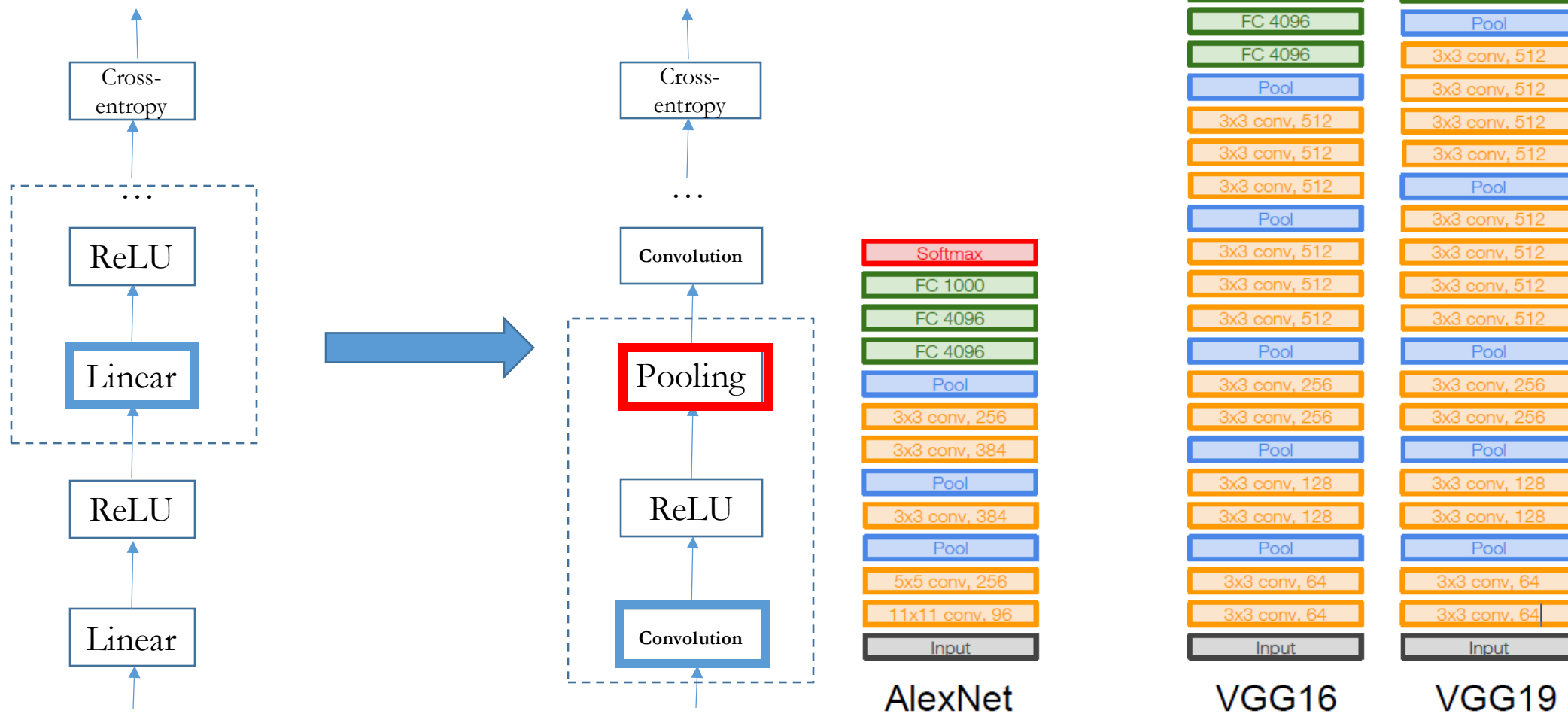  - Regularization:
    - Early stopping

Slide credit: Wang Wei

# Agenda

- Convolution basics

- 1D and 2D convolution operations

- Pooling operations

# Convolutional Neural Networks

- Referred to in short form as ConvNets or CNNs

- Most frequently used for image(-related data):
  - Image classification
  - Object detection
  - (Medical) image segmentation
  - [Face recognition](#)
  - [Image generation](#)
  - [Art composition](#)

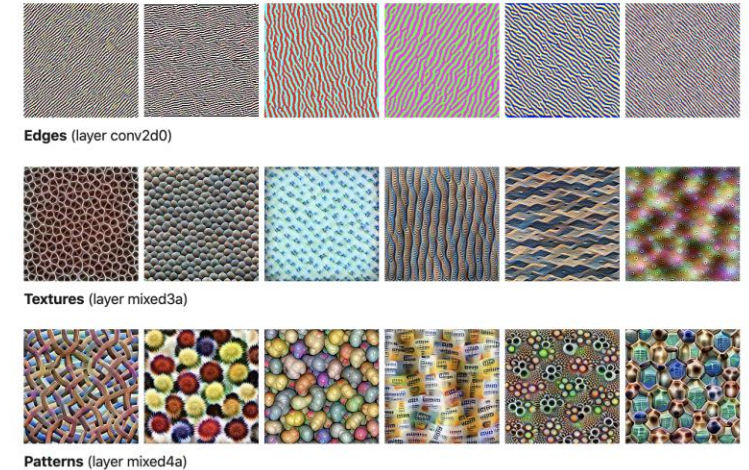Slide credit: Wang Wei

# From MLP to CNN

# Convolution

How can we get the most valuable information from the input (e.g. image)?

- A linear transformation

- Since output is a "feature" of the input, convolution can be considered feature extraction

- 1D / 2D / 3D
  - 1D: Text processing
  - 2D: Image processing
  - 3D: 3D data, CT, microscopy, etc.

- mD convolution
  - "mD" comes from "m" dimensions of the source data
  - Can apply 1D convolution over 1D or 2D data;     More easily, we can convert a matrix to a vector
  - Can apply 2D convolution over 2D or 3D data;     More easily, we can convert a tensor to a matrix



