

Deep Learning Seminar

7. Segmentation

HA SEUNG HYUN

Contents

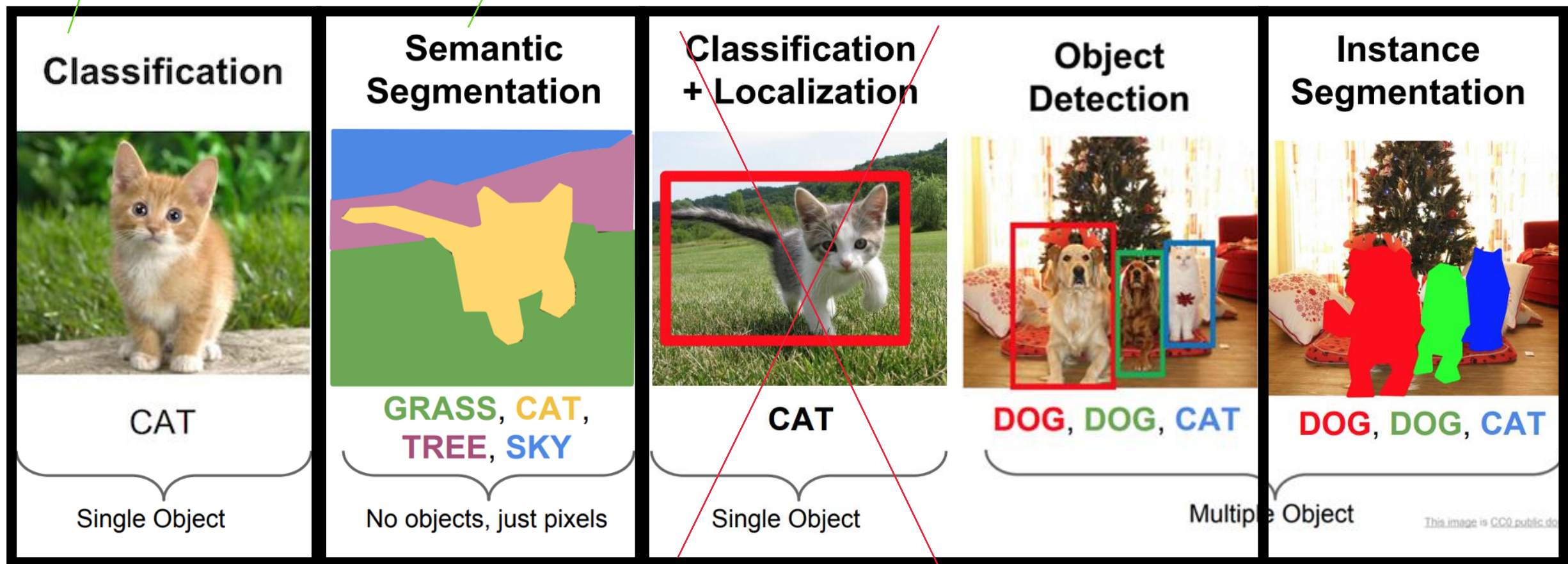
1. Overview

2. Classification

3. Semantic Segmentation

1. Overview

Overview



object detection single

multi

1. Computation Cost 3, 2, 1

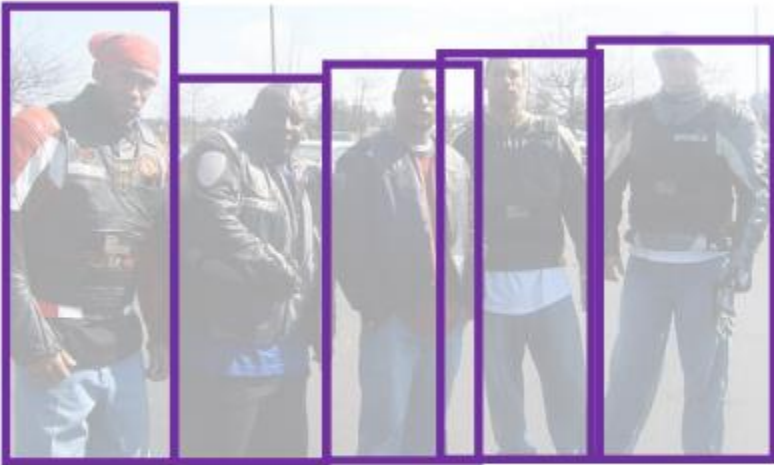
Object Detection 1

Semantic Segmentation 2 - 3

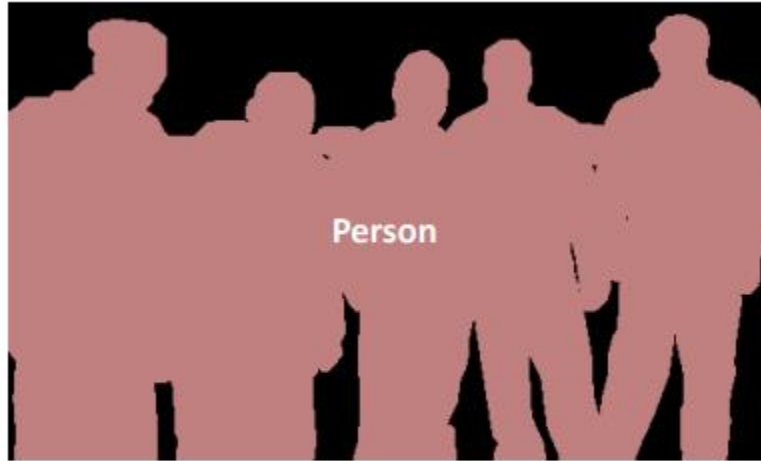
Instance Segmentation 3 - 4

가 . 가 Object Detection

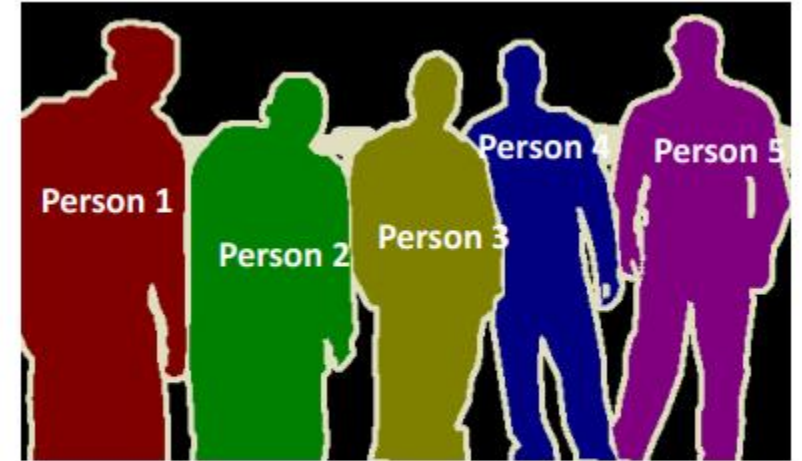
가 . Object Detection



Object Detection



Semantic Segmentation



Instance Segmentation

Overview

FPS (Frame Per Second)

, 1

Frame

가

. 1 8

FPS가

가

.

가

가

Computation Cost가

FPS가

1

가

Object Detection Example

가

. layer가

...

