# Fully Convolutional Networks

Machine Perception

Siyu Tang

02 April 2019
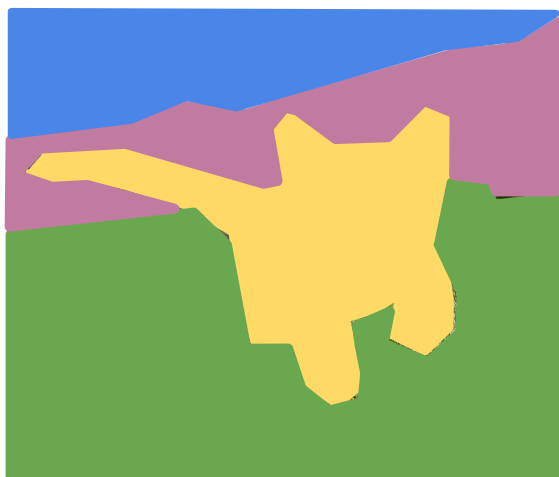
Advanced Interactive Technologies

ETH

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

# Pixel-wise Computer Vision Tasks

# Semantic Segmentation

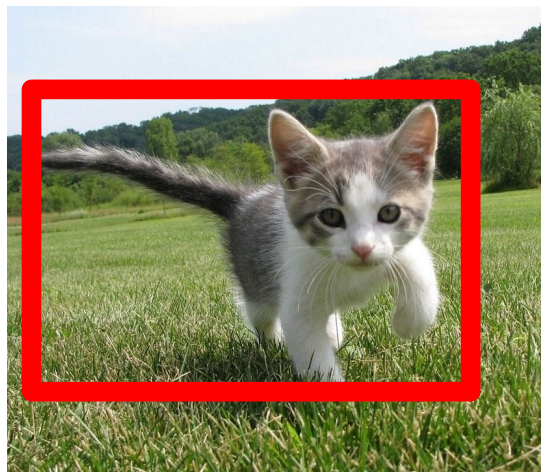# Semantic Segmentation



**Semantic Segmentation**
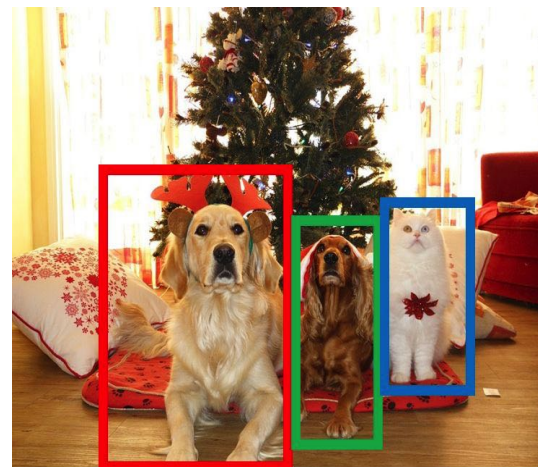
GRASS, CAT, TREE, SKY

No objects, just pixels

**Classification + Localization**

CAT

Single Object

**Object Detection**

DOG, DOG, CAT

Multiple Object

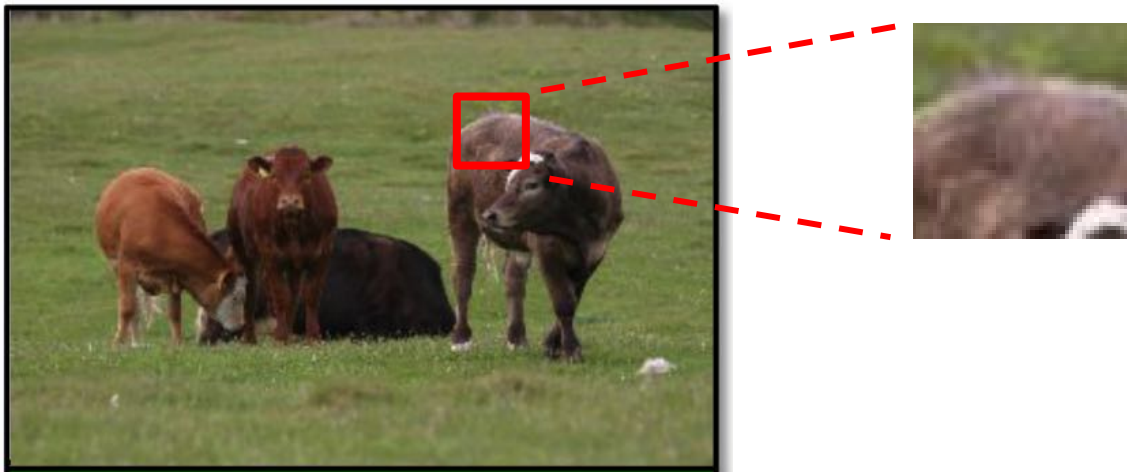**Instance Segmentation**