**Feature Visualization**
How neural networks build up their understanding of images

**Edges** (layer conv2d0)

**Textures** (layer mixed3a)

**Patterns** (layer mixed4a)

Feature visualization



Input → convolution → feature

Slide credit: Wang Wei

# 1D convolution
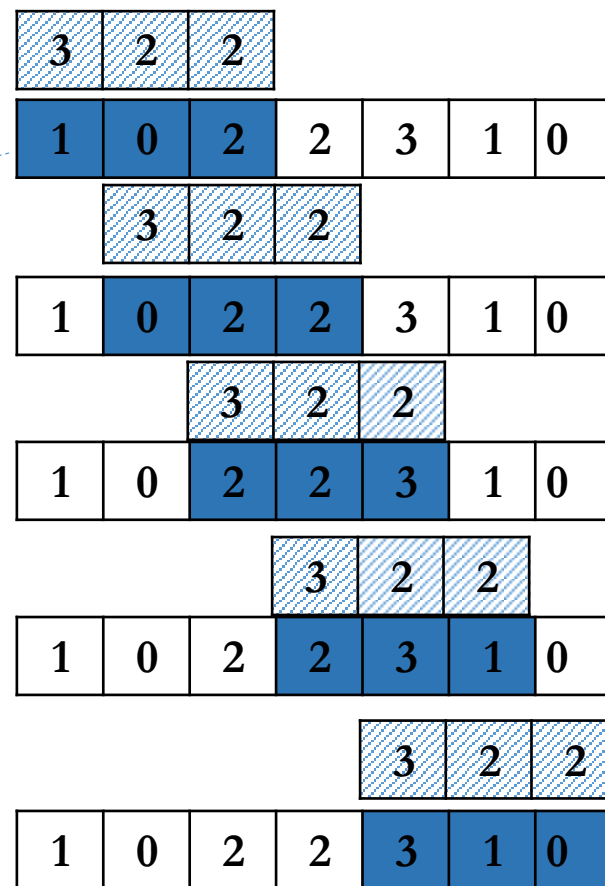
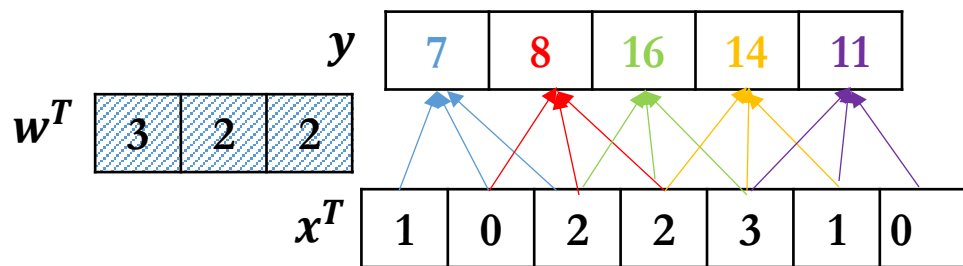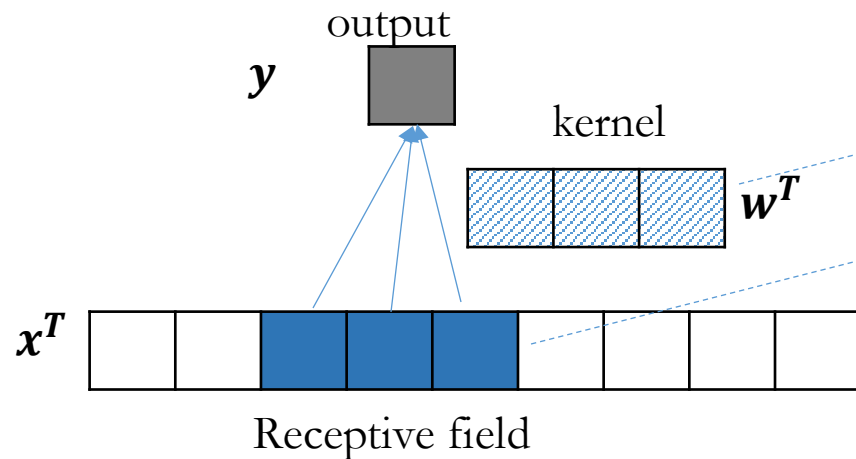$$y_t = \sum_{i=0}^{k-1} w_i \times x_{t+i}$$

Cross-correlation operation [link]

- **In CNNs, convolution refers to cross-correlation**

- $\boldsymbol{w}$ is called the kernel/filter; length $k$
  - "weights" or parameters to be trained;
- $\boldsymbol{x}$ is the input; length $n$
- the applied or input area, i.e. t, t+1, …, t+k-1 is called the receptive field
  - one receptive field generates one output value
- $\boldsymbol{y}_t$ is the output feature; length $o$
  - yt should be y, including all the output values

In signal processing, cross-correlation is a measure of similarity of two series as a function of the displacement of one relative to the other.
It is also known as a sliding dot product.

# 1D convolution



$3 \times 1 + 2 \times 0 + 2 \times 2 = 7$
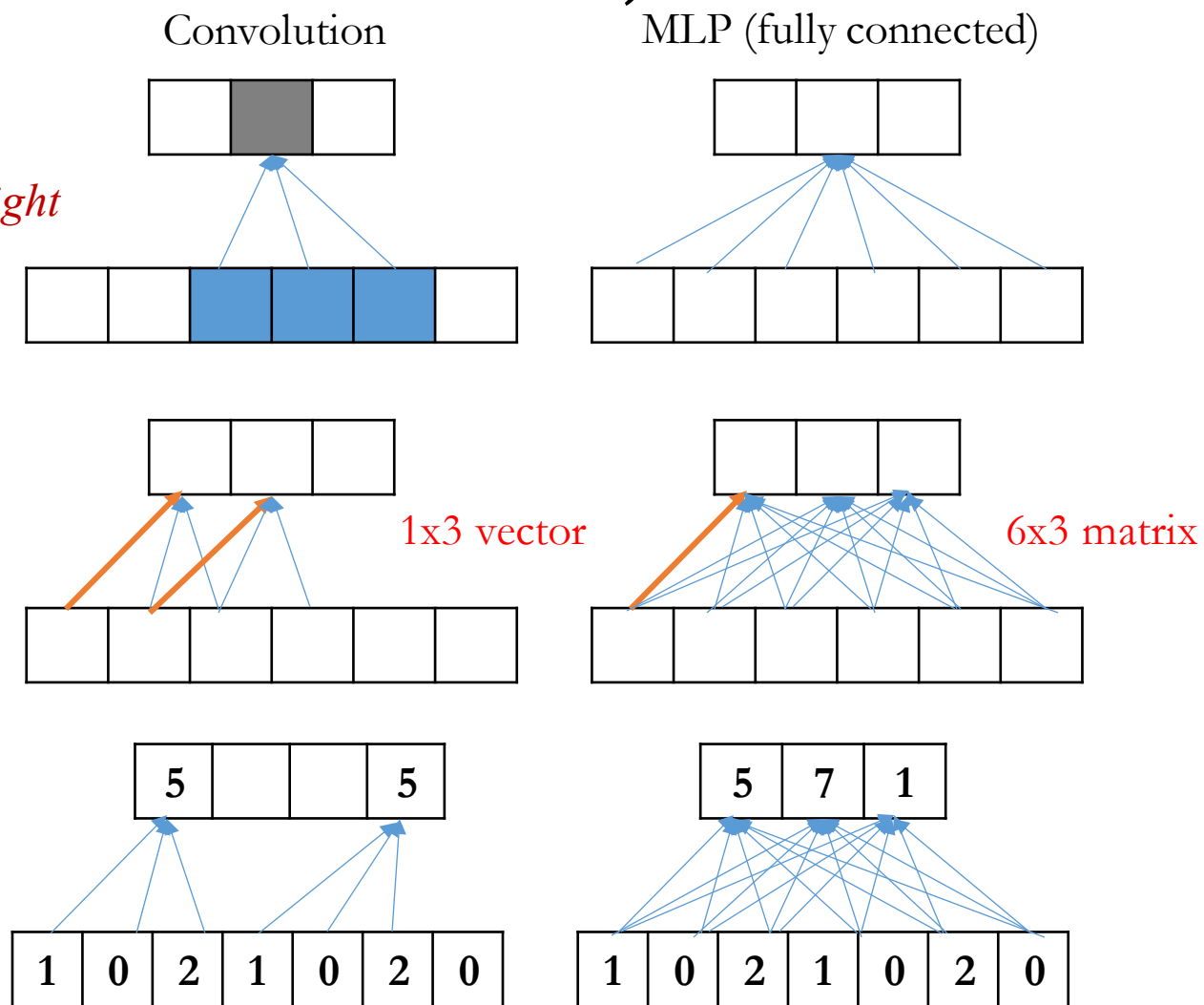
$3 \times 0 + 2 \times 2 + 2 \times 2 = 8$

$3 \times 2 + 2 \times 2 + 2 \times 3 = 16$

$3 \times 2 + 2 \times 2 + 2 \times 1 = 14$

$3 \times 3 + 2 \times 1 + 2 \times 0 = 11$

# Properties (Why is convolution better?)

- Sparse connection:
  - each output is connected only to inputs within receptive field vs. all inputs
  - → fewer parameters (each output needs 3 vs. 6 in example)
  - → less overfitting
- Weight sharing vs. unique weights
  - Regularization
  - Less overfitting
- Location or Spatial invariant
  - Function transformations should not depend on the location within the image, i.e.
  - Make the same prediction no matter where the object is in the image