,

가

Forwarding

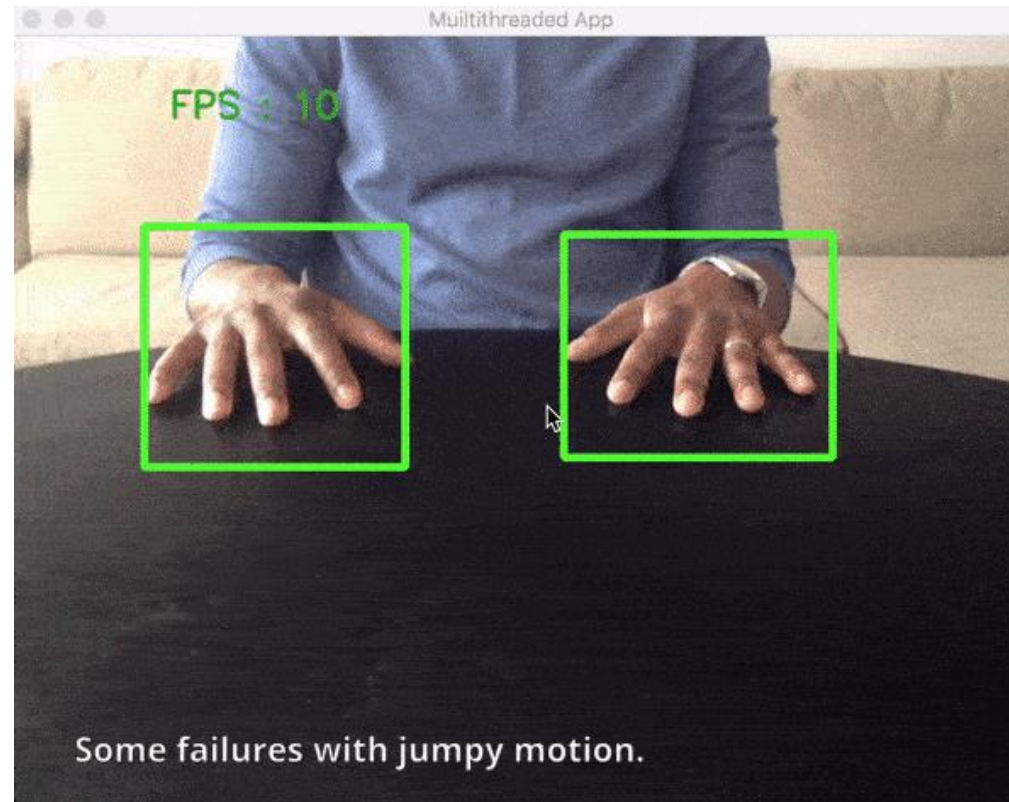
~

가

- Low And Model

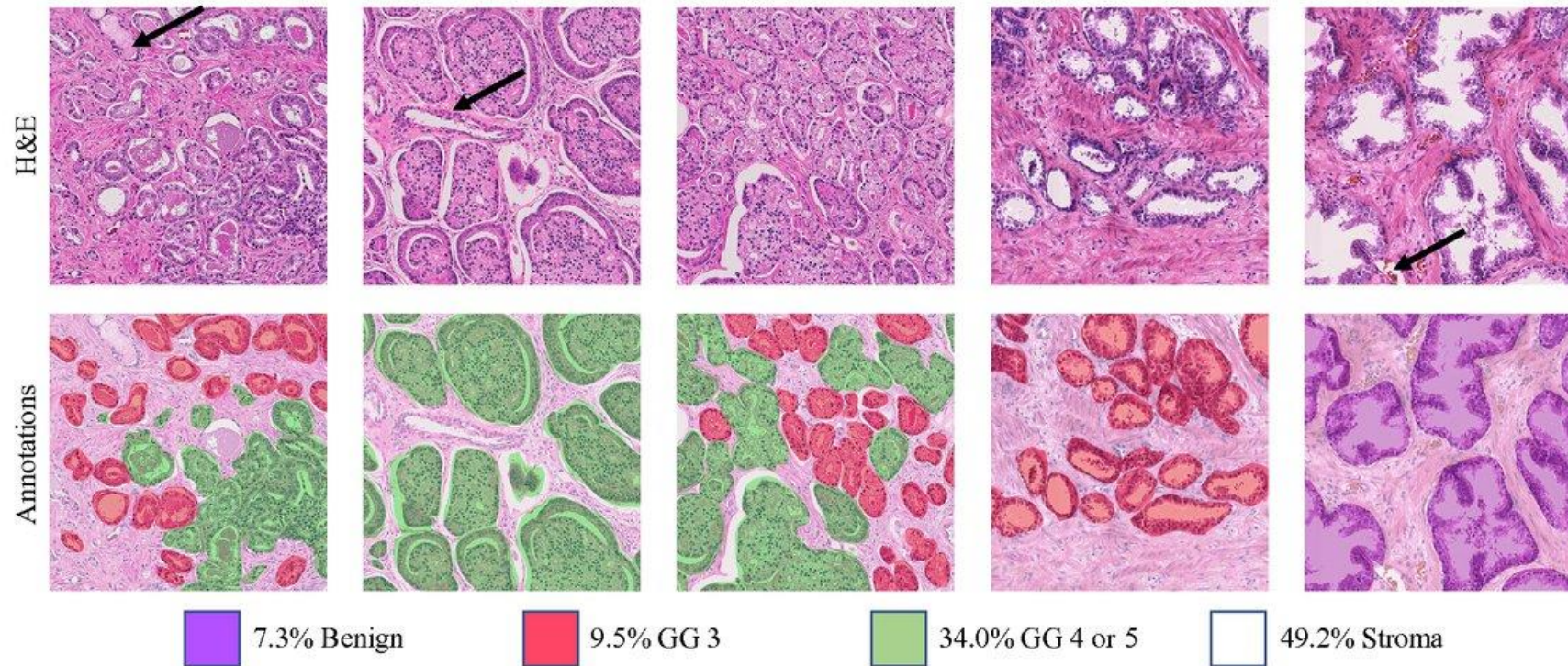
High And Model

Trade - Off



Overview

Semantic Segmentation Example



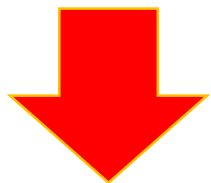
Overview

Instance Segmentation Example



2. Classification

Classification



Classification



CAT

Single Object

Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Classification + Localization



CAT

Single Object

Object Detection



DOG, DOG, CAT

Instance Segmentation



DOG, DOG, CAT

Multiple Object

This image is CC0 public domain

Classification

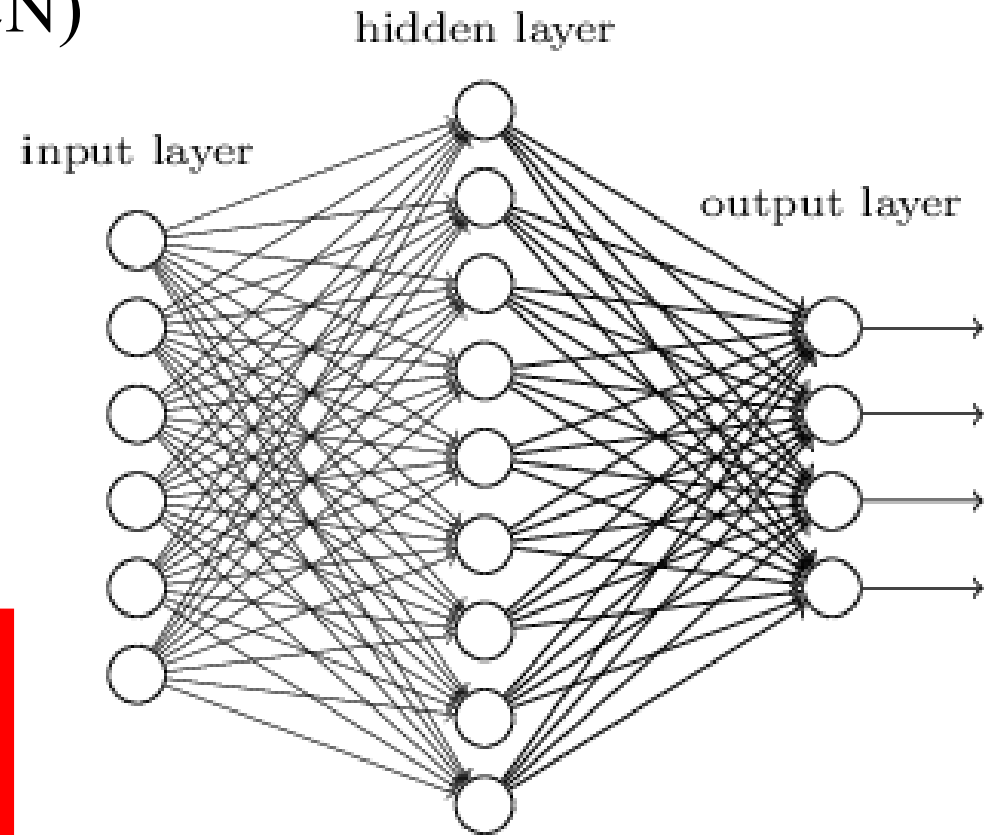
- Fully Connected Network (FCN)

1) 1



- 1) Disappear spatial information
- 2) Computationally Expensive

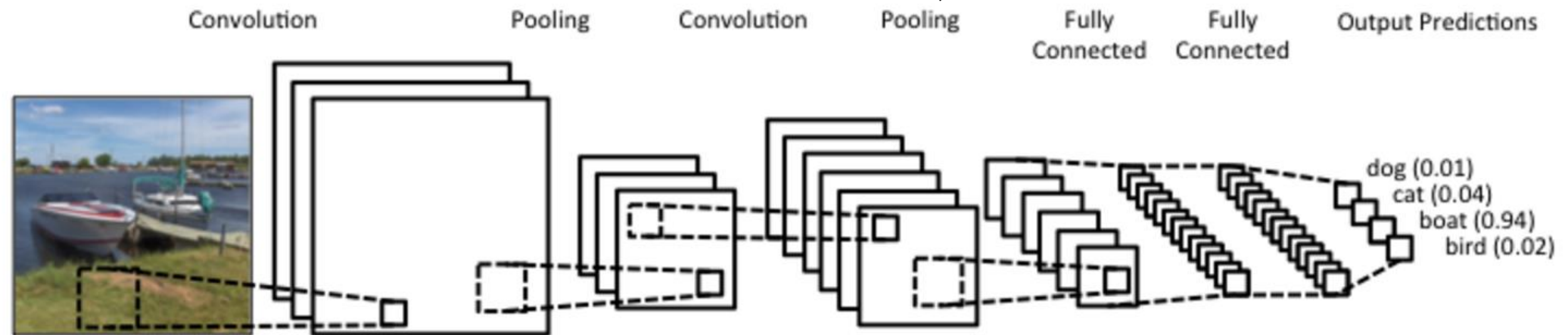
Input



Fully Connected Network

Classification

- Convolutional Neural Network (CNN)



Convolutional Neural Network

Classification

- Popular Model

- 1) VGG
- 2) GoogLeNet (Inception)
- 3) ResNet
- 4) DenseNet

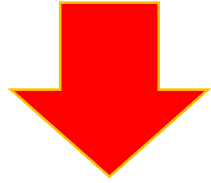
3. Semantic Segmentation


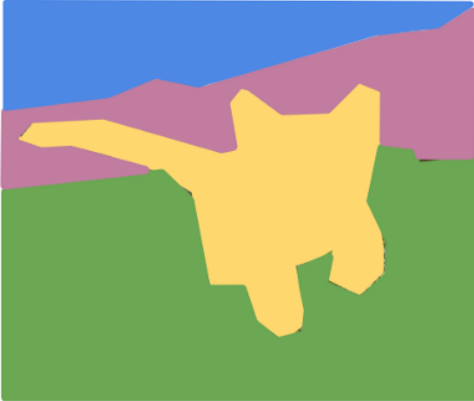
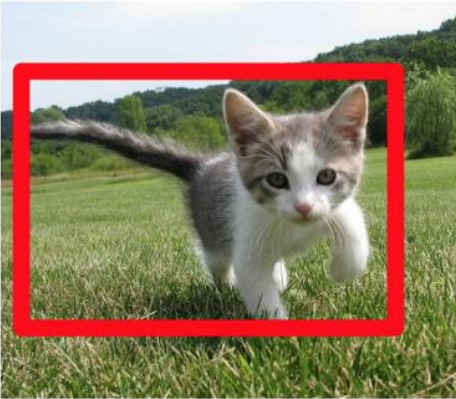


- 3-1. Overview of semantic segmentation
- 3-2. Upsampling & Convolutions
- 3-3. U-Net
- 3-4. Evaluation Matrix

FPS

Object Detection
Low And Model

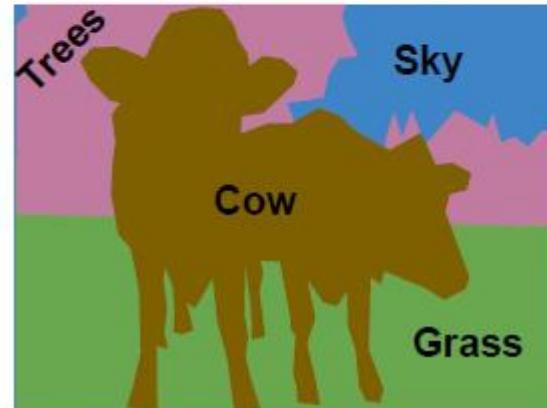
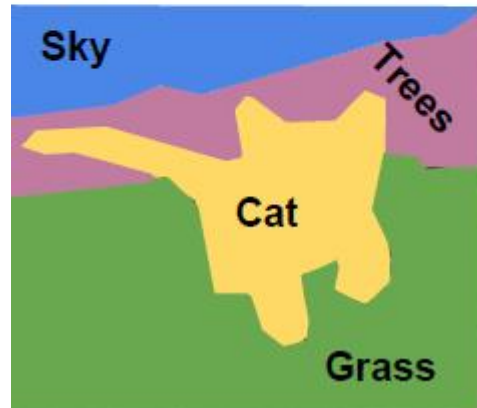
Semantic Segmentation



Classification	Semantic Segmentation	Classification + Localization	Object Detection	Instance Segmentation
				
CAT	GRASS, CAT, TREE, SKY	CAT	DOG, DOG, CAT	DOG, DOG, CAT
Single Object	No objects, just pixels	Single Object	Multiple Object	

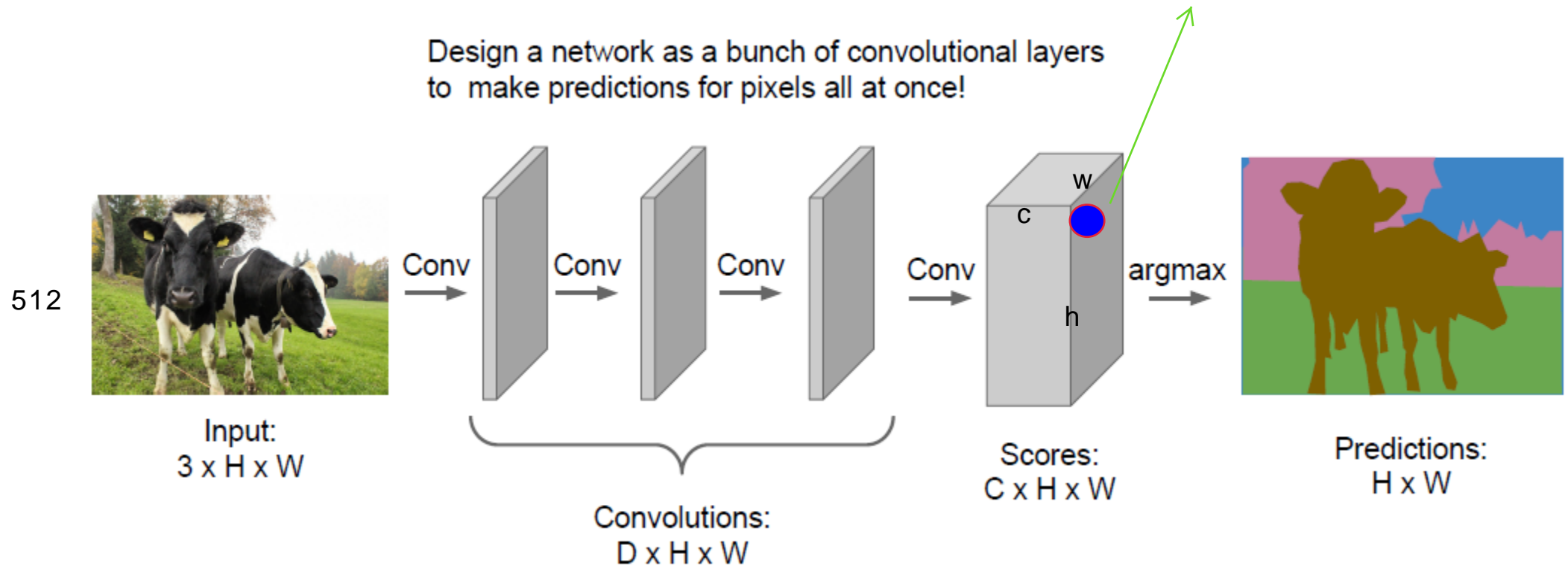
Semantic Segmentation

- Label each pixel in the image with category label



Semantic Segmentation

CNN 가 Semantic Segmentation 가



ConV

. filter 3x3

stride=1, padding=1

(CNN)

Semantic Segmentation



Input:
 $3 \times H \times W$

Conv

Conv

Conv

Conv

argmax



Predictions:
 $H \times W$

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!

