

# **Linear Algebra Preparation for Emergent Neural Network Architectures (LAPENNA)**

## **Unit 8**

### **Convolutional neural Network III**

**Kwai Wong, Stan Tomov,  
Julian Halloy, Stephen Qiu, Eric Zhao**

**March 24, 2022**

**University of Tennessee, Knoxville**

# Acknowledgements:

- Support from NSF, UTK, JICS, ICL, NICS
- LAPENNA, [www.jics.utk.edu/lapenna](http://www.jics.utk.edu/lapenna), NSF award #202409
- [www.icl.utk.edu](http://www.icl.utk.edu), [cfdlab.utk.edu](http://cfdlab.utk.edu), [www.xsede.org](http://www.xsede.org),  
[www.jics.utk.edu/recsem-reu](http://www.jics.utk.edu/recsem-reu),
- MagmaDNN is a project grown out from the RECSEM REU Summer program supported under NSF award #1659502
- Source code: [www.bitbucket.org/icl/magmadnn](http://www.bitbucket.org/icl/magmadnn)
- [www.bitbucket.org/cfdl/opendnnwheel](http://www.bitbucket.org/cfdl/opendnnwheel)



INNOVATIVE  
COMPUTING LABORATORY

THE UNIVERSITY OF  
TENNESSEE  
KNOXVILLE

**JICS**  
Joint Institute for  
Computational Sciences  
ORNL  
Computational Sciences

OAK RIDGE  
National Laboratory

**The major goal of The Linear Algebra Preparation for Emergent Neural Network Architectures (LAPENNA) program is to provide essential knowledge to advance literacy in AI to sustain the growth and development of the workforce in the cyberinfrastructure (CI) ecosystem. This program aims to prepare college faculty, researchers, and industrial practitioners to design, enable and direct their own course curricula, collaborative projects, and training programs for in-house data-driven sciences programs. The LAPENNA program focuses on delivering computational techniques, numerical algorithms and libraries, and implementation of AI software on emergent CPU and GPU platforms.**

Ecosystem Linear Algebra Preparation for Emergent Neural Network Architectures (LAPENNA)

Modeling, Numerical Linear Algebra, Data Analytics, Machine Learning, DNN, GPU, HPC

Session 1	Session 2	Session 3	Session 4
7/2020 - 12/2020	1/2021- 6/2021	7/2021 - 1/2022	1/2022 - 6/2022
16 participants	16 participants	16 participants	16 participants
Faculty/Students	Faculty/Students	Faculty/Students	Faculty/Students
10 webinars	10 webinars	10 webinars	10 webinars
4 Q & A	4 Q & A	4 Q & A	4 Q & A

Colleges courses, continuous integration, online courses, projects, software support

Web-based resources, tutorials, webinars, training, outreach

- ✓ PIs : **Kwai Wong (JICS), Stan Tomov (ICL), University of Tennessee, Knoxville**
  - Stephen Qiu, Julian Halloy, Eric Zhao (Students)
  
- ✓ Team : Clemson University
- ✓ Teams : University of Arkansas
- ✓ Team : University of Houston, Clear Lake
- ✓ Team : Miami University, Ohio
- ✓ Team : Boston University
- ✓ Team : West Virginia University
- ✓ Team : Louisiana State University, Alexandria
- ✓ Teams : Jackson Laboratory
- ✓ Team : Georgia State University
- ✓ Teams : University of Tennessee, Knoxville
- ✓ Teams : Morehouse College, Atlanta
- ✓ Team : North Carolina A & T University
- ✓ Team : Clark Atlanta University, Atlanta
- ✓ Team : Alabama A & M University
- ✓ Team : Slippery Rock University
- ✓ Team : University of Maryland, Baltimore County

- ✓ **Webinar Meeting time. Thursday 8:00 – 10:00 pm ET,**
- ✓ **Tentative schedule, [www.jics.utk.edu/lapenna](http://www.jics.utk.edu/lapenna) --> Spring 2022**

Topic: LAPENNA Spring 2022 Webinar

Time: Feb 3, 2022 08:00 PM Eastern Time (US and Canada)

Every week on Thu, 12 occurrence(s)

Feb 3, 2022 07:30 PM

Feb 10, 2022 07:30 PM

Feb 17, 2022 07:30 PM

Feb 24, 2022 07:30 PM

Mar 3, 2022 07:30 PM

Mar 10, 2022 07:30 PM

Mar 17, 2022 07:30 PM

Mar 24, 2022 07:30 PM

Mar 31, 2022 07:30 PM

Apr 7, 2022 07:30 PM

Apr 14, 2022 07:30 PM

Apr 21, 2022 07:30 PM

Join from PC, Mac, Linux, iOS or Android: <https://tennessee.zoom.us/j/94140469394>

Password: 708069

# Schedule of LAPENNA Spring 2022

## Thursday 8:00pm -10:00pm Eastern Time



The goal of The Linear Algebra Preparation for Emergent Neural Network Architectures (LAPENNA) program is to provide essential knowledge to advance literacy in AI to sustain the growth and development of the workforce in the cyberinfrastructure (CI) ecosystem for data-driven applications.