Convolution          MLP (fully connected)

*output*

*kernel/weight*

*input*

1x3 vector          6x3 matrix

| 5 | | | 5 |

| 5 | 7 | 1 |

| 1 | 0 | 2 | 1 | 0 | 2 | 0 |

| 1 | 0 | 2 | 1 | 0 | 2 | 0 |

Slide credit: Wang Wei

# Location or Spatial Invariant

- You can recognize an object even its appearance varies in some way

- Convolution operator commutes with respect to translation
  - If you convolve f with g, it doesn't matter if you translate the convolved output f*g, or you translate f or g first, then convolve them.
  - https://en.wikipedia.org/wiki/Convolution

- Location or Spatial invariant
  - Function transformations should not depend on the location within the image, i.e.
  - Make the same prediction no matter where the object is in the image



**Translation Invariance**

**Rotation/Viewpoint Invariance**

**Size Invariance**

Slide credit: Matt Krause

# Perceptron, MLP and Convolution

| Perceptron | Perceptron is too simple |
|---|---|

Perceptron is too simple
→ underfitting
→ add more layers
→ MLP

| MLP | MLP has too many parameters |
|---|---|

MLP has too many parameters
→ High dimension
→ difficult to optimize and overfitting
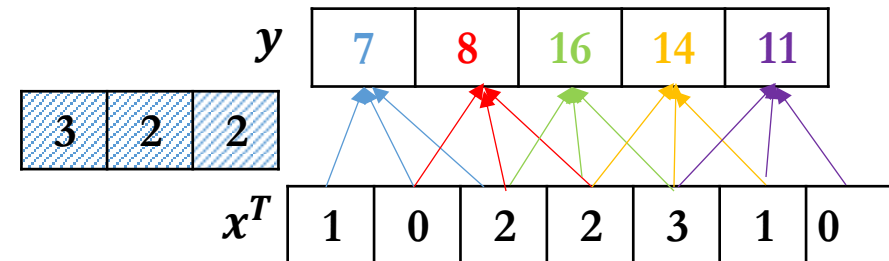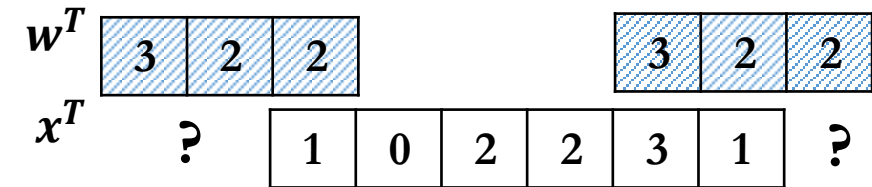→ CNN (with more regularization)

| CNN |
|---|

# Padding

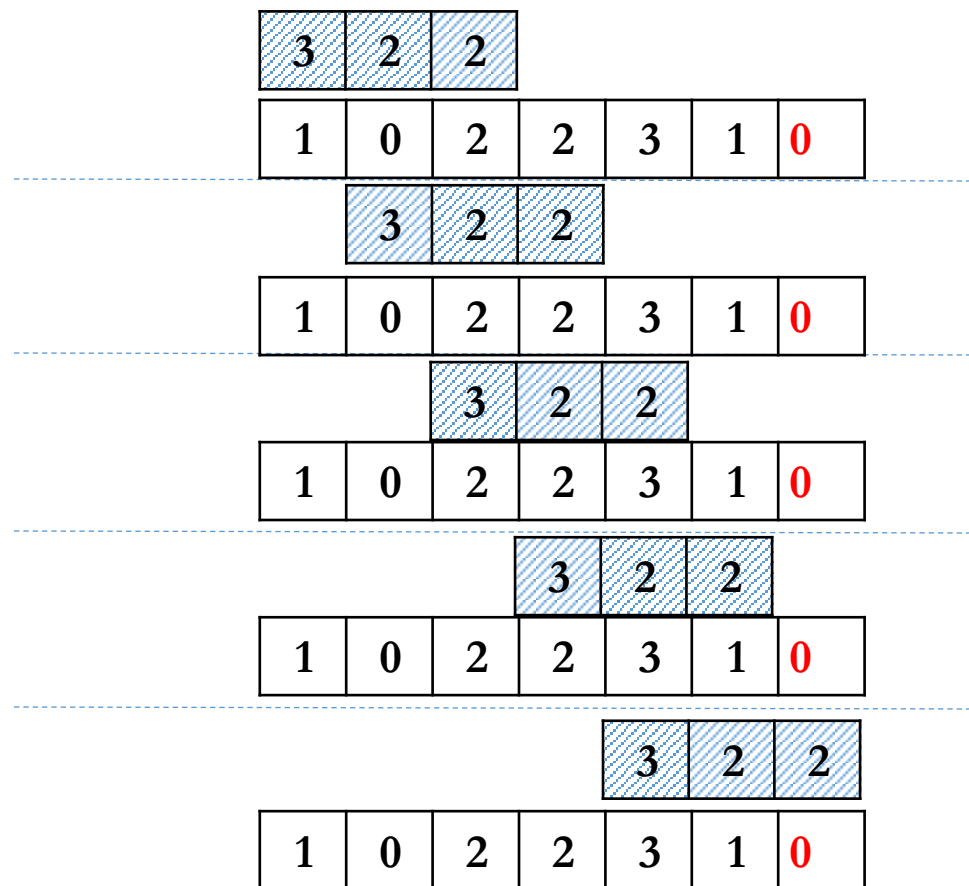- How to determine the edge values?
  Ignore non-valid regions?



We want the first element of input to be processed by all the parameters in the kernel.

- $o = n - k + 1$

  - n: input length, k: kernel length,
    o: output length
  - Output is shorter than input

- To retain the resolution/size

  - Pad with extra values (usually 0s)

# Padding

- Manual padding ($p$)
  - Output feature values for $p = 1$
    - 3x1+2x0+2x2=7
    - 3x0+2x2+2x2=8
    - 3x2+2x2+2x3=16
    - 3x2+2x2+2x1=14
    - 3x3+2x1+2x0=11
- What value to pad with?
  - usu. k << n so value doesn't matter too much; so don't bother tuning
  - 0 picked for convenience
- Operation supported in many deep learning libraries, e.g.
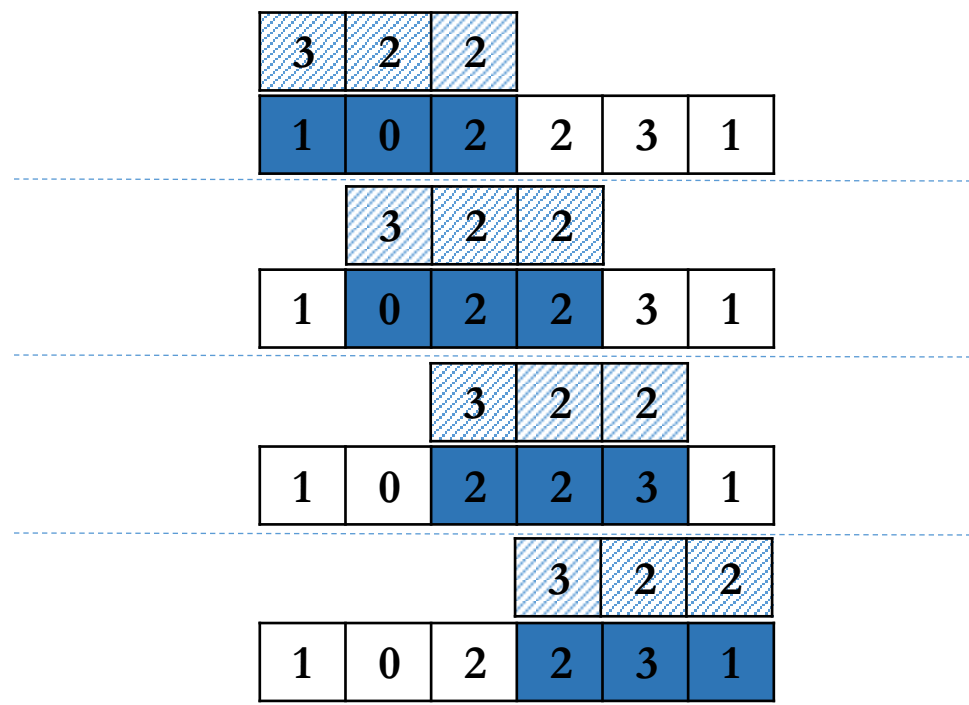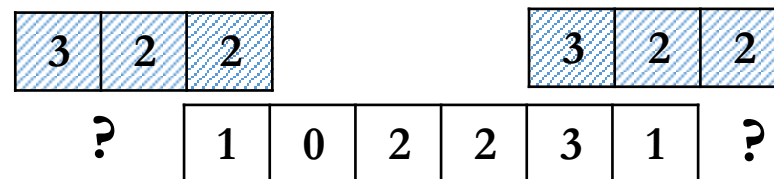  - Torch, PyTorch, Caffe, SINGA