Convolutions:
 $D \times H \times W$

Scores:
 $C \times H \times W$

	C	H	W
input	3	512	512
Filter	3x3x3	padding 1	stride 1 16
layer1	16	512	512
Filter	3x3x3	padding 1	stride 1 32
layer2	32	512	512
Filter	3x3x3	padding 1	stride 1 64
layer3	64	512	512
Filter	3x3x3	padding 1	stride 1 128
layer4	128	512	512
layer5	64	512	512
layer6	32	512	512
output	10	512	512

Problem: convolutions at original image resolution will be very expensive ...

vgg size 2 channels 2

1/2

CIFAR10

32x32

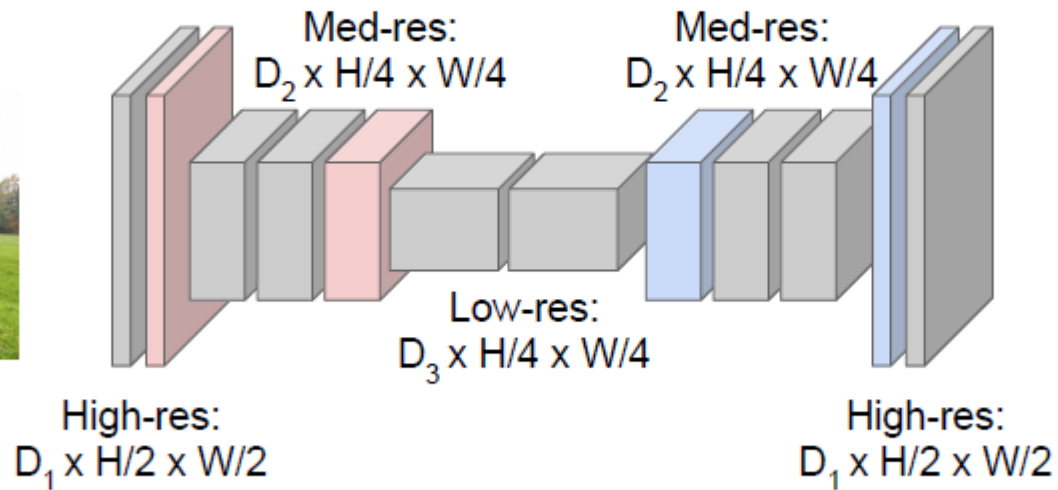
size

Semantic Segmentation

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



Input:
 $3 \times H \times W$



Predictions:
 $H \times W$

1. , output

2. 가(downsampling) (upsampling)

Semantic Segmentation

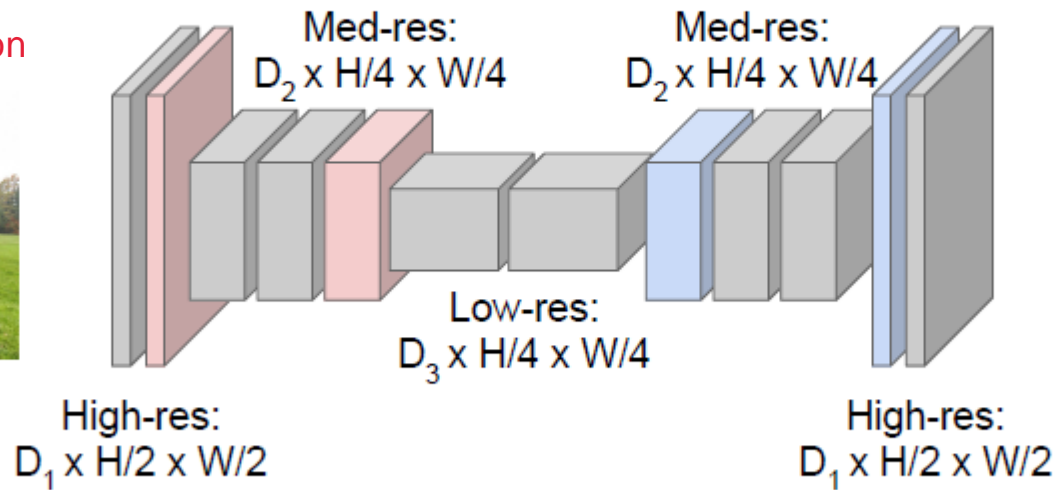
Downsampling:
Pooling, strided
convolution



Input:
 $3 \times H \times W$

가

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



. (down sampling) 1. stride = 2
2. max pooling

x.

Upsampling:
???

!!

: DeConvolution



Predictions:
 $H \times W$

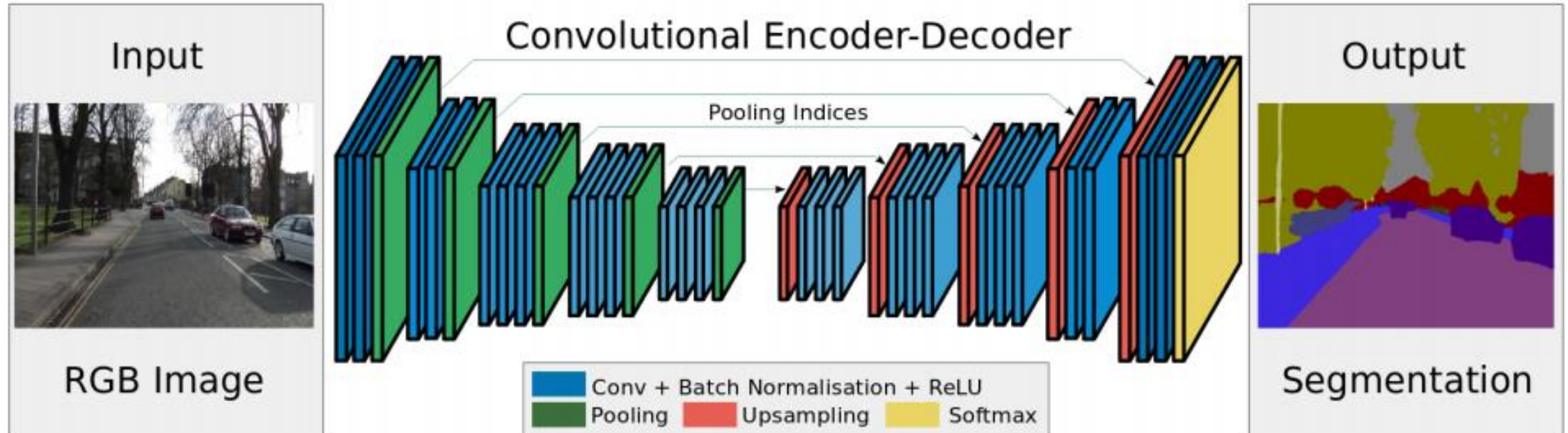
.

.

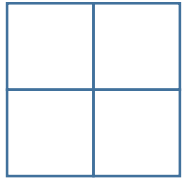
Semantic Segmentation

Batch Normalization
. batch size

w, b



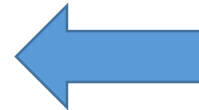
Upsampling



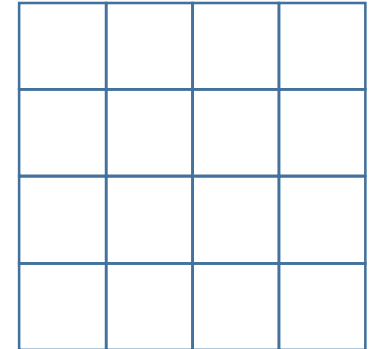
Input: 2 x 2

Transposed Convolution

Up-sampling
(ex. Deconv)



Down-sampling
(ex. Conv)



Output: 4 x 4

Unpooling (Upsampling)

- 1. 0 .
- Train .
- 0

가

“Bed of Nails”

1	2
3	4



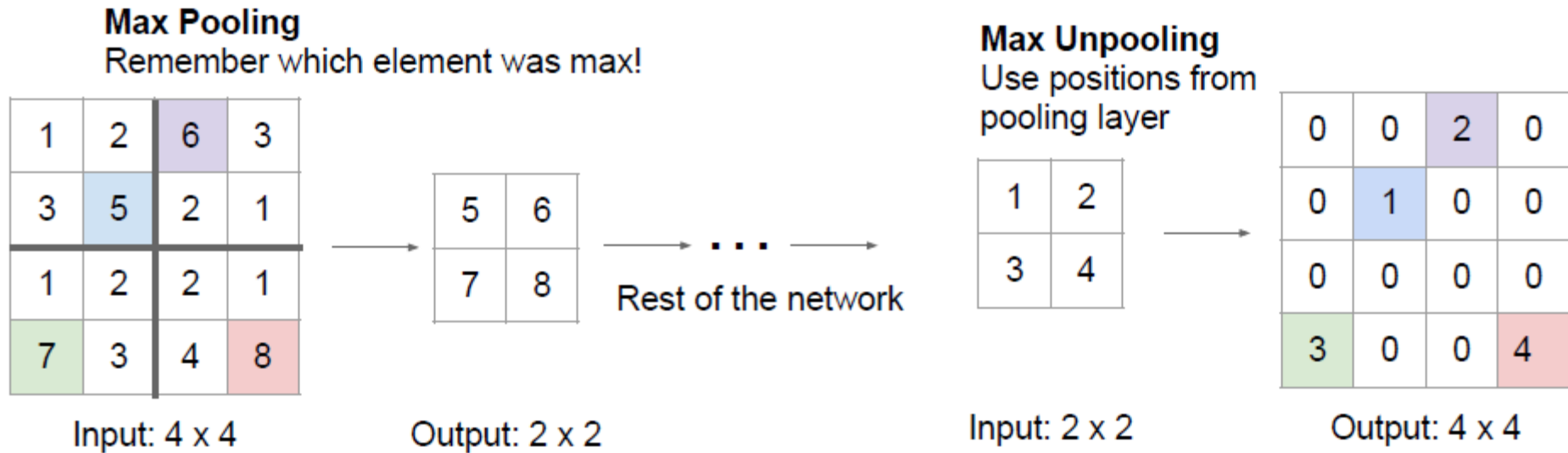
1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

Input: 2 x 2

Output: 4 x 4

Not Trainable

Unpooling (Upsampling)



Corresponding
downsampling
upsampling la

Not Trainable

Interpolation

(Upsampling)

- 가

ex) 1, 2 1 1 2 2
ex) 1, 2 1 1.5 2 2.5

Nearest Neighbor

1	2
3	4

Input: 2 x 2



1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

Output: 4 x 4

1.