<b>Month</b>	<b>Week</b>	<b>Date</b>	<b>Topics</b>
<b>February</b>	Week 01	3	Logistics, High Performance Computing
	Week 02	10	Computational Ecosystem, Linear Algebra
	Week 03	17	Introduction to DNN, Forward Path (MLP)
	Week 04	24	Backward Path (MLP), Math, Example
<b>March</b>	Week 05	3	Backpro (MLP), CNN Computation
	Week 06	10	CNN Backpropagation, Example
	Week 07	17	CNN Network, Linear Algebra
	Week 08	24	Segmentation, Unet,
<b>April</b>	Week 09	31	Object Detection, RC vehicle
	Week 10	7	RNN, LSTEM, Transformers
	Week 11	14	DNN Computing on GPU
<b>June or July</b>	Week 12	21	Overview, Closing
	Workshop	To be arranged	4 Days In Person at UTK

✓ **UNIT 7 : Convolutional Neural Network (CNN)**

- Overview Forward and Backpropagation of MLP
- CNN computation
- CNN forward path calculation
- MNIST TensorFlow code
- CIFAT 10 example
- AlexNet, number of parameters
- CNN Backward Calculations
- Homework exercises

# DNN Model

Applications

+

Data Ensemble + Input

+

**It's all about  
linear algebra calculations**

**DNN in particular**

+

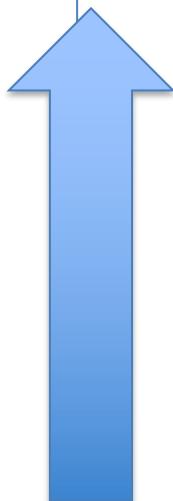
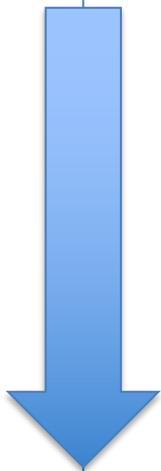
Algorithms, Software

+

Output + Data Analysis

+

Hardware



# Machine Learning – GPU acceleration



Image Classification

Object Detection

COMPUTER VISION

Voice Recognition

Translation

SPEECH AND AUDIO

Recommendation Engines

Sentiment Analysis

BEHAVIOR



cuDNN

DEEP LEARNING



cuBLAS

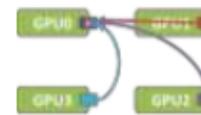


cuSPARSE



cuFFT

MATH LIBRARIES



NCCL

MULTI-GPU

**Supervised learning** is a type of machine learning which automate decision-making processes by generalizing from known examples. That is, the users provides the algorithm with pairs of inputs and desired outputs. The algorithm finds a way to produce the desired output given an input. In other words, we construct a model for the problem. The model is governed by some unknown parameters which we will learn from the data.

Supervised Learning Data:  
 $(x, y)$   $x$  is data,  $y$  is label

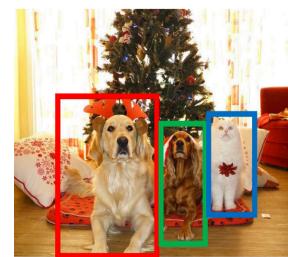
Goal: Learn a function to map  $x \rightarrow y$