DOG, DOG, CAT
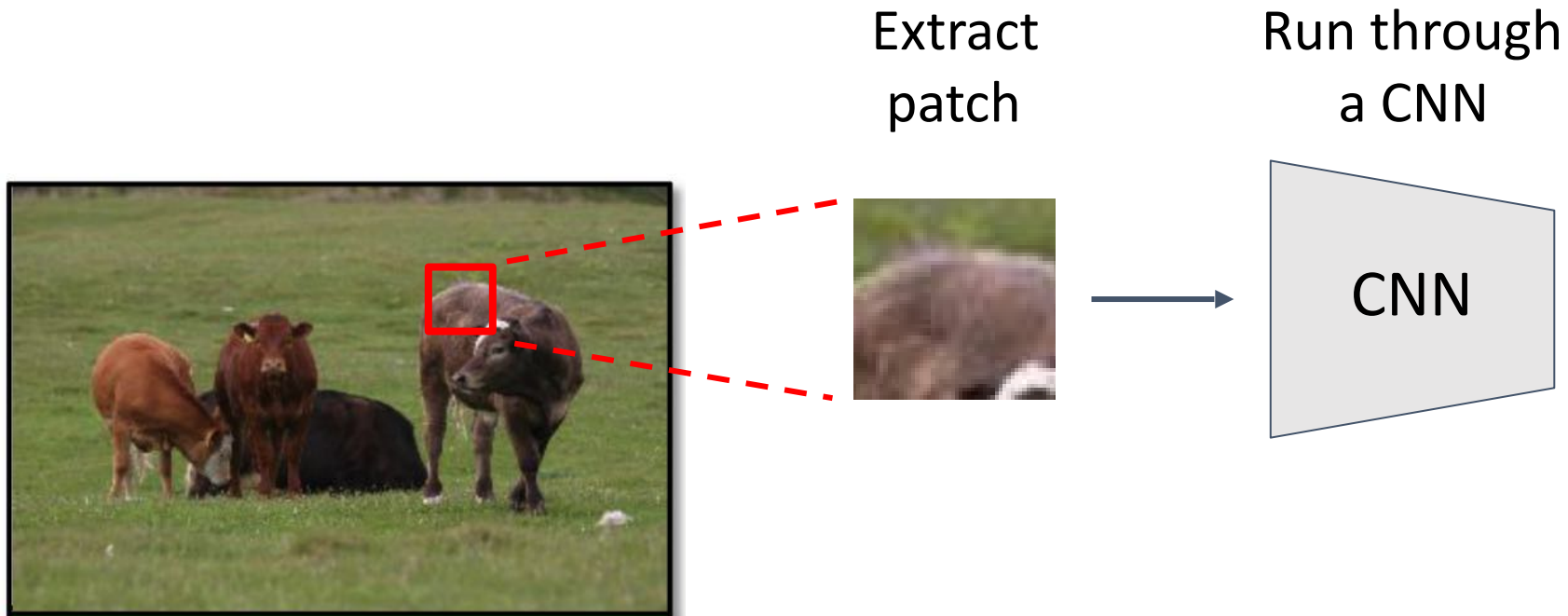
This image is CC0 public domain
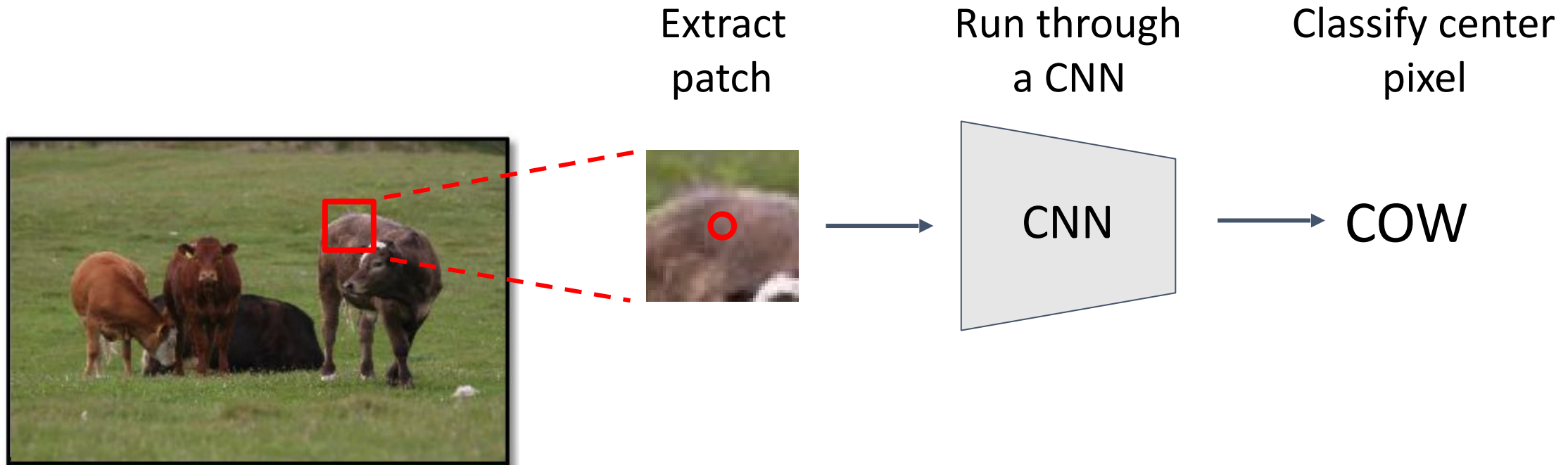
# Semantic Segmentation

# Semantic Segmentation

Extract
patch

# Semantic Segmentation

Extract patch

Run through a CNN



CNN

# Semantic Segmentation

**Extract patch**

**Run through a CNN**

**Classify center pixel**

CNN

→ COW

# Semantic Segmentation

Extract patch

Run through a CNN
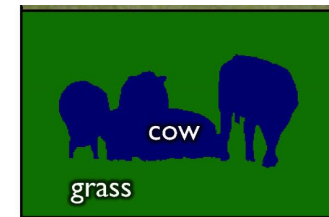
Classify center pixel



CNN

COW

Repeat for every pixel

cow

grass
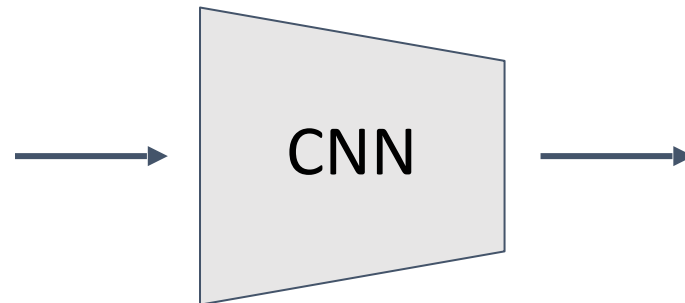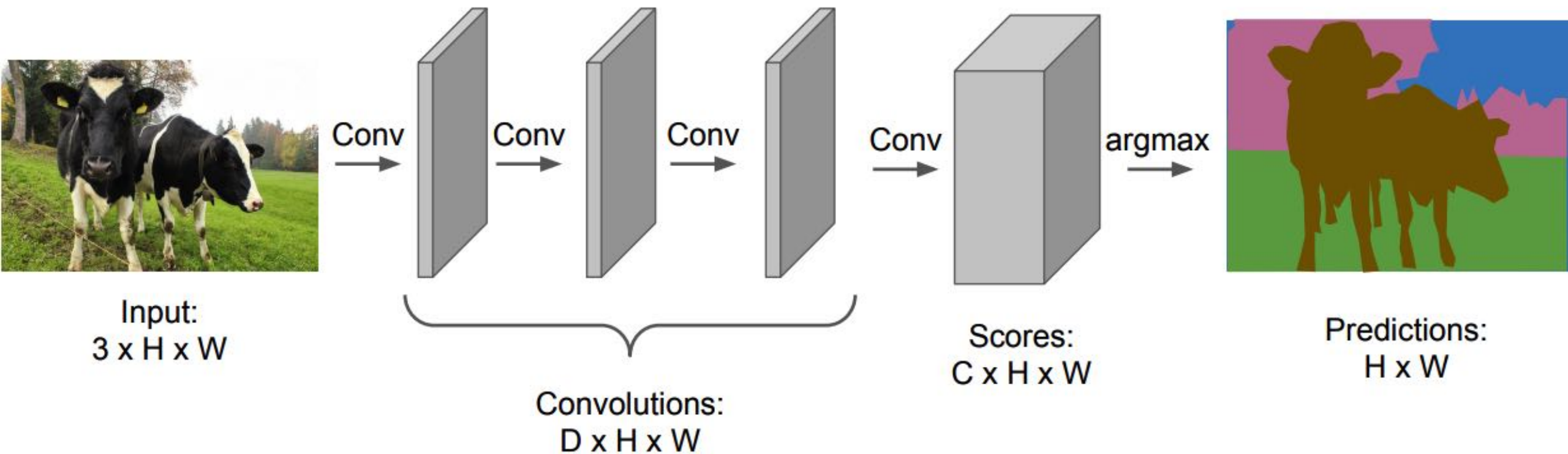
# Semantic Segmentation