2.

3.

-

Garbage

Not Trainable

filter

backward가

W

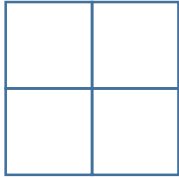
Filter

Not Trainable

Transpose Convolution

(Upsampling)

= Deconvolution



Input: 2 x 2

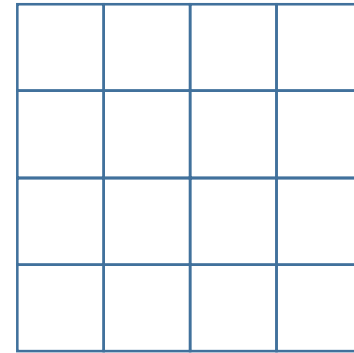
Upsampling



Weight

Transpose Filter

.



Output: 4 x 4

Trainable

Transpose Convolution

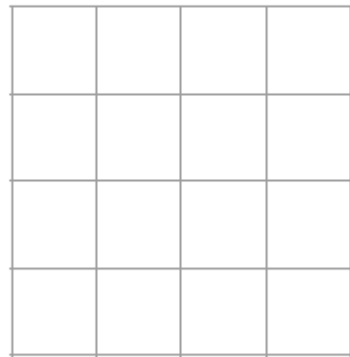
(Upsampling)

= Deconvolution

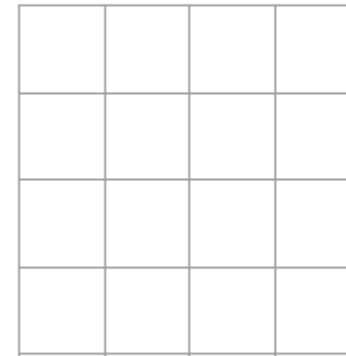
Hot

Interpolation VS Deconvolution

Recall: Typical 3 x 3 convolution, stride 1 pad 1



Input: 4 x 4



Output: 4 x 4

↓
"Trainable"

↓
"deconvolution checker border"

"

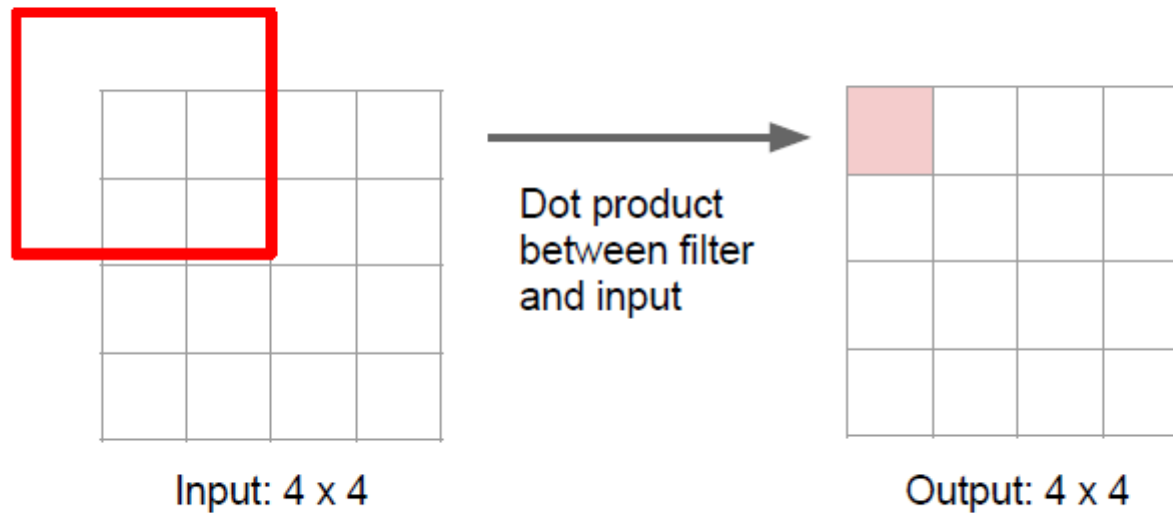
."

Interpolation

Transpose Convolution

(Upsampling)

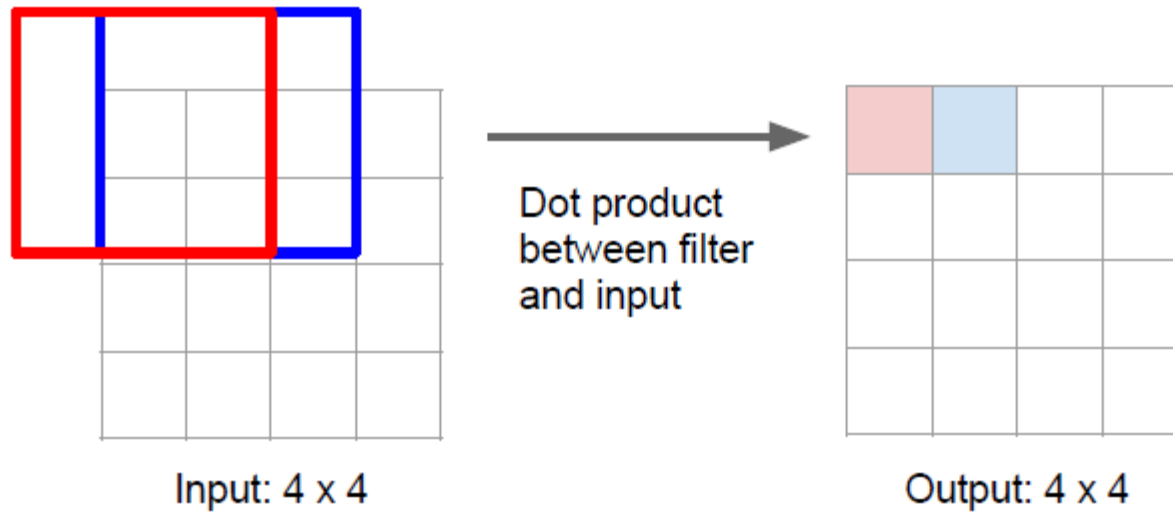
Recall: Normal 3 x 3 convolution, stride 1 pad 1



Transpose Convolution

(Upsampling)

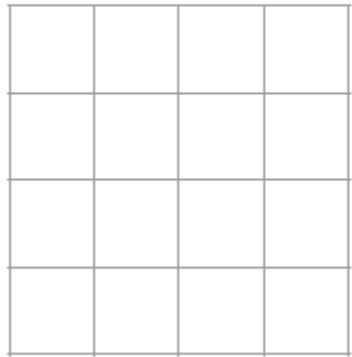
Recall: Normal 3 x 3 convolution, stride 1 pad 1



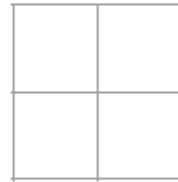
Transpose Convolution

(Upsampling)

Recall: Normal 3 x 3 convolution, stride 2 pad 1



Input: 4 x 4

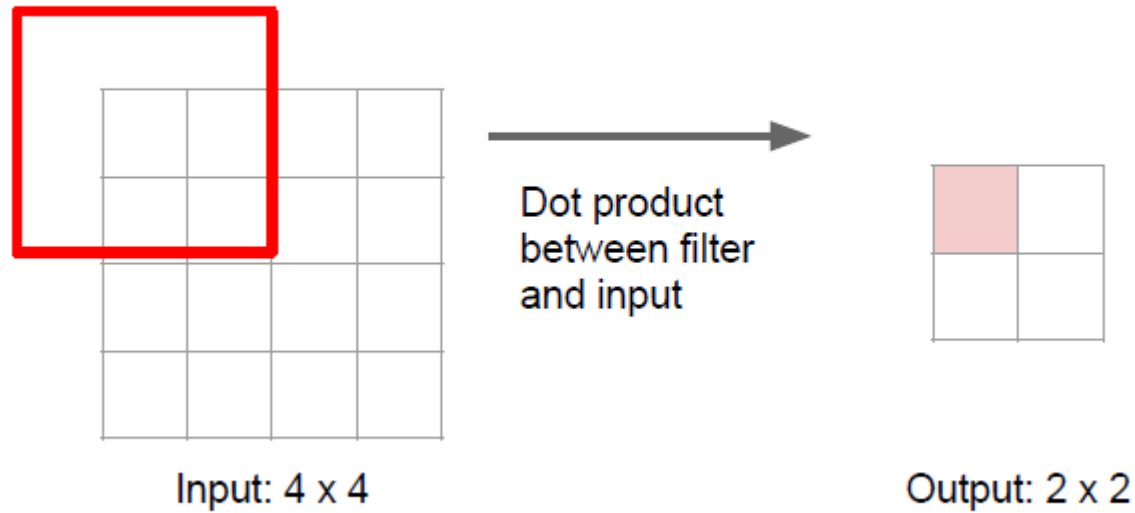


Output: 2 x 2

Transpose Convolution

(Upsampling)

Recall: Normal 3 x 3 convolution, stride 2 pad 1

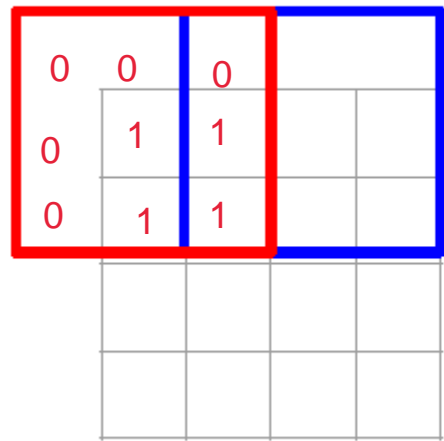


Transpose Convolution

(Upsampling)

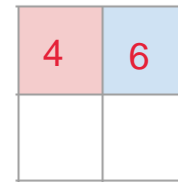
filter
1 1 1
1 1 1
1 1 1

Recall: Normal 3 x 3 convolution, stride 2 pad 1



Input: 4 x 4

Dot product
between filter
and input



Output: 2 x 2

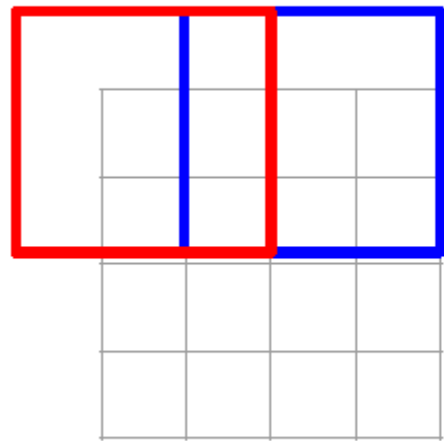
Filter moves 2 pixels in
the input for every one
pixel in the output

Stride gives ratio between
movement in input and
output

Transpose Convolution

(Upsampling)

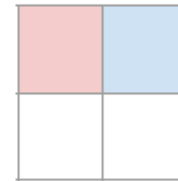
Recall: Normal 3 x 3 convolution, stride 2 pad 1



Input: 4 x 4



Dot product
between filter
and input



Output: 2 x 2

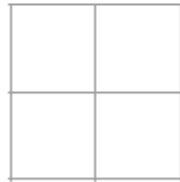
Filter moves 2 pixels in
the input for every one
pixel in the output