# Padding amount (p)

- Given padding ($p$), what will the output length be?
- Kernel size/length: $k$, input length: $n$
- Without padding:
  - \# outputs $o = n - k + 1$
- With padding:
  - \# outputs $o = (n + p) - k + 1$
- padding length ($p$) can be set manually or automatically
- 2 special automatic settings:
  - consider only "valid" convolutions → $p = 0$
  - same length output as input → "same"
    - $o = n$ → $p = f(k)$    $p = k\text{-}1$

Slide credit: Wang Wei

# "Valid" Convolution



- No padding ($p = 0$)
  - # inputs denoted as $n$
  - # outputs $o = n - k + 1 = $ 6-3+1=4
  - Output feature values
    - 3x1+2x0+2x2=7
    - 3x0+2x2+2x2=8
    - 3x2+2x2+2x3=16
    - 3x2+2x2+2x1=14
  - Outputs become shorter
- Is an option to be set in library

# Same Padding

- Same padding ($p$?)
  - $o = n = n + p - k + 1$
  - $p = k - 1$
    - Left padding = $\lfloor p/2 \rfloor$
    - Right padding = $\lceil p/2 \rceil$
  - Output values
    - 3x0+2x1+2x0=2
    - 3x1+2x0+2x2=7
    - 3x0+2x2+2x2=8
    - 3x2+2x2+2x3=16
    - 3x2+2x3+2x1=14
    - 3x3+2x1+2x0=11

If p is an odd number, libraries typically assign 1 less to the left than right or vice versa.

For $n$ an integer, $\lfloor n \rfloor = \lceil n \rceil = [n] = n$.

**Examples**

| $x$ | Floor $\lfloor x \rfloor$ | Ceiling $\lceil x \rceil$ | Fractional part $\{x\}$ |
|---|---|---|---|
| 2 | 2 | 2 | 0 |
| 2.4 | 2 | 3 | 0.4 |
| 2.9 | 2 | 3 | 0.9 |
| −2.7 | −3 | −2 | 0.3 |
| −2 | −2 | −2 | 0 |

Slide credit: Wang Wei

# Stride

- How many steps to move towards the next receptive field
  - so far, we have considered only a stride of 1, where the kernel is applied directly to the next element in the array
  - s>1, skip some elements
    - Faster to compute
    - *Effective receptive field size increases quickly*
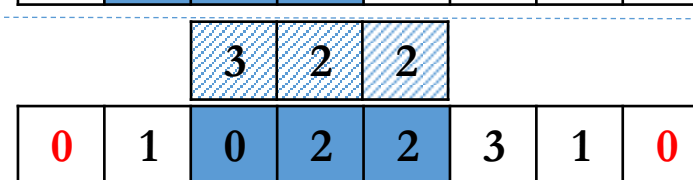
In s=2, the top element depends on the entire array (receptive field is 7), vs. 5 when s=1.

Layer 2

Layer 1

| | | | 3 | 2 | 1 | ? |
| | | | 3 | 2 | 1 | ? |

$w^T$ | 3 | 2 | 1 |
$x^T$ | 1 | 0 | 2 | 2 | 3 | 1 |

s=2

s=1

# Stride

- Increasing the stride is computationally faster, since we compute less convolutions

- Resulting output with higher stride is subsequently shorter

- $o = \left\lfloor \frac{n+p-k}{s} \right\rfloor + 1$

- Tunable hyperparameter, which may also vary depending on the layer

# Stride

- Exact matching
  - With padding p (=1)
  - $o = \left\lfloor \dfrac{n+p-k}{s} \right\rfloor + 1$
  - (6+1-3)/2+1=3

# Stride

- Not exact matching
  - With padding p (=2)
  - $o = \left\lfloor \dfrac{n+p-k}{s} \right\rfloor + 1$
  - (6+2-3)/2+1=3

This last computation is not valid.

This equation works in both situations.

Slide credit: Wang Wei

# Stride (for Tensorflow)

- When stride > 1
  - Valid padding, $p = 0$
  - Same padding, ?

- output length cannot be equal to the input, since a stride greater than 1 will shorten the output

- "same" is defined as:
  - $o = \left\lceil \dfrac{n}{s} \right\rceil, p = ?$

  ceiling operation.

Slide credit: Wang Wei

# Stride (for Tensorflow)

- When stride > 1
  - Valid padding, $p = 0$
  - Same padding, $o = \left\lceil \frac{n}{s} \right\rceil$, $p =?$

$$o = \left\lfloor \frac{n+p-k}{s} \right\rfloor + 1 \quad \text{<= (n+p-k)/s + 1}$$

$$\frac{n+p-k}{s} \geq o - 1$$

$$p \geq s(o-1) + k - n$$

Tensorflow internally computes this p value when we set "same" padding for strides greater than 1.

$$p = \max(s(o-1) + k - n, 0)$$

# Summary by an Example (2-minute Quiz)

- Input Length = 13; Stride = 5; Kernel Length = 6

- "valid" method: no padding, drop the non-valid region

  How many elements in input vector will be dropped?

- "same" method: padding at both ends

  $p = \max(s(o-1) + k - n, 0)$

  What is the padding length? Where?

# Summary by an Example

- Input Length = 13; Stride = 5; Kernel Length = 6

- "valid" method: no padding, drop the non-valid region

```
inputs:            1  2  3  4  5  6  7  8  9  10 11 (12 13)
                  |_____|                 dropped
                        |_____|
```

$p = \max(s(o-1) + k - n, 0)$

$p = \max(s(\text{ceiling}(n/s)-1) + k - n, 0)$

$p = \max(5(\text{ceiling}(13/5)-1) + 6 - 13, 0)$

- "same" method: padding at both ends

$p = \max(5(\text{ceiling}(2.6)-1) + 6 - 13, 0)$

```
             pad|                              |pad
inputs:      0 |1  2  3  4  5  6  7  8  9  10 11 12 13|0  0
               |_____|
                     |_____|
                           |_____|
```

$p = \max(5(3-1) + 6 - 13, 0)$

$p = \max(10 + 6 - 13, 0) = 3$

# Computing

- Conv1D
  - Forward(x, w)

$$w^T \quad \boxed{3 \quad 2 \quad 1}$$

$$x^T \quad \boxed{1 \quad 0 \quad 2 \quad 2 \quad 3 \quad 1}$$

S=2, k=3, p=1

| 1 | 0 | 2 | 2 | 3 | 1 | 0 |
|---|---|---|---|---|---|---|
| 1 | 0 | 2 | 2 | 3 | 1 | 0 |
| 1 | 0 | 2 | 2 | 3 | 1 | 0 |