Run "fully convolutional" network
to get all pixels at once



Smaller output
due to pooling

# Semantic Segmentation



Input:
3 x H x W

Conv → Conv → Conv → Conv → argmax

Convolutions:
D x H x W

Scores:
C x H x W

Predictions:
H x W

# Semantic Segmentation



Input: 3 x H x W

Conv → Conv → Conv → Conv

Convolutions: D x H x W

argmax

Scores: C x H x W
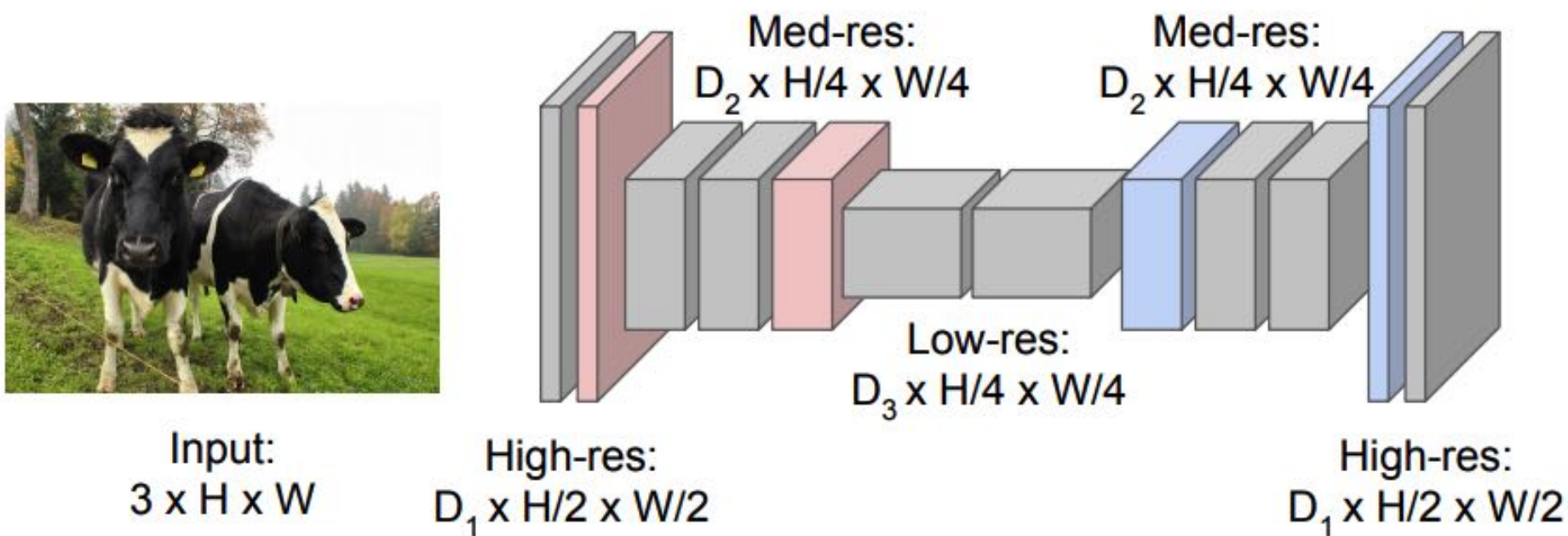
Predictions: H x W

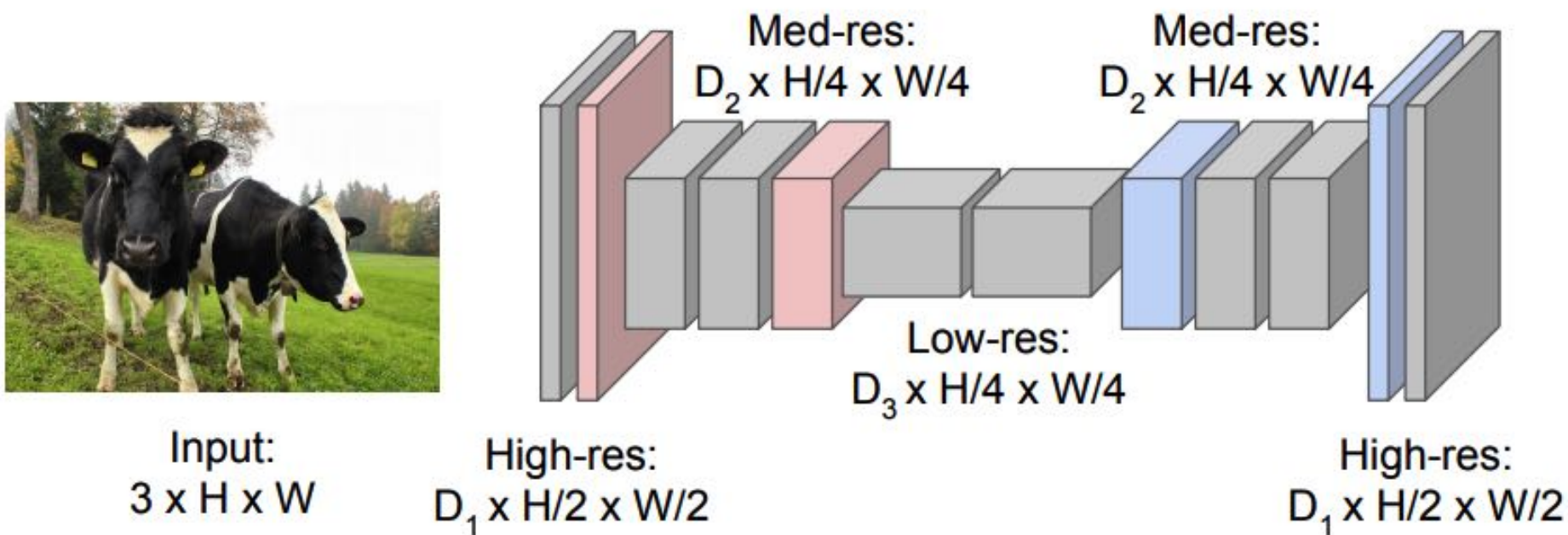Problem: convolutions at original image resolution is very expensive!

# Semantic Segmentation



Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!

Input: $3 \times H \times W$

High-res: $D_1 \times H/2 \times W/2$

Med-res: $D_2 \times H/4 \times W/4$

Low-res: $D_3 \times H/4 \times W/4$

Med-res: $D_2 \times H/4 \times W/4$

High-res: $D_1 \times H/2 \times W/2$

Predictions: $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

# Semantic Segmentation



Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!

Input:
$3 \times H \times W$

High-res:
$D_1 \times H/2 \times W/2$

Med-res:
$D_2 \times H/4 \times W/4$

Low-res:
$D_3 \times H/4 \times W/4$

Med-res:
$D_2 \times H/4 \times W/4$

High-res:
$D_1 \times H/2 \times W/2$

Predictions:
$H \times W$

**Downsampling**:
Pooling, strided convolution

**Upsampling**:
???

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

# In-Network upsampling: Unpooling

**Nearest Neighbor**

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

→

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

Input: 2 x 2          Output: 4 x 4

**"Bed of Nails"**

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

→

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

Input: 2 x 2          Output: 4 x 4

# In-Network upsampling: Max Unpooling

**Max Pooling**
Remember which element was max!

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

Input: 4 x 4

→

| 5 | 6 |
|---|---|
| 7 | 8 |

Output: 2 x 2

→ **...** →

Rest of the network

**Max Unpooling**
Use positions from pooling layer

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

Input: 2 x 2

→

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

Output: 4 x 4

Corresponding pairs of downsampling and upsampling layers

# Learnable Upsampling

Normal 3X3 convolution, stride 1 and pad 1



Input: 4 x 4

Output: 4 x 4

# Learnable Upsampling

Normal 3X3 convolution, stride 1 and pad 1



Dot product between filter and input

Input: 4 x 4

Output: 4 x 4

# Learnable Upsampling

Normal 3X3 convolution, stride 1 and pad 1



Dot product between filter and input

Input: 4 x 4

Output: 4 x 4

# Learnable Upsampling

Normal 3X3 convolution, stride 2 and pad 1

Input: 4 x 4

Output: 2 x 2

# Learnable Upsampling

Normal 3X3 convolution, stride 2 and pad 1



Dot product
between filter
and input

Input: 4 x 4

Output: 2 x 2

# Learnable Upsampling

Normal 3X3 convolution, stride 2 and pad 1



Dot product between filter and input

Input: 4 x 4

Output: 2 x 2

Filter moves 2 pixels in the input for every one pixel in the output

Stride gives ratio between movement in input and output

# Transposed Convolutions

Transposed Convolutions are used to upsample the input feature map to a desired output feature map using some learnable parameters.