Stride gives ratio between
movement in input and
output

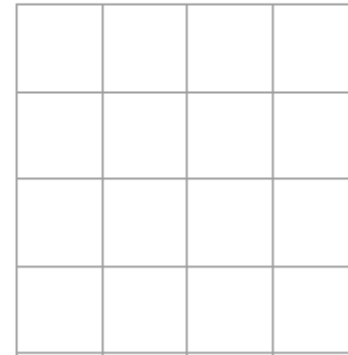
Transpose Convolution

(Upsampling)

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



Input: 2 x 2

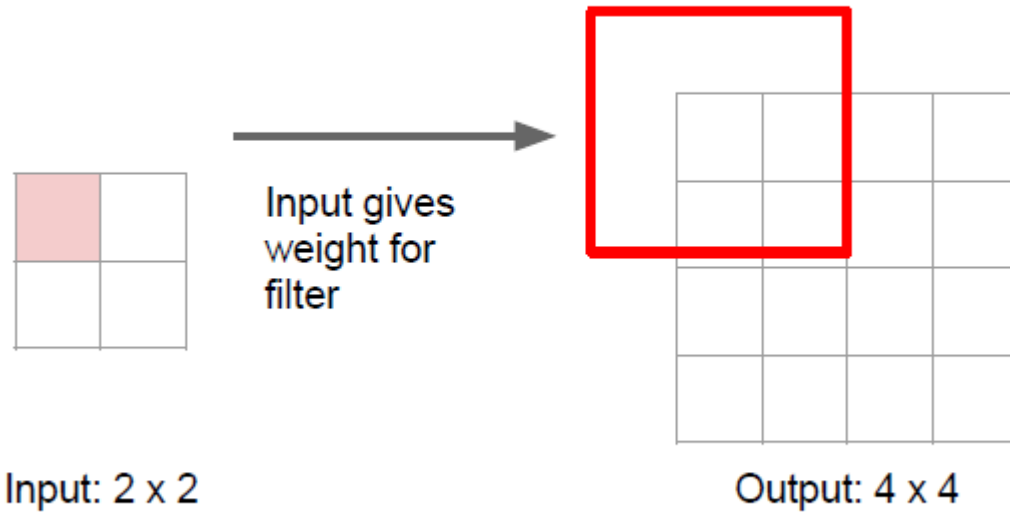


Output: 4 x 4

Transpose Convolution

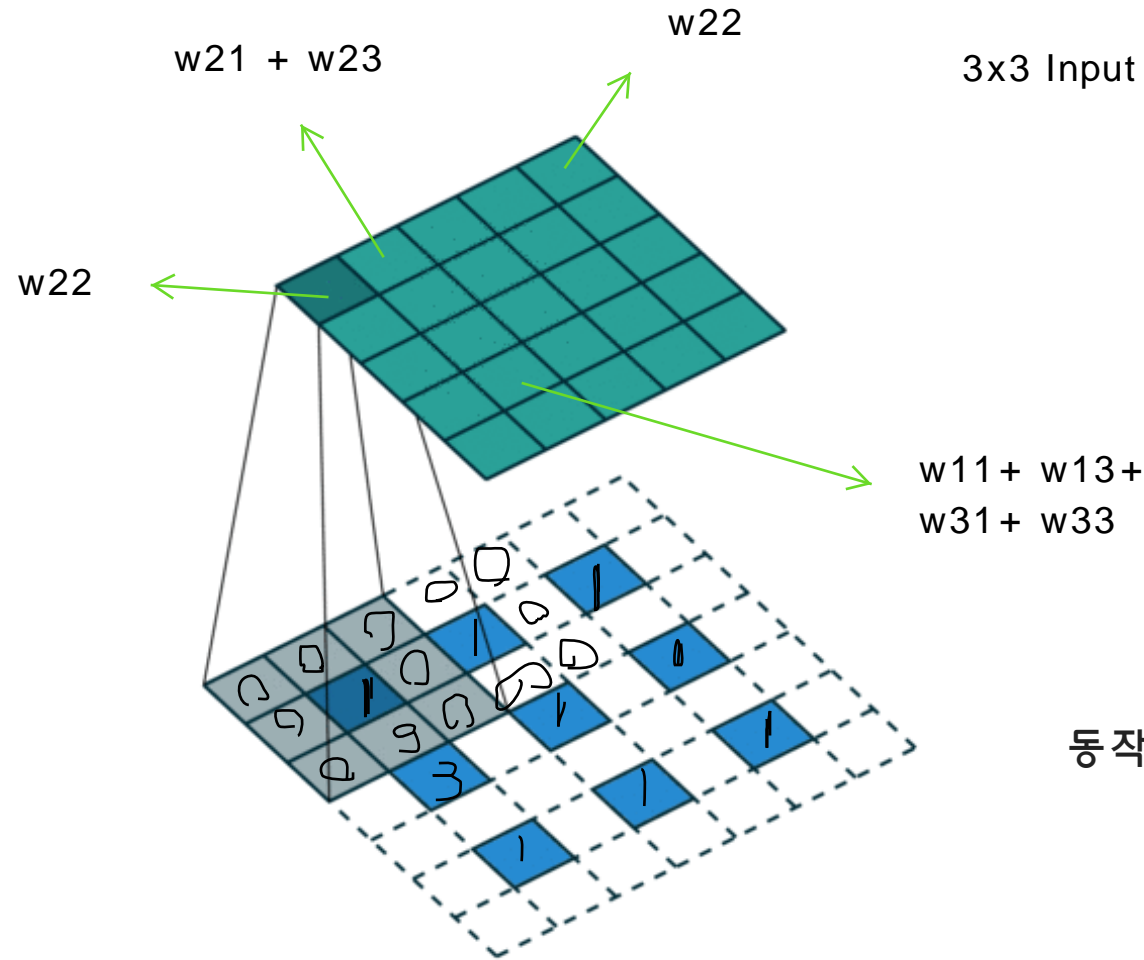
(Upsampling)

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



Transpose Convolution

(Upsampling)



Input

Kernel

Output

w_{11} w_{12} w_{13}
 w_{21} w_{22} w_{23}
 w_{31} w_{32} w_{33}

filter가

input

동작 방식

filter가
w

filter
output

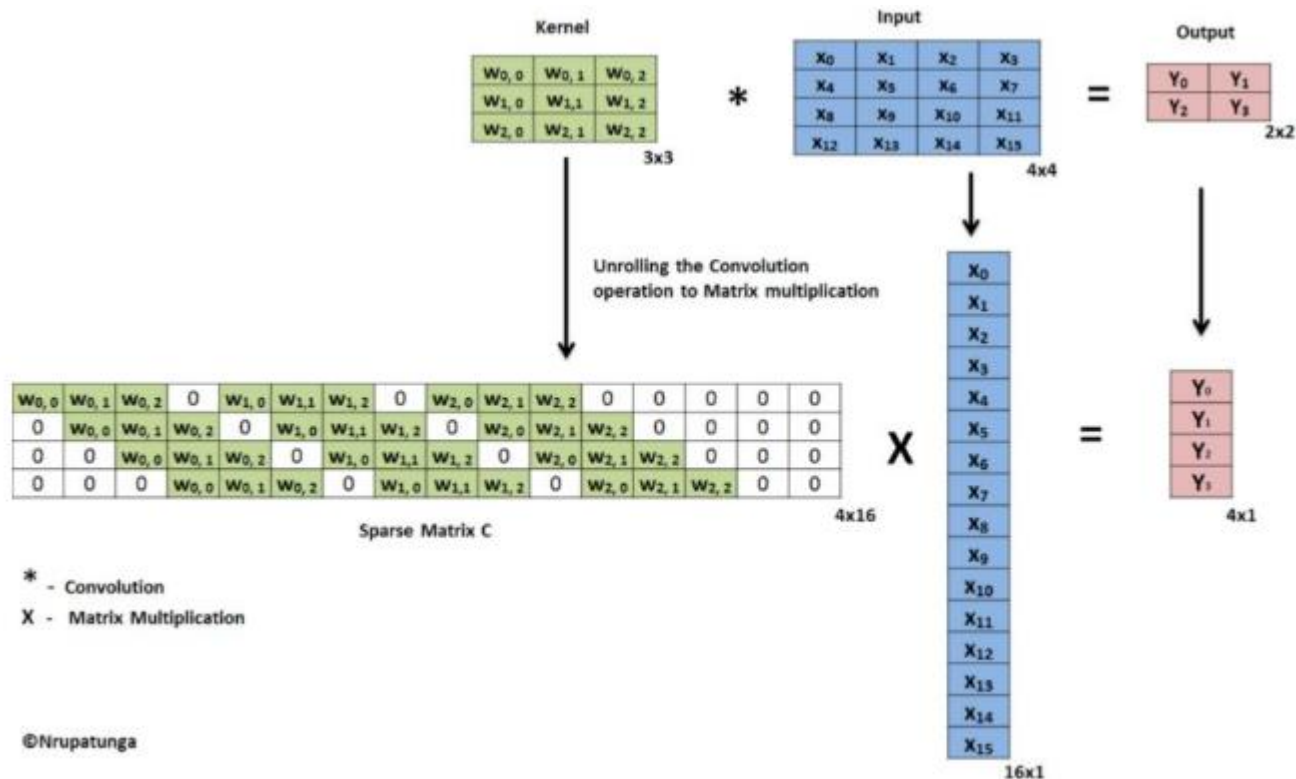
input

- 1) Input 내 각 픽셀 주위에 zero-padding
- 2) Zero-padding된 input에 convolution 연산

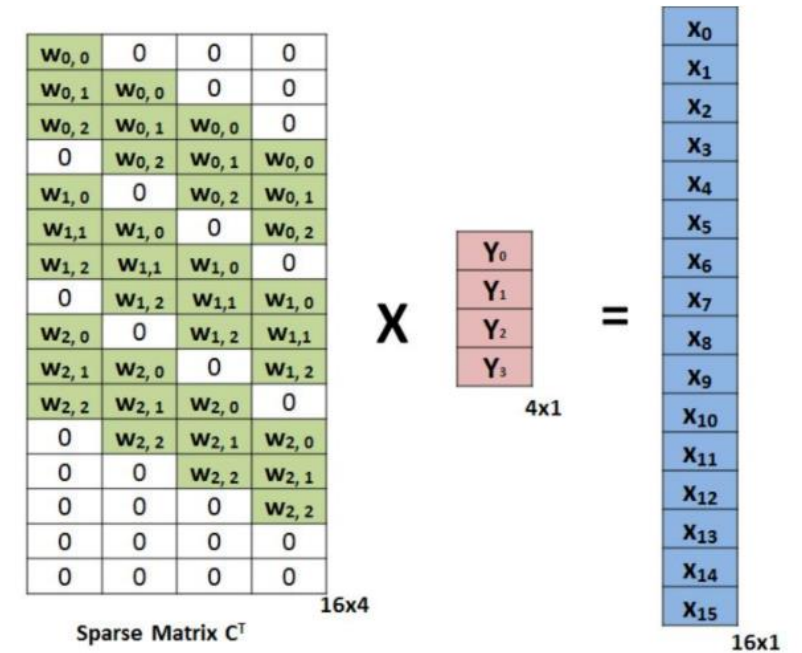
Transpose Convolution

(Upsampling)