Slide credit: Wang Wei

# Forward

$$w^T \quad \boxed{3 \mid 2 \mid 1}$$

$$x^T \quad \boxed{1 \mid 0 \mid 2 \mid 2 \mid 3 \mid 1}$$

S=2, k=3, p=1

Convolution via dot products in a for loop.

```
x_pad ← x
for t in range(o):
    y[t] = dot(w, x_pad[t*s:t*s+k])
```

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

receptive field to column

Vectorization (img2col)

More efficient vectorized implementation; convert each receptive field as a single column.

```
for t in range(o):
    X[:, t] = x_pad[t*s:t*s+k]
y = dot(w, X)
```

Try out in Google colab!

$$w^T \qquad k$$
$$\boxed{3 \mid 2 \mid 1}$$

$$k \quad \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 0 & 2 & 1 \\ \hline 2 & 3 & 0 \\ \hline \end{array} \quad \mathbf{X}$$

The result y is a row vector; if we compute y as $X^T w$, then the result y is a column vector

$dot()$

Converting many dot products to a vector-matrix multiplication

Slide credit: Wang Wei

# 2D Convolution

- 2D convolution follows the same principles, but the inputs, kernels and outputs are generalized into 2D matrices instead of 1D vectors



$$Y_{i,j} = \sum_{a=0}^{k_h-1} \sum_{b=0}^{k_w-1} X_{i*s_h+a,\,j*s_w+b} \times W_{a,b}$$

The operation between the receptive field and the kernel.

Is this a matrix-matrix multiplication?

Slide credit: Wang Wei

# 2D Convolution

- 2D convolution follows the same principles, but the inputs, kernels and outputs are generalized into 2D matrices instead of 1D vectors



No. It is not a matrix-matrix multiply.

$$Y_{i,j} = \sum_{a=0}^{k_h-1} \sum_{b=0}^{k_w-1} X_{i*s_h+a,\, j*s_w+b} \times W_{a,b}$$

Is this a matrix-matrix multiplication?

Element-wise operation: O(n^2) not O(n^3)

Slide credit: Wang Wei

# 2D Convolution

**Y**

**X**



k=3, p=0, s=1 (Valid)          k=3, p=2, s=1 (Same)          k=3, p=4, s=1 (Full)          k=3, p=2, s=2

Source: http://deeplearning.net/software/theano/tutorial/conv_arithmetic.html

# 2D Convolution

http://setosa.io/ev/image-kernels/

In deep learning, we learn the kernels or weights which gives us good predictions e.g. classification, regression etc.

| -1 | 1 |
|----|---|
| -1 | 1 |

vertical gradients

| 1 | 1 |
|---|---|
| 1 | 1 |

We get different feature maps depending on the kernel used.

(v. slight) smoothing

# 2D Convolution

- $Y_{i,j} = \sum_{a=0}^{k_h-1} \sum_{b=0}^{k_w-1} X_{i*s_h+a,j*s_w+b} \times W_{a,b}$



| 1 | 2 | 3 | 2 |
|---|---|---|---|
| 2 | 3 | 2 | 5 |
| 1 | 2 | 2 | 1 |
| 0 | 2 | 1 | 2 |

| 1 | 2 | 3 | 2 |
|---|---|---|---|
| 2 | 3 | 2 | 5 |
| 1 | 2 | 2 | 1 |
| 0 | 2 | 1 | 2 |

| -1 | 1 |
|----|---|
| -1 | 1 |

| 2 | | |
|---|---|---|
| | | |
| | | |

# 2D Convolution

$$Y_{i,j} = \sum_{a=0}^{k_h-1} \sum_{b=0}^{k_w-1} X_{i*s_h+a, j*s_w+b} \times W_{a,b}$$



| 1 | 2 | 3 | 2 |
|---|---|---|---|
| 2 | 3 | 2 | 5 |
| 1 | 2 | 2 | 1 |
| 0 | 2 | 1 | 2 |

| 1 | 2 | 3 | 2 |
|---|---|---|---|
| 2 | 3 | 2 | 5 |
| 1 | 2 | 2 | 1 |
| 0 | 2 | 1 | 2 |

| -1 | 1 |
|----|---|
| -1 | 1 |

→

| 2 | | |
|---|---|---|
| | | |
| | | |

| 1 | 2 | 3 | 2 |
|---|---|---|---|
| 2 | 3 | 2 | 5 |
| 1 | 2 | 2 | 1 |
| 0 | 2 | 1 | 2 |

| -1 | 1 |
|----|---|
| -1 | 1 |

→

| 2 | 0 | |
|---|---|---|
| | | |
| | | |

| 1 | 2 | 3 | 2 |
|---|---|---|---|
| 2 | 3 | 2 | 5 |
| 1 | 2 | 2 | 1 |
| 0 | 2 | 1 | 2 |

| -1 | 1 |
|----|---|
| -1 | 1 |

→

| 2 | 0 | 2 |
|---|---|---|
| | | |
| | | |

...

| 1 | 2 | 3 | 2 |
|---|---|---|---|
| 2 | 3 | 2 | 5 |
| 1 | 2 | 2 | 1 |
| 0 | 2 | 1 | 2 |

| -1 | 1 |
|----|---|
| -1 | 1 |

→

| 2 | 0 | 2 |
|---|---|---|
| 2 | -1 | 2 |
| 3 | -1 | 0 |

Slide credit: Wang Wei

# Implementation

- Img2Col
  - Convert each receptive field into a column

- $Y_{i,j} = \sum_{a=0}^{k_h-1} \sum_{b=0}^{k_w-1} X_{i*s_h+a,\,j*s_w+b} \times W_{a,b}$

**W**

| -1 | 1 |
|----|---|
| -1 | 1 |

| -1 | 1 | -1 | 1 |
|----|---|----|---|

$k_h k_w$

$\hat{X}$

$o_h o_w$

| 1 | 2 | | 2 |
|---|---|---|---|
| 2 | 3 | … | 1 |
| 2 | 3 | | 1 |
| 3 | 2 | | 2 |

$k_h k_w$

**X**

| 1 | 2 | 3 | 2 |
|---|---|---|---|
| 2 | 3 | 2 | 5 |
| 1 | 2 | 2 | 1 |
| 0 | 2 | 1 | 2 |

| 1 | 2 | 3 | 2 |
|---|---|---|---|
| 2 | 3 | 2 | 5 |
| 1 | 2 | 2 | 1 |
| 0 | 2 | 1 | 2 |

…

| 1 | 2 | 3 | 2 |
|---|---|---|---|
| 2 | 3 | 2 | 5 |
| 1 | 2 | 2 | 1 |
| 0 | 2 | 1 | 2 |

$dot()$

$1 \times o_h o_w$    Vector to matrix

$o_h$   $o_w$

**Convolution -> affine transformation**

Slide credit: Wang Wei

# Statistics

- Shape:
  - Input $\boldsymbol{X} \in R^{n_h \times n_w}$
  - Kernel $\boldsymbol{W} \in R^{k_h \times k_w}$
  - Output $\mathbf{Y} \in R^{o_h \times o_w}$

- Parameter size
  - $k_h \times k_w$

- Output shape
  - $(o_h, o_w) = (\left\lfloor \frac{n_h + p_h - k_h}{s_h} \right\rfloor + 1, \left\lfloor \frac{n_w + p_w - k_w}{s_w} \right\rfloor + 1)$

- Computation cost
  - $O(k_h \times k_w \times o_h \times o_w)$ (float multiplication ops, FLOP)