The basic operation that goes in a transposed convolution is explained below:

1. Consider a 2x2 encoded feature map which needs to be upsampled to 3x3 feature map.



Figure 7. Input Feature Map

**Output**

Figure 8. Output Feature Map

3. Now we take the upper left element of the input feature map and multiply it with every element of the kernel as shown in figure 10.



Figure 10.

2. We take a kernel of size 2x2 with unit stride and zero padding.



Figure 9. Kernel of size 2x2

4. Similarly, we do it for all the remaining elements of the input feature map as depicted in figure 11.



5. As you can see, some of the elements of the resulting upsampled feature maps are over-lapping. To solve this issue, we simply add the elements of the over-lapping positions.

# Learnable Upsampling
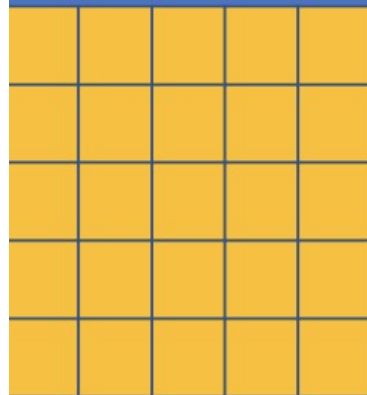
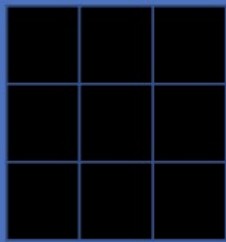3 x 3 **transpose** convolution, stride 2 pad 1
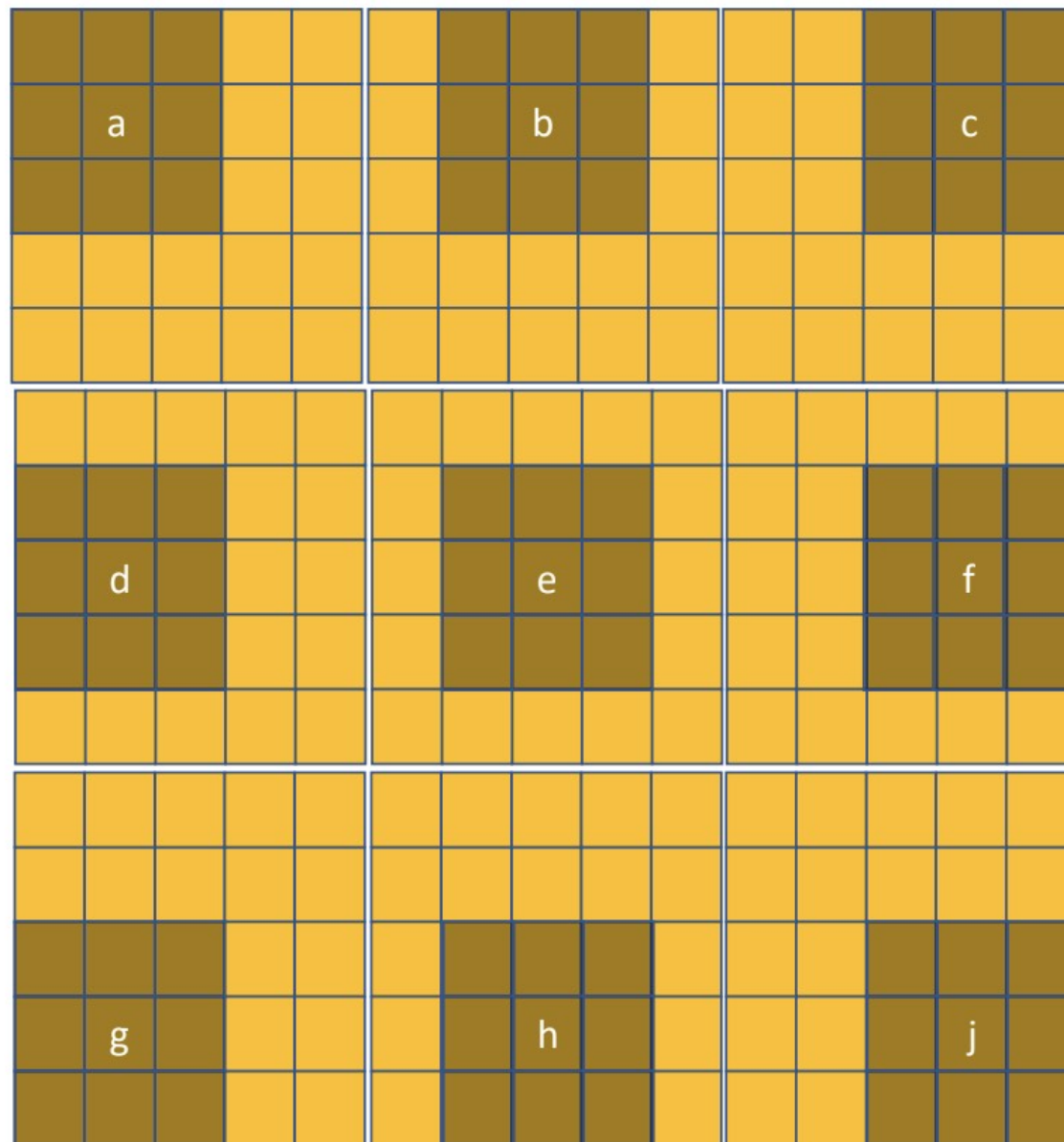
Input: 2 x 2