• Convolution



• Deconvolution



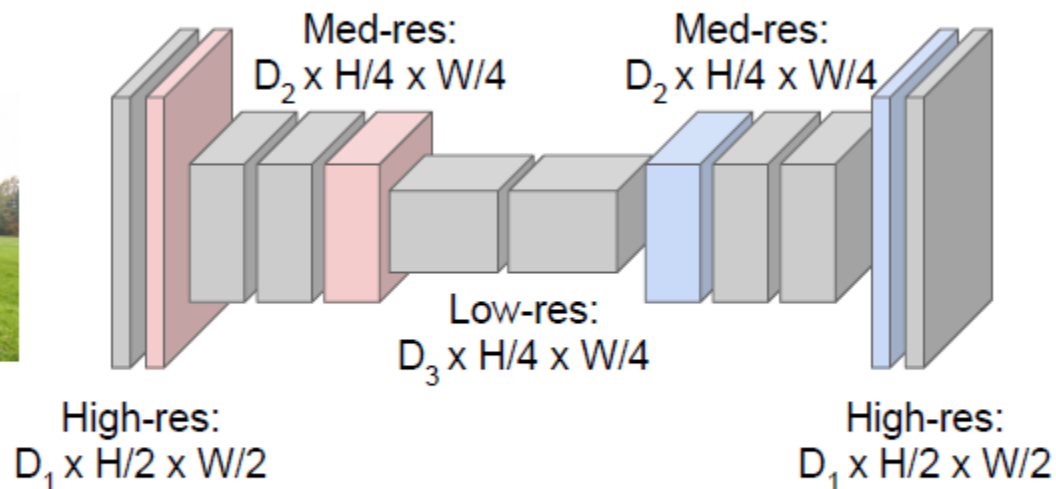
Semantic Segmentation

Downsampling:
Pooling, strided
convolution



Input:
 $3 \times H \times W$

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



Upsampling:
Unpooling or strided
transpose convolution



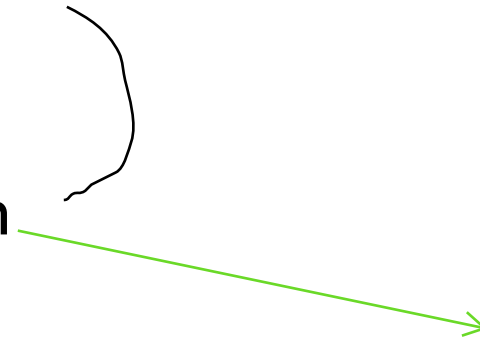
Predictions:
 $H \times W$

Several types of convolution

1) Convolution

1, 3, 4

2) Transpose Convolution



3) Atrous Convolution

4) Separable Convolution

가 . 가 .

ex) mobileNet - 100

가
main.py - 가
train.py - github +
Object Detection - Apple AI

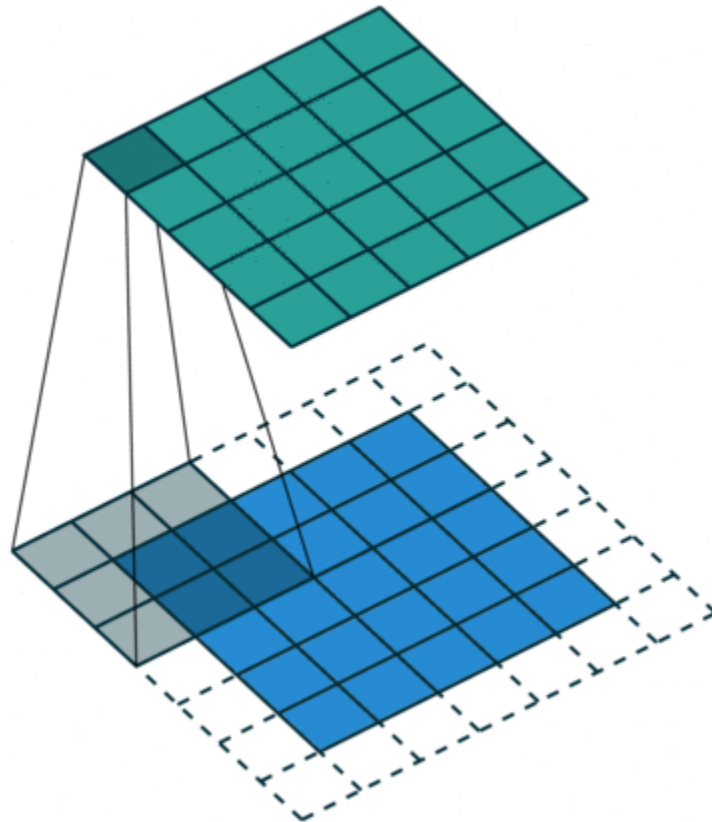
Several types of convolution

1) Convolution

- Down-sampling

DeConvolution

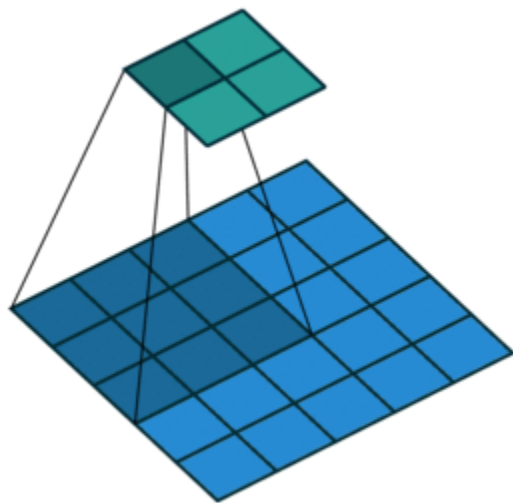
가 .



Several types of convolution

2) Transpose Convolution

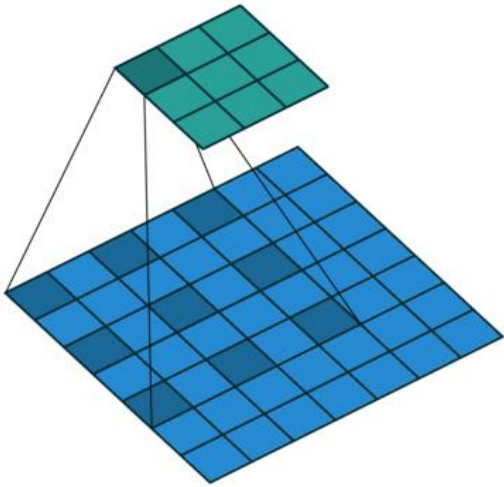
- Up-sampling
- Checker boarder Issue



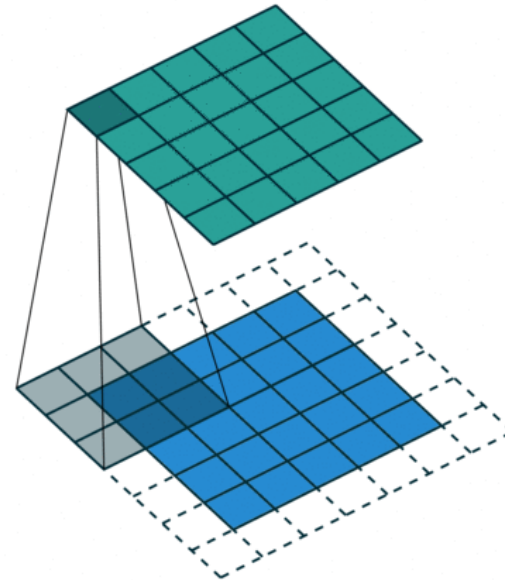
Several types of convolution

3) Atrous Convolution

- Down-sampling
- Wider Field of view at the same computational cost
- Use it when you need a wide field of view (Not good if field of view is too small)



Atrous convolution

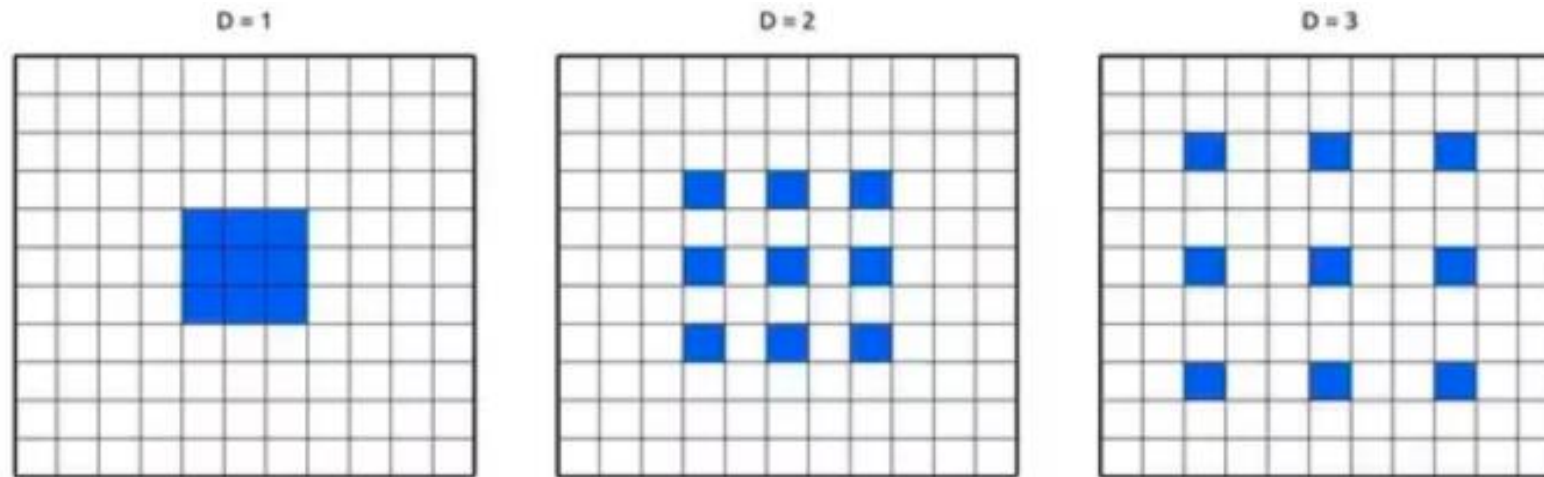


Convolution

Several types of convolution

3) Atrous Convolution

다양한 Dilated rate를 이용하여 병렬적으로 사용해서 더 많은 특징 추출가능
(ex. Deeplab v3+)



Atrous Convolution. 왼쪽부터 dilation rate: 1, 2, 3

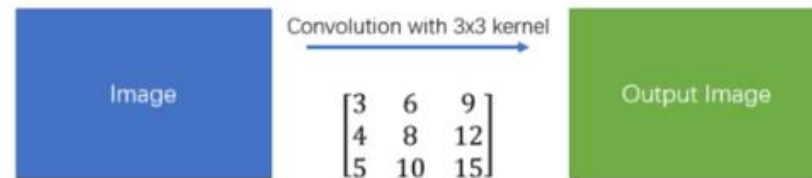
Several types of convolution

4) Separable Convolution

기존의 Convolution과 다르게 공간과 관련된 Convolution과 채널과 관련된 Convolution을 따로 적용하여 기존 Convolution을 표현할 수 있으면서도 파라미터 수를 낮추는 Convolution

(Guo et al. Network Decoupling: From Regular to Depthwise Separable Convolutions)