Operations for each output or each Kernel computation is K_h x K_w

# 2D Convolution

- Multiple kernels/filters



$$Y_{l,i,j} = \sum_{a=0}^{k_h-1} \sum_{b=0}^{k_w-1} X_{i*s_h+a, j*s_w+b} \times W_{l,a,b} , l \in [0, c_o)$$

*Please note, $c_o$ is not C0*

# Implementation

$$Y_{l,i,j} = \sum_{a=0}^{k_h-1} \sum_{b=0}^{k_w-1} X_{i*s_h+a,j*s_w+b} \times W_{l,a,b}, \, l \in [0, c_o)$$

$X$

| 1 | 2 | 3 | 2 |
|---|---|---|---|
| 2 | 3 | 2 | 5 |
| 1 | 2 | 2 | 1 |
| 0 | 2 | 1 | 2 |

$\hat{X}$

$o_h o_w$

| 1 | 2 | 3 | 2 |
|---|---|---|---|
| 2 | 3 | 2 | 5 |
| 1 | 2 | 2 | 1 |
| 0 | 2 | 1 | 2 |

$W$

| -1 | 1 |
|----|---|
| -1 | 1 |

$c_o$ kernels

| 1 | 1 |
|---|---|
| 1 | 1 |

| 0 | 1 |
|---|---|
| 0 | 1 |

$k_h k_w$

| -1 | 1 | -1 | 1 |
|----|---|----|---|
| 1 | 1 | 1 | 1 |
| 0 | 1 | 0 | 1 |

$c_o$

$k_h k_w$

| 1 | 2 | 2 |
|---|---|---|
| 2 | 3 | 1 |
| 2 | 3 | 1 |
| 3 | 2 | 2 |

$\cdots$

$\cdots$

| 1 | 2 | 3 | 2 |
|---|---|---|---|
| 2 | 3 | 2 | 5 |
| 1 | 2 | 2 | 1 |
| 0 | 2 | 1 | 2 |

$dot()$

Vector to matrix

$Y$ $o_w$ $o_h$ $c_o$

**Convolution -> affine transformation**

# Statistics

- applying multiple kernels (filters) $c_o$, all of the same stride and padding
- Parameter size
    - $c_o \times k_h \times k_w$
- Output shape
    - $(c_o, o_h, o_w) = (c_o, \left\lfloor \frac{n_h + p_h - k_h}{s_h} \right\rfloor + 1, \left\lfloor \frac{n_w + p_w - k_w}{s_w} \right\rfloor + 1)$
- Computation cost
    - O( $(c_o \times k_h \times k_w) \times (o_h \times o_w)$ )  (float multiplication ops, FLOP)

Slide credit: Wang Wei

# 2D Convolution

With multiple ($c_i$) input channels and kernels (filters)



$k_h$

$c_i$ $k_w$

$n_h$

$o_h$

$o_w$

$n_w$

$c_o$

$c_i$

$c_o$ filters/kernels

** The convolution results across all input channels are summed.

$$Y_{l,i,j} = \sum_{d=0}^{c_i-1} \sum_{a=0}^{k_h-1} \sum_{b=0}^{k_w-1} X_{d,i+a,j+b} \times W_{l,d,a,b} + b_l, l \in [0, c_o)$$

$\boldsymbol{b}$ is a bias vector

l,i,j have noting to do with C_i

** this is still a 2D convolution because the kernel is moved only across the horizontal and vertical dimensions (as indexed by a, b).

# 2D Convolution (w/ 3D kernels and data)



Input Volume (+pad 1) (7x7x3)

x[:,:,0]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 1 | 1 | 0 |
| 0 | 1 | 2 | 1 | 1 | 0 | 0 |
| 0 | 0 | 2 | 1 | 2 | 2 | 0 |
| 0 | 2 | 2 | 0 | 2 | 2 | 0 |
| 0 | 0 | 1 | 2 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,1]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 0 | 2 | 0 |
| 0 | 0 | 1 | 1 | 1 | 0 | 0 |
| 0 | 2 | 0 | 1 | 2 | 0 | 0 |
| 0 | 0 | 0 | 1 | 1 | 2 | 0 |
| 0 | 0 | 0 | 1 | 0 | 2 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,2]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 0 | 1 | 1 | 0 | 2 | 0 |
| 0 | 0 | 0 | 2 | 1 | 2 | 0 |
| 0 | 0 | 0 | 2 | 2 | 2 | 0 |
| 0 | 1 | 1 | 2 | 1 | 2 | 0 |
| 0 | 2 | 2 | 1 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Filter W0 (3x3x3)

w0[:,:,0]

| 1 | 1 | 0 |
|---|---|---|
| 0 | -1 | -1 |
| -1 | -1 | 1 |

w0[:,:,1]

| 0 | -1 | -1 |
|---|---|---|
| -1 | 0 | -1 |
| 0 | 0 | -1 |

w0[:,:,2]

| 0 | 1 | 0 |
|---|---|---|
| -1 | 1 | 0 |
| -1 | 0 | 1 |

Bias b0 (1x1x1)

b0[:,:,0]

| 1 |
|---|

Filter W1 (3x3x3)

w1[:,:,0]

| -1 | -1 | 0 |
|---|---|---|
| 0 | 1 | -1 |
| 0 | -1 | -1 |

w1[:,:,1]

| 0 | -1 | 0 |
|---|---|---|
| 1 | 0 | 0 |
| 0 | 0 | 0 |

w1[:,:,2]

| 1 | 1 | 1 |
|---|---|---|
| 0 | -1 | -1 |
| 0 | -1 | -1 |

Bias b1 (1x1x1)

b1[:,:,0]

| 0 |
|---|

Output Volume (3x3x2)

o[:,:,0]

| 0 | -3 | 0 |
|---|---|---|
| 0 | 0 | -5 |
| 5 | 0 | 5 |

o[:,:,1]

| -3 | -6 | -3 |
|---|---|---|
| -9 | -11 | 0 |
| -5 | 2 | -3 |

Demo [link]
Visualization [link]

toggle movement

Slide credit: Wang Wei

# Implementation

- Forward
  - Convert input feature maps $\boldsymbol{X}$ into matrix $\hat{\boldsymbol{X}}$ (img2col) of size $(c_i k_w k_h \times o_h o_w)$
  - Reshape the filters $\boldsymbol{W}$ to $(c_o \times c_i k_h k_w)$
  - $\boldsymbol{Y} = \boldsymbol{W}\hat{\boldsymbol{X}} + \boldsymbol{b}$
  - Computational cost, $O(c_o \times c_i \times k_w \times k_h \times o_h \times o_w)$

- Backward
  - Given $\dfrac{\partial L}{\partial \boldsymbol{Y}}$
  - Compute $\dfrac{\partial L}{\partial \boldsymbol{W}} = \dfrac{\partial L}{\partial \boldsymbol{Y}} \hat{\boldsymbol{X}}^T$ , $\dfrac{\partial L}{\partial \hat{\boldsymbol{X}}} = \boldsymbol{W}^T \dfrac{\partial L}{\partial \boldsymbol{Y}}$
  - Column to receptive field transformation to get gradient wrt original X

# Img2Col

- Slide the window from left to right, top to bottom

- Copy the values from the receptive field into a column of $\hat{X}$
  - Receptive fields across feature maps are concatenated into to one column

- Reshape $\hat{X}$ into $(c_i k_h k_w \times o_h o_w)$

$$o_h o_w = 2 \cdot 3$$

Feature map 0

| 1 | 2 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| 2 | 2 | 3 | 1 | 2 | 0 |
| 1 | 1 | 2 | 1 | 0 | 1 |
| 2 | 1 | 2 | 1 | 1 | 3 |

Feature map 1

| 0 | 2 | 3 | 2 | 0 | 1 |
|---|---|---|---|---|---|
| 1 | 0 | 1 | 1 | 2 | 0 |
| 2 | 0 | 1 | 1 | 1 | 1 |
| 1 | 1 | 2 | 1 | 1 | 2 |

$c_i k_h k_w$

$= 2 \cdot 2 \cdot 2$

| 1 | 1 | 3 | 1 | 2 | 0 |
|---|---|---|---|---|---|
| 2 | 2 | 4 | 1 | 1 | 1 |
| 2 | 3 | 2 | 2 | 2 | 1 |
| 2 | 1 | 0 | 1 | 1 | 3 |
| 0 | 3 | 0 | 2 | 1 | 1 |
| 2 | 2 | 1 | 0 | 1 | 1 |
| 1 | 1 | 2 | 1 | 2 | 1 |
| 0 | 1 | 0 | 1 | 1 | 2 |

# 2D convolution

Filter 1

$c_i\ k_h\ k_w$

Filter 2

...