Output: 4 x 4

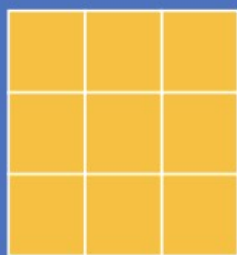Computer Vision

Detection

Image

kernel

Convolution

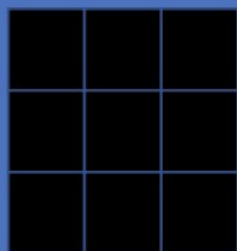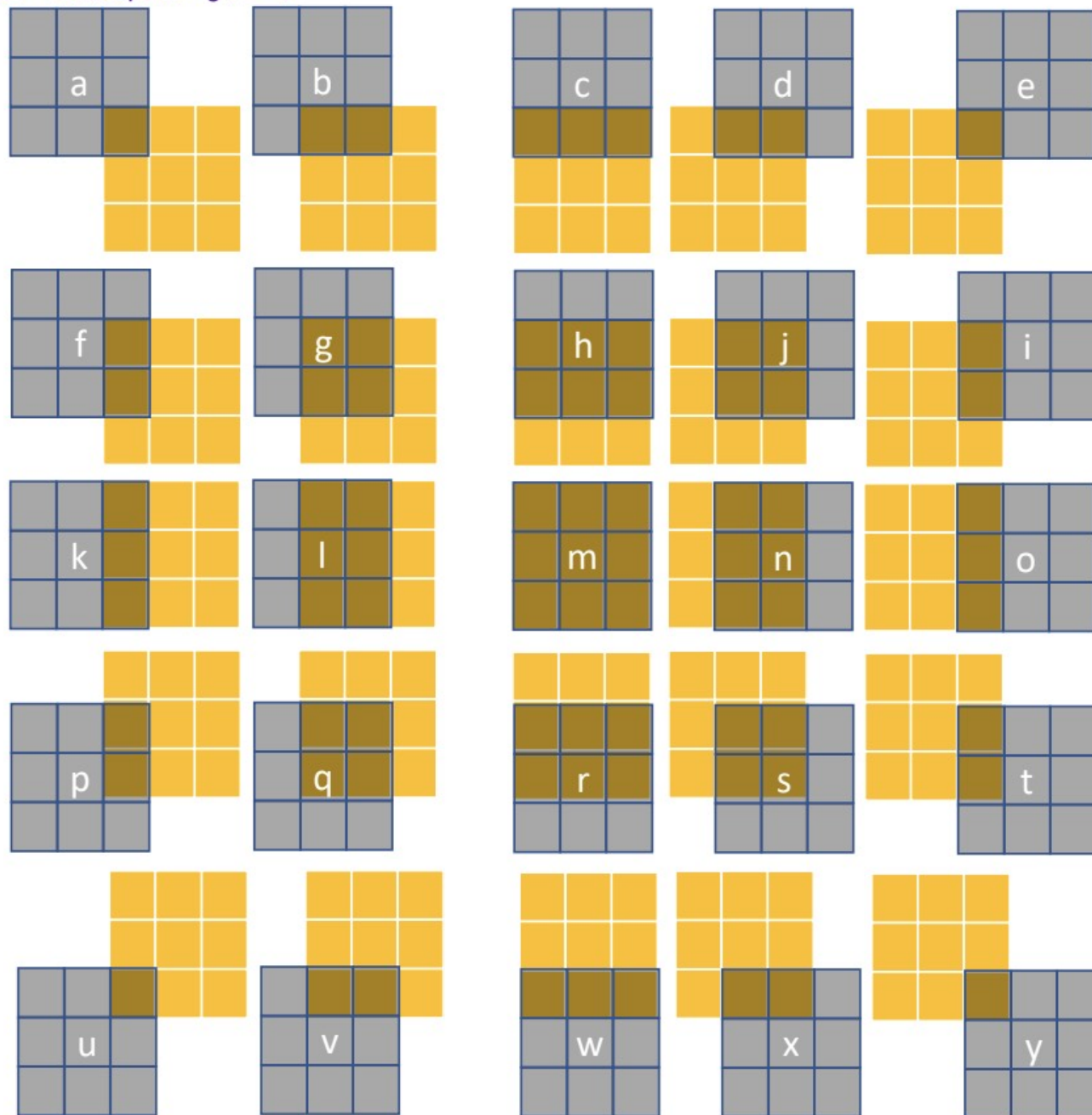| a | b | c |
|---|---|---|
| d | e | f |
| g | h | j |

# Computer Vision

Detection

Segmentation

Image

kernel

这里可能是padding了一下

Transposed Convolution

Sometimes referred to as deconvolution but that is not correct terminology.

有时称为反卷积，但这不是正确的术语

# Learnable Upsampling

3 x 3 **transpose** convolution, stride 2 pad 1

pad 2x2 to 4x4



Input gives
weight for
filter

Input: 2 x 2

Output: 4 x 4

# Learnable Upsampling



3 x 3 **transpose** convolution, stride 2 pad 1

Sum where output overlaps

Input gives weight for filter

Input: 2 x 2

Output: 4 x 4

Filter moves 2 pixels in the <u>output</u> for every one pixel in the <u>input</u>
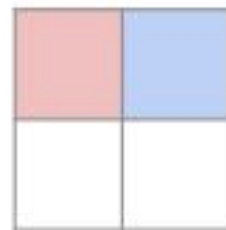
Stride gives ratio between movement in output and input
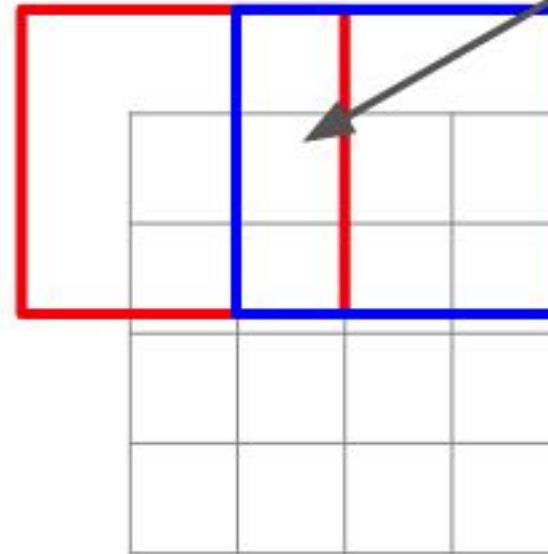
# Learnable Upsampling



3 x 3 **transpose** convolution, stride 2 pad 1

Input gives weight for filter

Input: 2 x 2

Output: 4 x 4

Sum where output overlaps

Filter moves 2 pixels in the <u>output</u> for every one pixel in the <u>input</u>
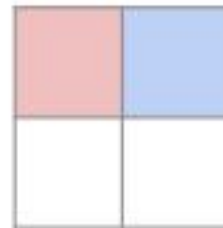
Stride gives ratio between movement in output and input
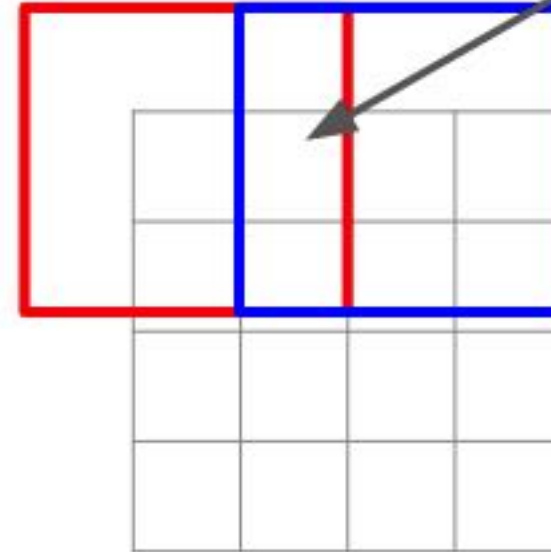
# Learnable Upsampling