Simple Convolution



Spatial Separable Convolution



Several types of convolution

4) Separable Convolution

- Normal Convolution

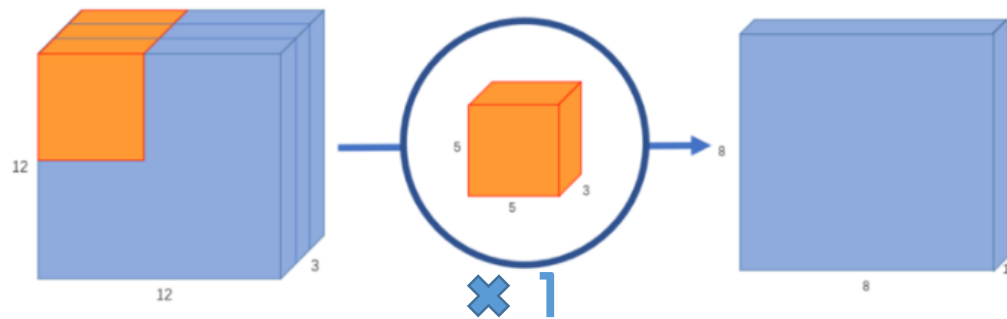


Image 4: A normal convolution with 8x8x1 output

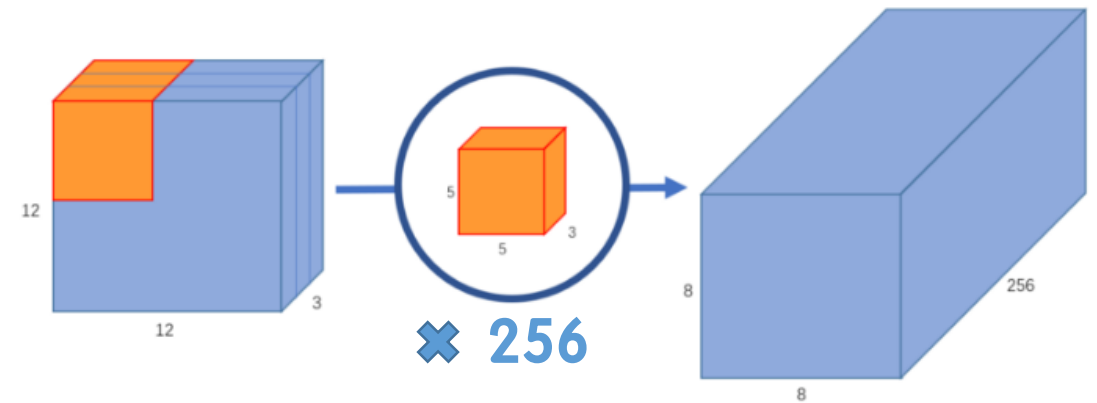
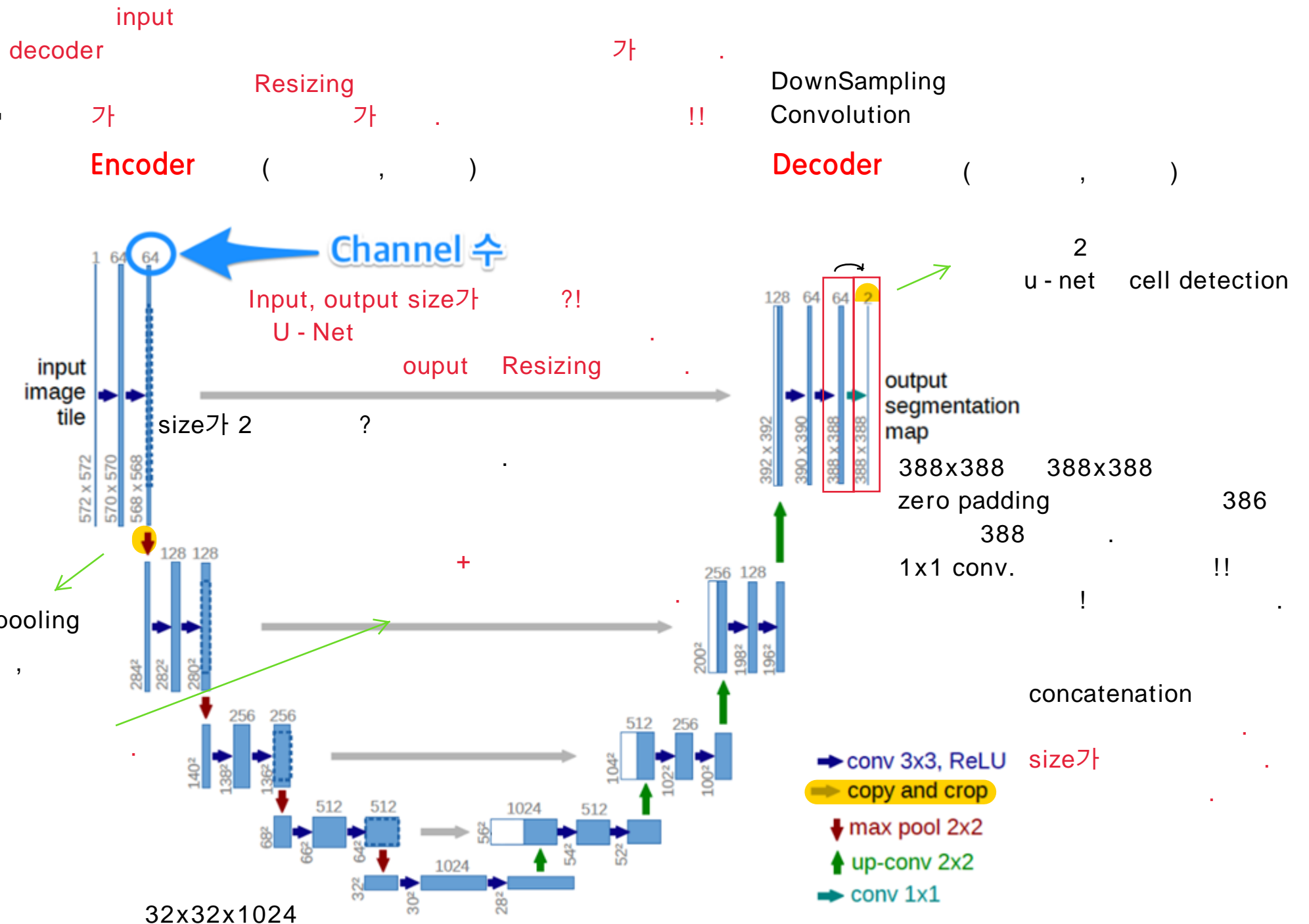




Image 5: A normal convolution with 8x8x256 output

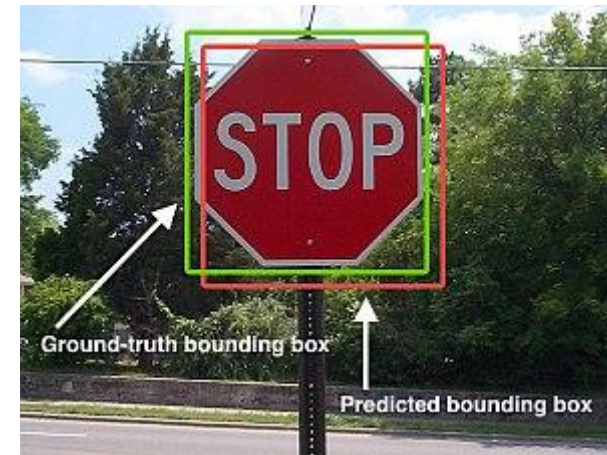
U-Net



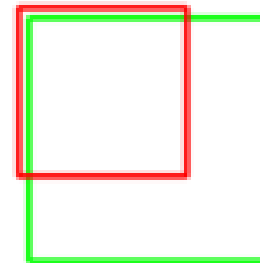
U-Net

- IOU (Intersection over union)


$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$


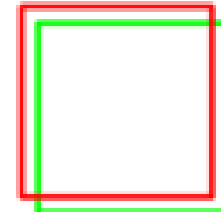


IoU: 0.4034



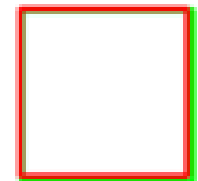
Poor

IoU: 0.7330



Good

IoU: 0.9264



Excellent

U-Net

- Checkerboard Artifacts on deconvolution



U-Net

- Deconv vs. Interpolation



Using deconvolution.
Heavy checkerboard artifacts.

deconvolution



Using resize-convolution.
No checkerboard artifacts.

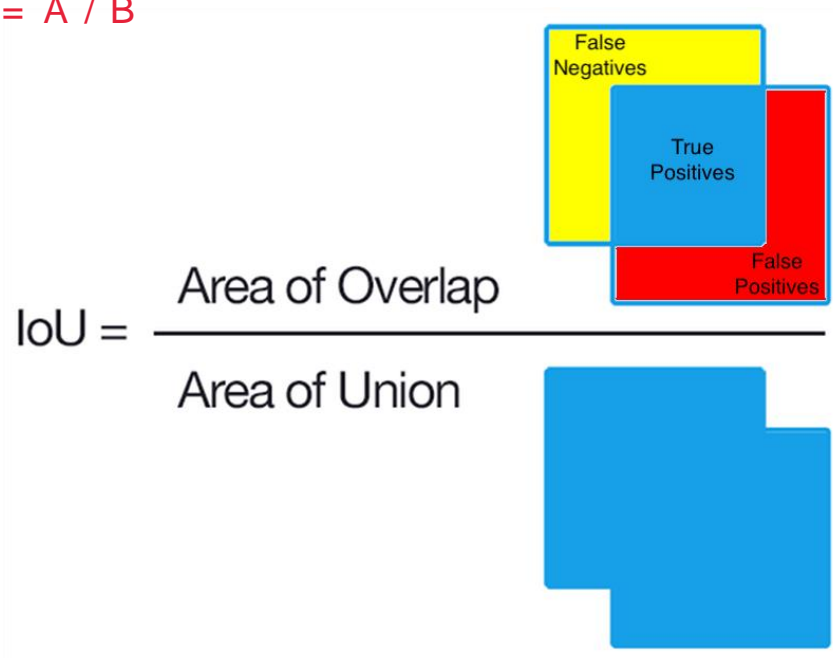
Interpolation

Evaluation Matrix

- IOU (Intersection over union)

$$DCE = 2A / (A+B)$$

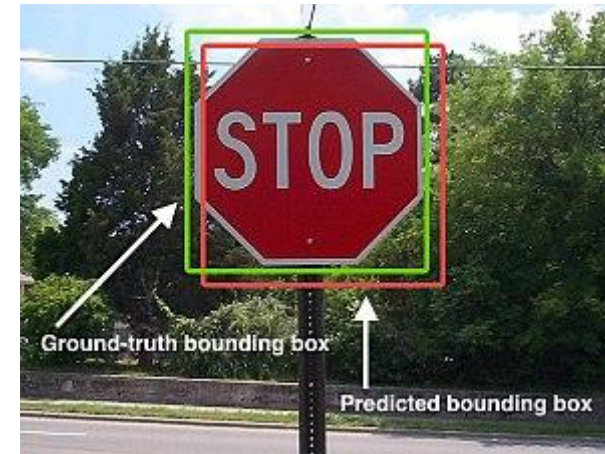
$$IOU = A / B$$



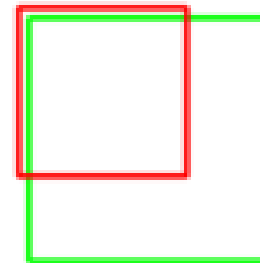
0.7

가 .