Reshape

**W**

$c_i\ k_h\ k_w$

$c_o$

...

$dot()$

Matrix to tensor

$\hat{X}$

$o_h o_w$

| 0 | |
|---|---|
| 3 | |
| 1 | |
| 4 | |
| 1 | |
| 0 | |
| 1 | |
| 1 | |
| 1 | |
| 2 | |
| 2 | |
| 3 | |

$c_i\ k_h\ k_w$

...

**X**

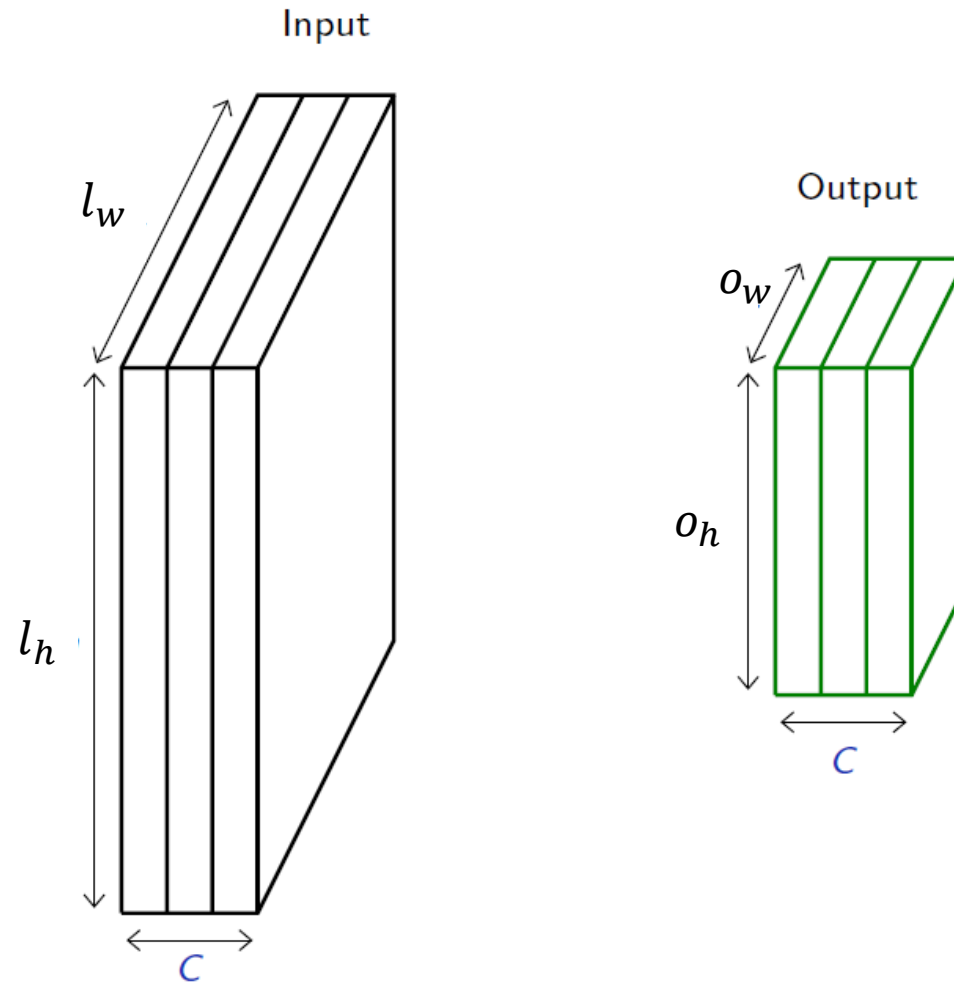| 1 | 2 | 3 | 2 |
|---|---|---|---|
| 2 | 3 | 2 | 5 |
| 1 | 2 | 2 | 1 |
| 0 | 2 | 1 | 2 |

$o_w$
$o_h$
$c_o$

# Statistics

- Parameter size
  - Weights: $c_o \times (c_i \; k_h \; k_w)$
  - Bias: $c_o$

- Output shape
  - $(c_o, o_h, o_w) = (c_o, \left\lfloor \frac{n_h + p_h - k_h}{s_h} \right\rfloor + 1, \left\lfloor \frac{n_w + p_w - k_w}{s_w} \right\rfloor + 1)$

- Computation cost
  - $O(c_o \times \; c_i \; k_h \; k_w \; \times \; o_h \; o_w)$
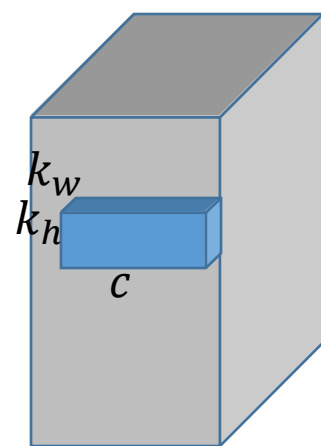
Slide credit: Wang Wei

# Pooling

- Aggregate information from each receptive field
  - Max
  - Average

- No parameters

- Applied for each channel respectively
  - #input channels = # output channels, i.e., $c_i = c_o = c$

- Padding and stride can be applied

Slide credit: Wang Wei

# Pooling Visualization



Input

$l_w$

$l_h$

$C$

Output

$o_w$

$o_h$

$C$

Slide credit: Wang Wei

# Max Pooling

$c$ feature maps (take <u>max</u> of each one)

| 1 | 2 | 1 |
|---|---|---|
| **3** | 1 | 1 |
| -1 | 0 | 0 |

| 1 | 1 | 1 |
|---|---|---|
| **2** | 0 | 1 |
| 1 | 0 | 0 |

$\cdots$

| 0 | -2 | 1 |
|---|---|---|
| 0 | 1 | **1** |
| 0 | 0 | 0 |

3          2          1

| 3 | 2 | $\cdots$ | 1 |
|---|---|---|---|

$o_h$

$o_w$

$c$

$k_w$

$k_h$

$c$

# Average Pooling

*c* feature maps (take <u>average</u> of each feature map)



| 1 | 2 | 1 |
|---|---|---|
| 3 | 1 | 1 |
| -1 | 0 | 0 |

8/9

| 1 | 1 | 1 |
|---|---|---|
| 2 | 0 | 1 |
| 1 | 0 | 0 |

7/9

...

| 0 | -2 | 1 |
|---|---|---|
| 0 | 1 | 1 |
| 0 | 0 | 0 |

1/9

Often, in neural networks that use average pooling, the pooling is "skipped" because it can be included directly into the convolution. If you are given original kernel weights **w**, what are the equivalent weights **w\*** which incorporate average pooling directly?