**Other names:**
- Deconvolution (bad)
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution

3 x 3 **transpose** convolution, stride 2 pad 1

Input gives weight for filter

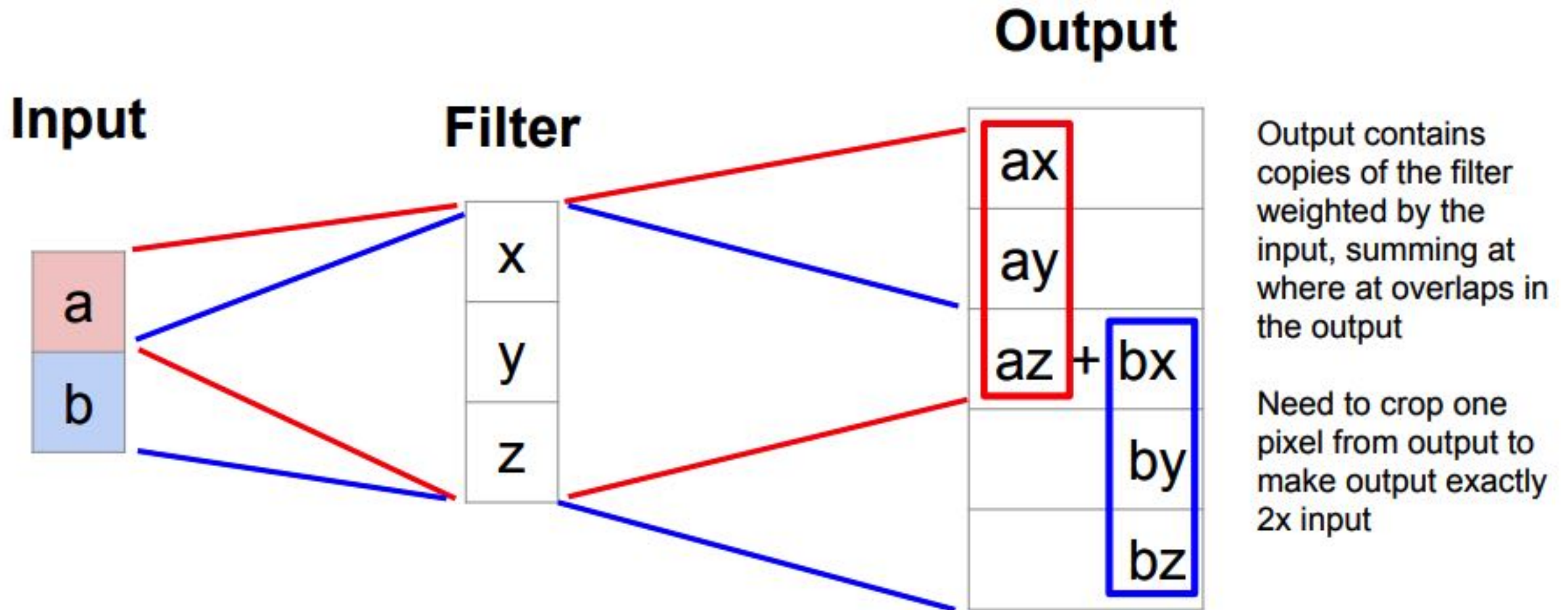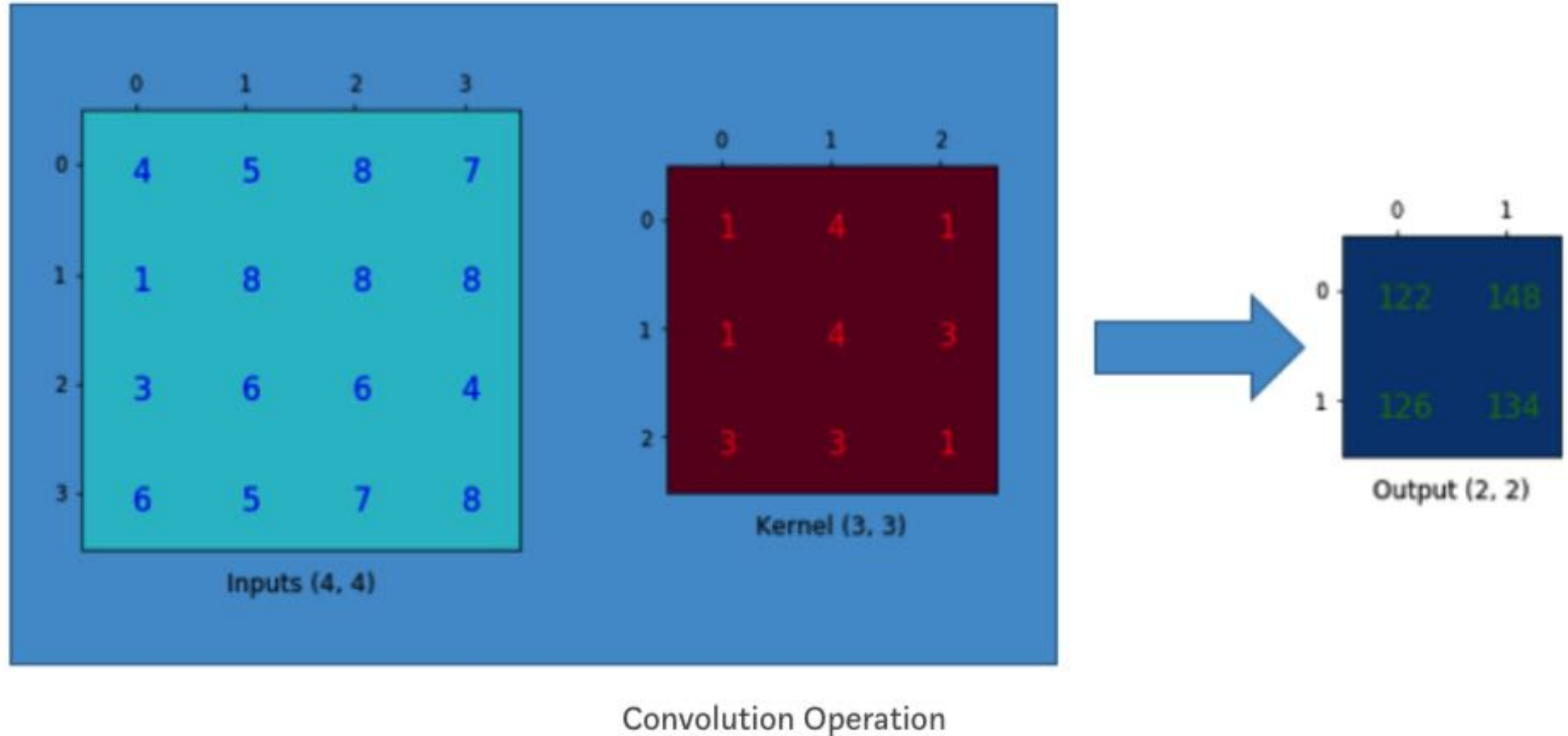Input: 2 x 2

Sum where output overlaps

Filter moves 2 pixels in the <u>output</u> for every one pixel in the <u>input</u>