가

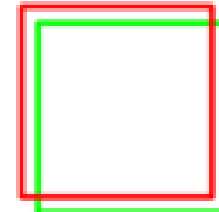


IoU: 0.4034



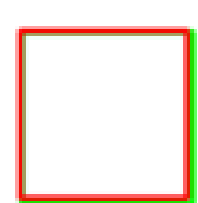
Poor

IoU: 0.7330



Good

IoU: 0.9264



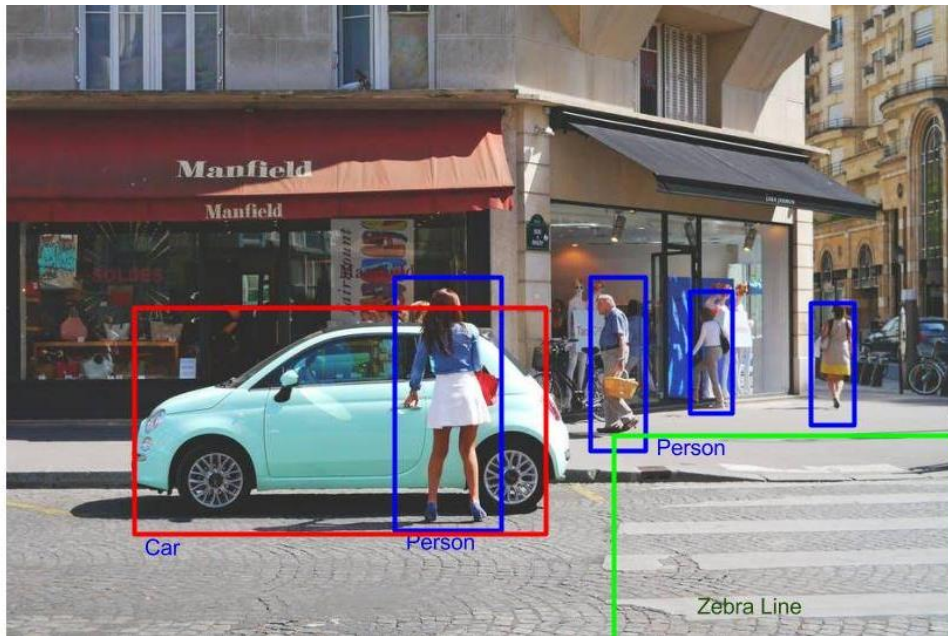
Excellent

Evaluation Matrix

- mIOU (Mean intersection over union)

IOU

Segmentation 가



Evaluation Matrix

- Dice Coefficient

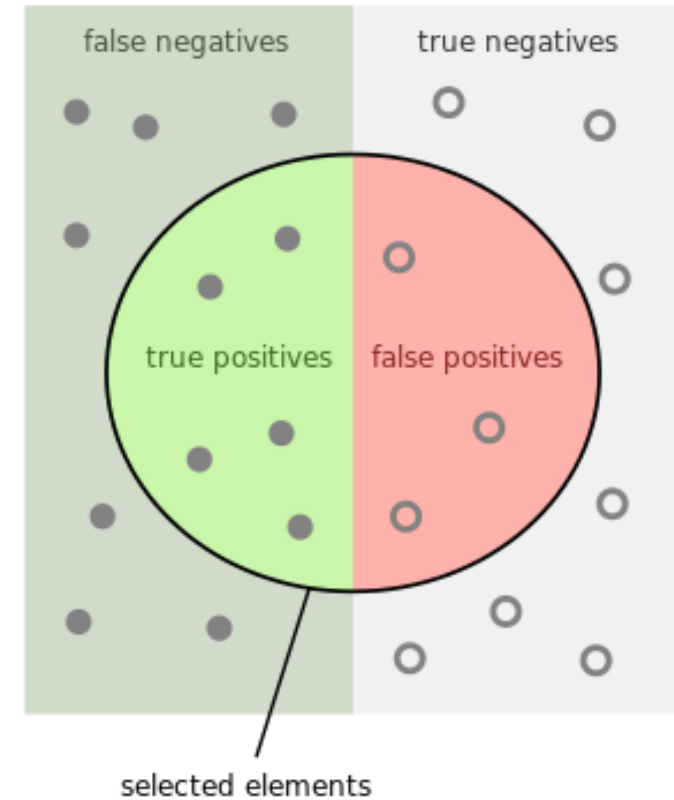
$$DSC = \frac{2|X \cap Y|}{|X| + |Y|}$$

$$DSC = \frac{2TP}{2TP + FP + FN}$$

+ : mIOU
+ : DCE

가

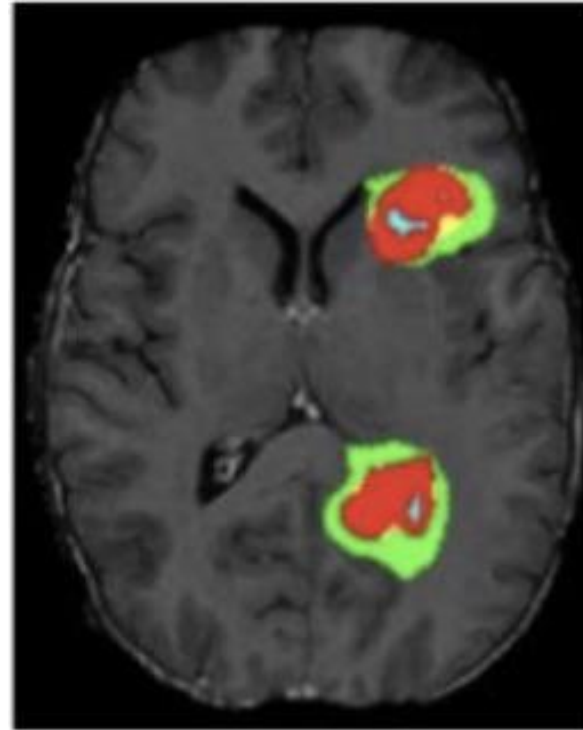
.!!



Evaluation Matrix

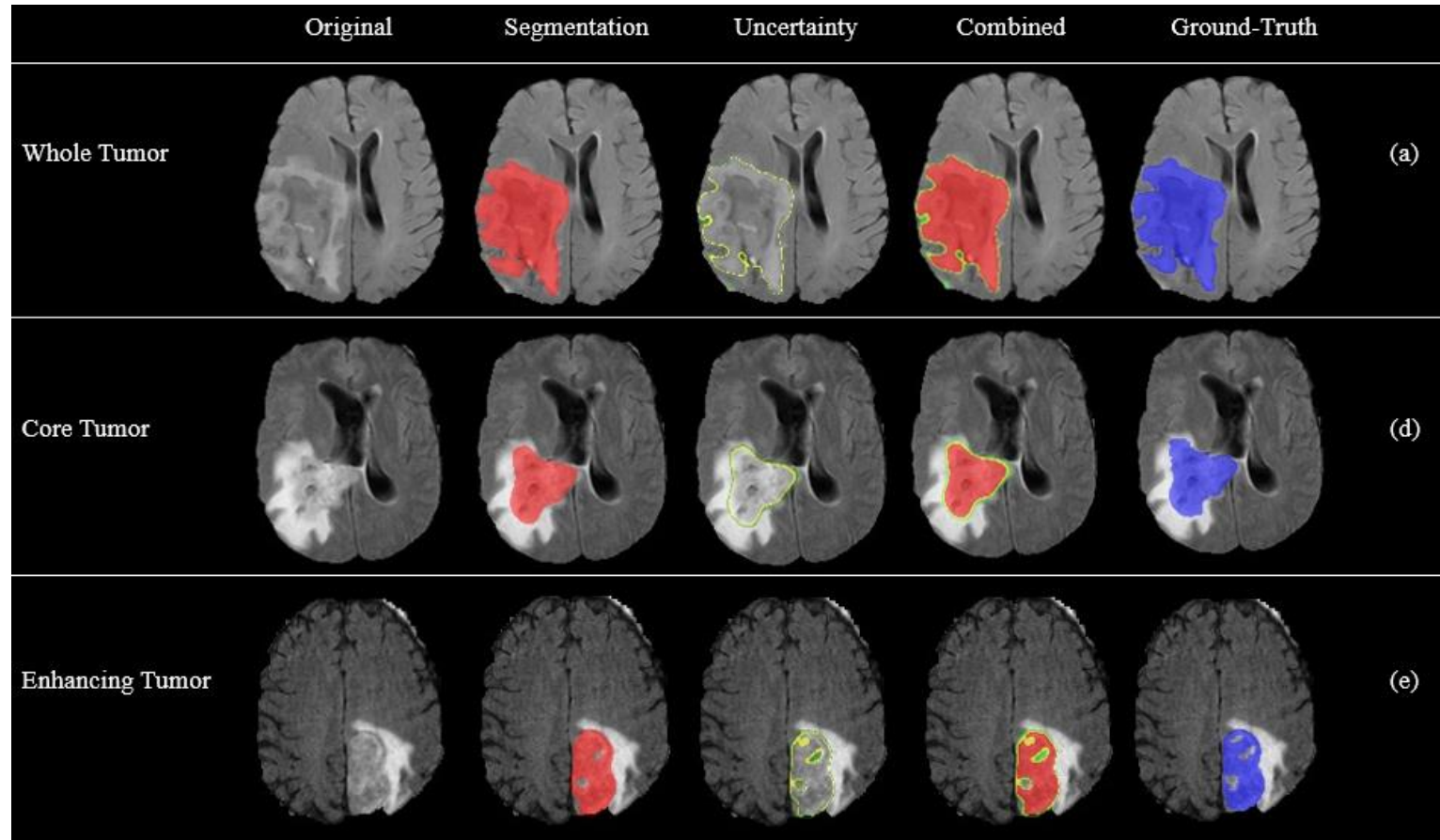
- Dice Coefficient

$$DSC = \frac{2TP}{2TP + FP + FN}.$$

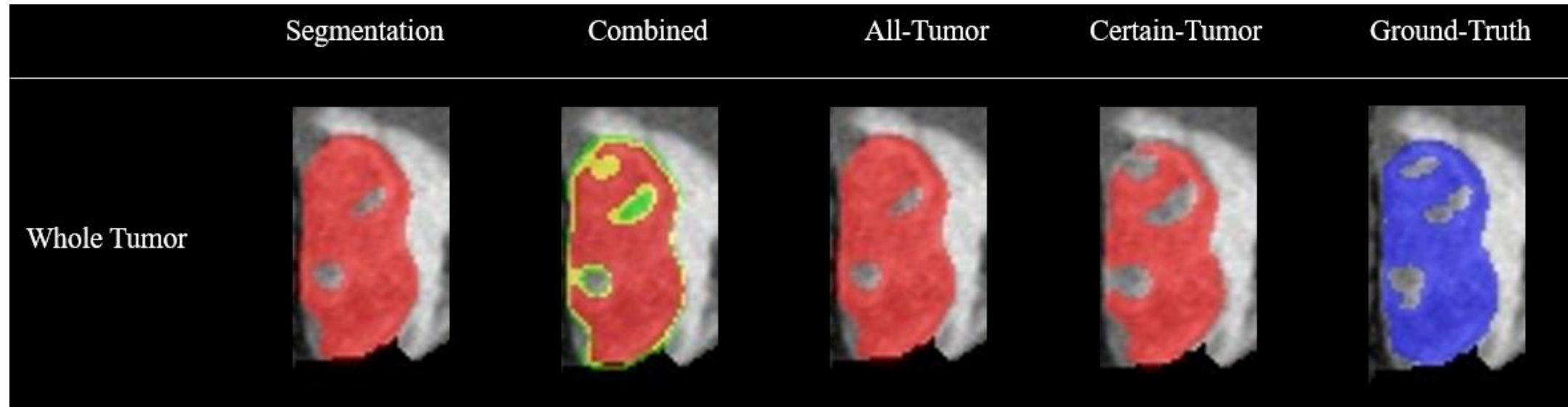


Very Effective to train Imbalanced dataset

Uncertainty Quantification



Uncertainty Quantification



Let's try it !

(code)