# Average Pooling

$k_w$
$k_h$

$l_h$

$l_w$

$c$

$c$ feature maps (take <u>average</u> of each feature map)

| 1 | 2 | 1 |
|---|---|---|
| 3 | 1 | 1 |
| -1 | 0 | 0 |

| 1 | 1 | 1 |
|---|---|---|
| 2 | 0 | 1 |
| 1 | 0 | 0 |

...

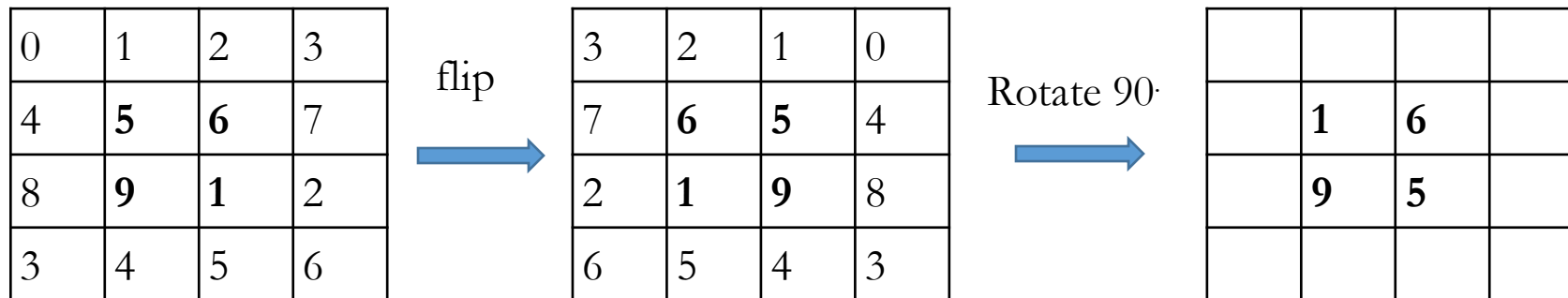| 0 | -2 | 1 |
|---|---|---|
| 0 | 1 | 1 |
| 0 | 0 | 0 |

8/9

7/9

1/9

...

$o_h$

$o_w$

$c$

Equivalent to convolution with kernel values = 1/kernel_size

Often, in neural networks that use average pooling, the pooling is "skipped" because it can be directly into the convolution. If you are given original kernel weights **w**, what are the equivalent weights **w*** which incorporate average pooling directly?

# Effect of Pooling

- Reduces the feature size and model size

- Information aggregation
    - Max pooling: invariant to rotation of the **input** image
    - Average pooling: can be replaced by convolution; much cheaper (no weights)

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | **5** | **6** | 7 |
| 8 | **9** | **1** | 2 |
| 3 | 4 | 5 | 6 |

flip

| 3 | 2 | 1 | 0 |
|---|---|---|---|
| 7 | **6** | **5** | 4 |
| 2 | **1** | **9** | 8 |
| 6 | 5 | 4 | 3 |

Rotate 90

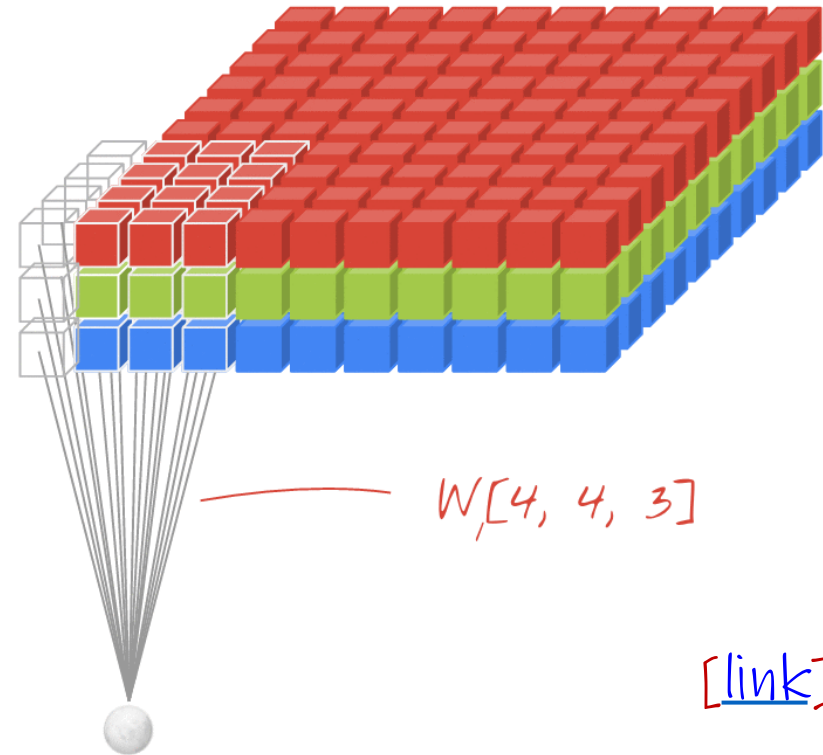|  |  |  |  |
|---|---|---|---|
|  | **1** | **6** |  |
|  | **9** | **5** |  |
|  |  |  |  |

# Multiple input channels, multiple filters

Configuration?

2-minute quiz:

How many input channels?
How many kernels are there?
What is the kernel size?
What is the padding?
What is the stride?

W[4, 4, 3]
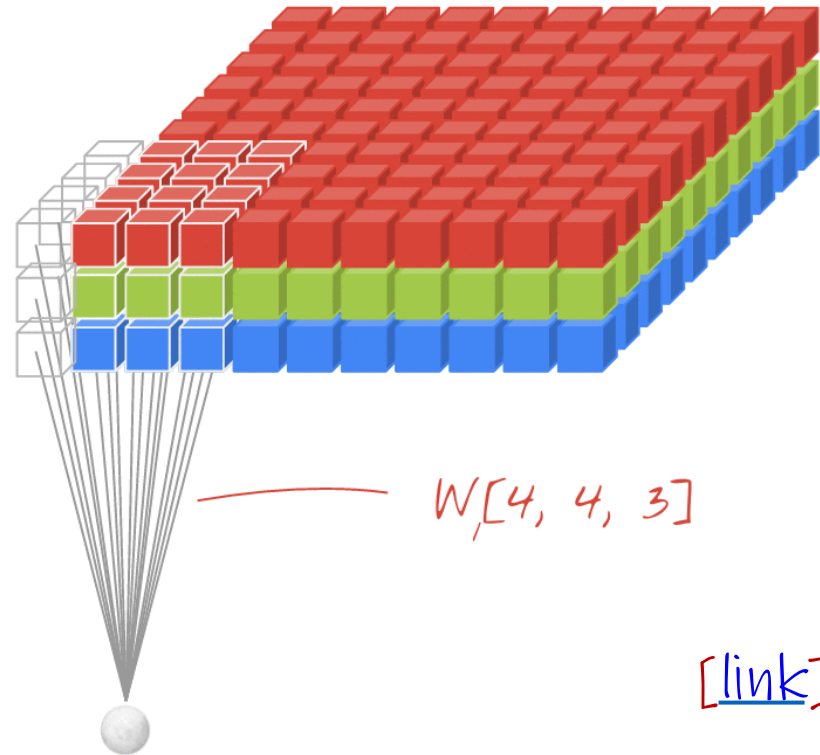
[link]

# Multiple input channels, multiple filters

Configuration?

Kernel size: 3 x 4 x 4
3 input channels
2 filters/kernels
Padding = 3
Stride =1



W[4, 4, 3]

[link]

Slide credit: Wang Wei

# Summary

- Convolution
  - Cross-correlation
  - Kernel, receptive field, padding, stride
  - VS MLP

- 2D convolution
  - Single channel, single kernel
  - Single channel, multiple kernels
  - Multiple channels, multiple kernels

- Pooling
  - Max and Average pooling

Slide credit: Wang Wei