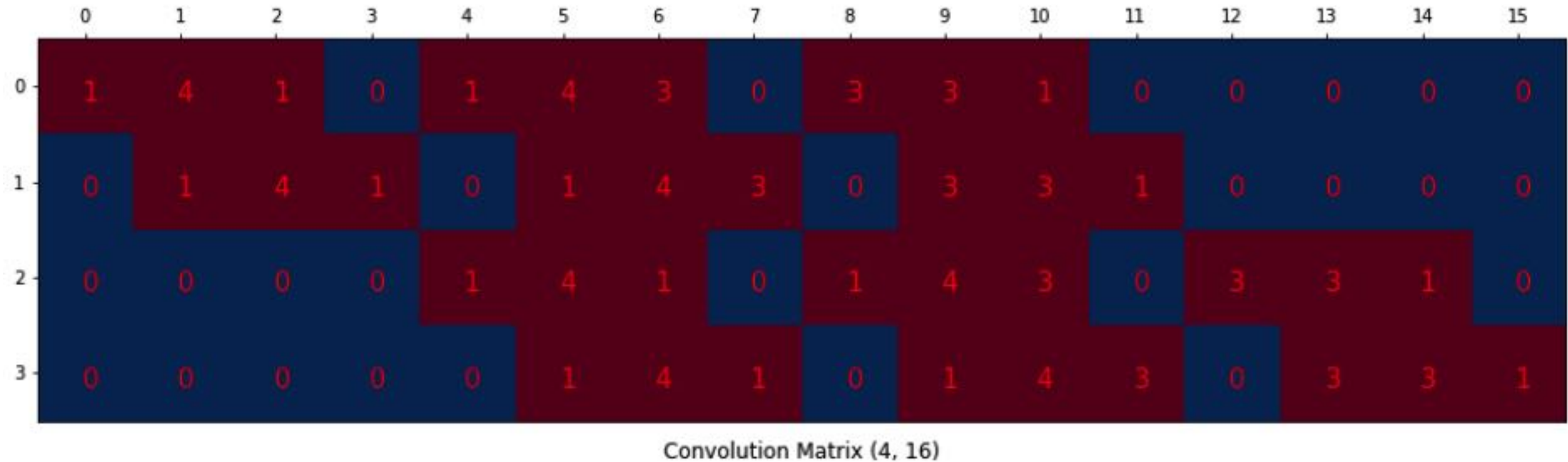Stride gives ratio between movement in output and input

Output: 4 x 4

# Learnable Upsampling: 1D Example



**Input**

a
b

**Filter**

x
y
z

**Output**

ax
ay
az + bx
by
bz

Output contains copies of the filter weighted by the input, summing at where at overlaps in the output

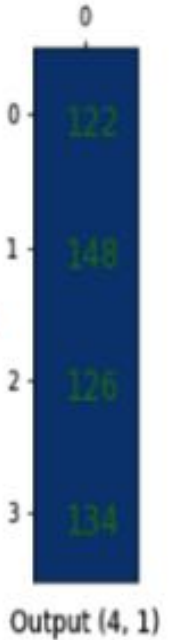Need to crop one pixel from output to make output exactly 2x input

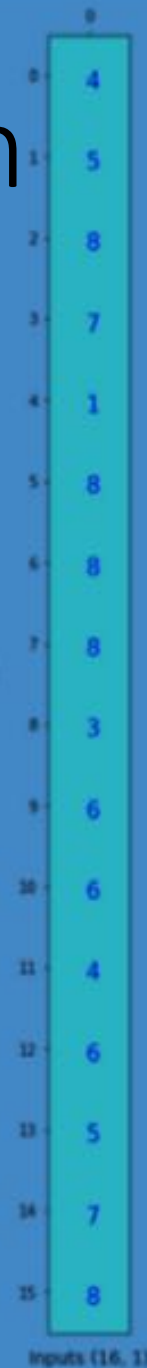# Convolution as Matrix Operation
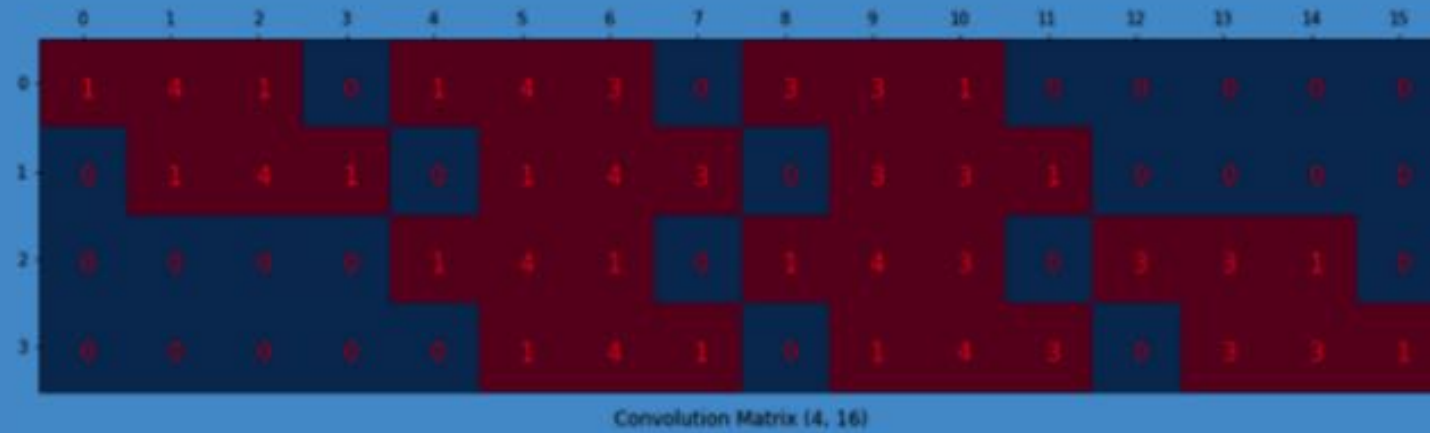


Convolution Operation
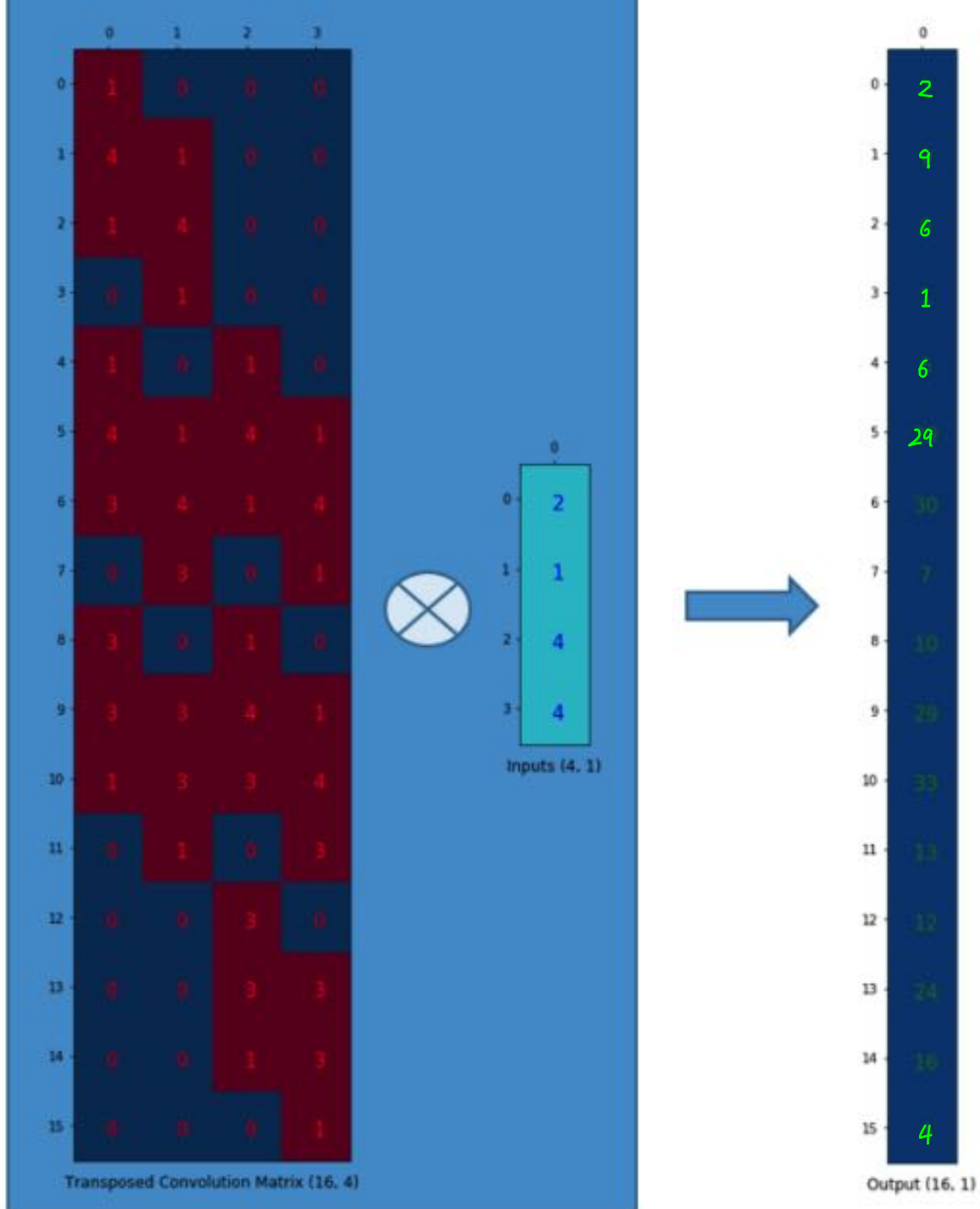
# Convolution as Matrix Operation



Convolution Matrix (4, 16)
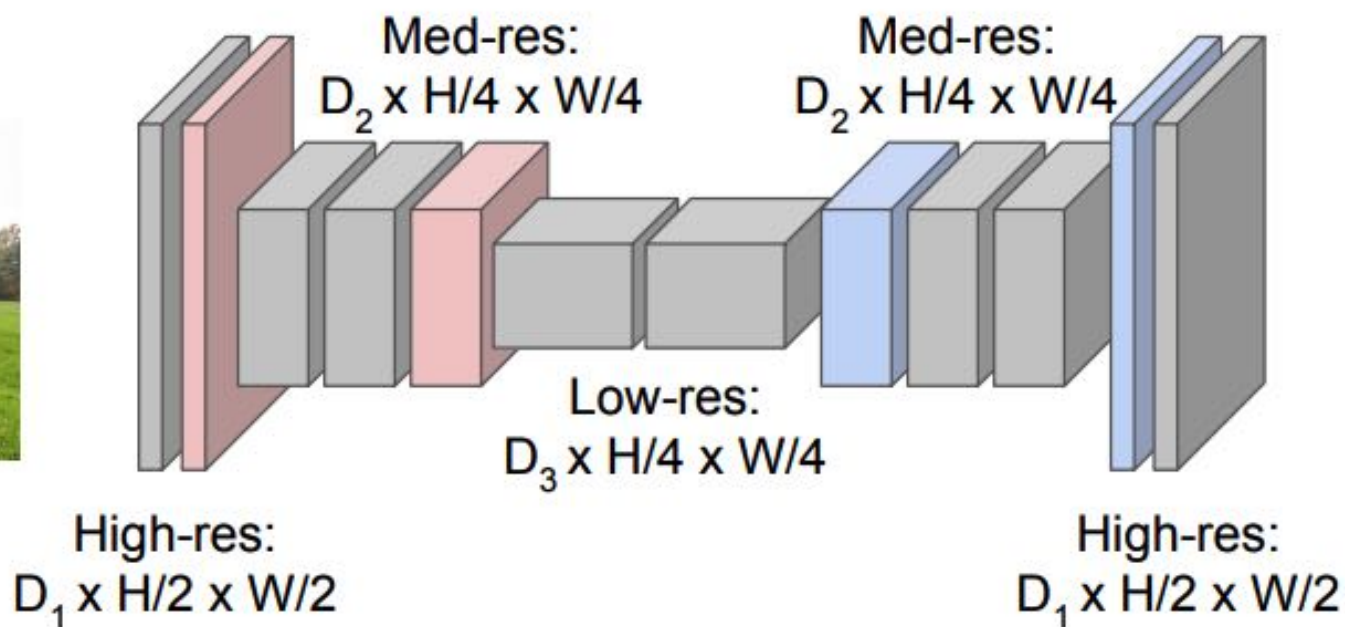
# Convolution as Matrix Operation

# Transpose Convolution

# Semantic Segmentation



Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!

Med-res:
$D_2 \times H/4 \times W/4$

Med-res:
$D_2 \times H/4 \times W/4$

Low-res:
$D_3 \times H/4 \times W/4$

Input:
$3 \times H \times W$

High-res:
$D_1 \times H/2 \times W/2$

High-res:
$D_1 \times H/2 \times W/2$

Predictions:
$H \times W$

**Downsampling**:
Pooling, strided convolution
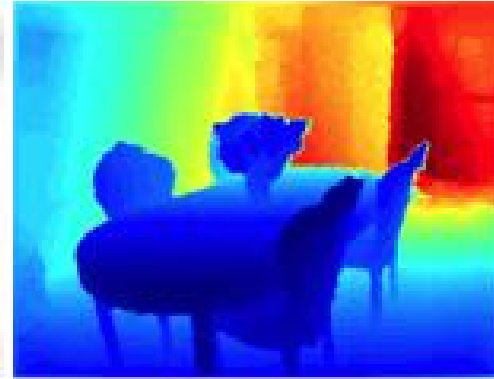
**Upsampling**:
Unpooling or strided transpose convolution

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

# Beyond Semantic Segmentation

semantic segmentation

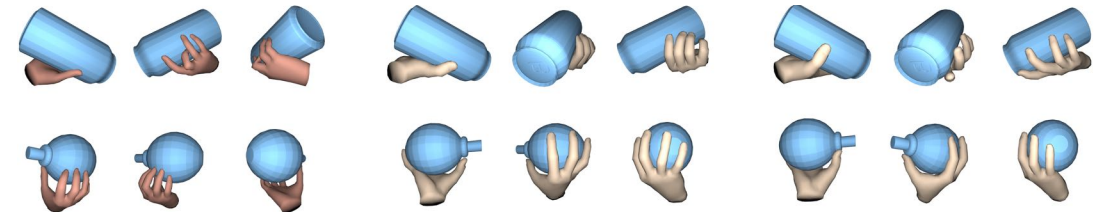monocular depth estimation (Liu et al. 2015)

boundary prediction (Xie & Tu 2015)

Slide credit: Jonathan Long

# Computer Vision and Learning Group (VLG)



Grasping generation

Human Scene Interaction

Generative human modelling

Human pose estimation