Examples:  
Classification (discrete values),  
regression (numerical values),



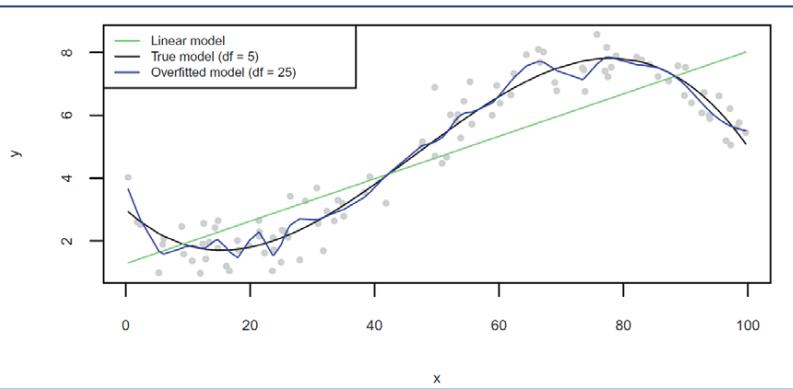
CAT

Classification



DOG, DOG, CAT

Object Detection



GRASS, CAT, TREE, SKY

Semantic Segmentation

# Basic Ideas

## Typical Neural Network – MLP

## Convolutional Neural Network

STEP 1 : Model Definition

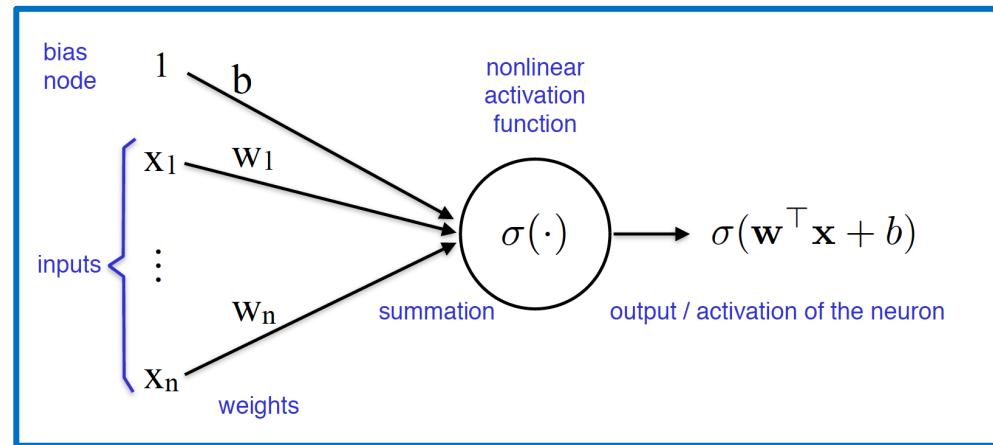
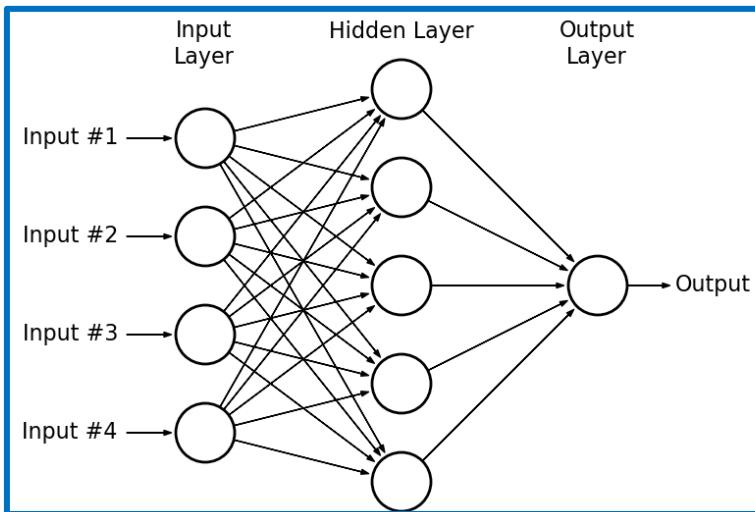
STEP 2 : Cost Function

Step 3 : Optimization Scheme

Step 4 : Numerical Implementation (fitting)

Step 5 : Evaluation

# Parametric Model : NN Modeling

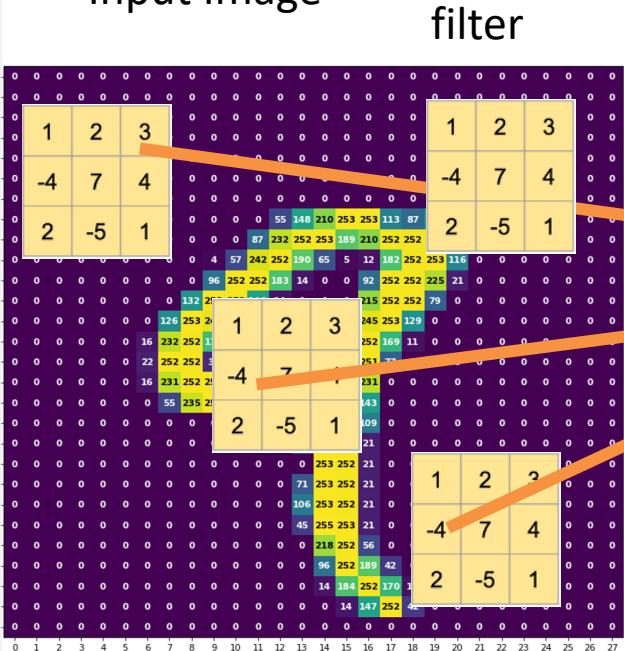


- ✓ A node in the neural network is a mathematical function or activation function which maps input to output values.
- ✓ Inputs represent a set of vectors containing weights ( $w$ ) and bias ( $b$ ). They are the sets of parameters to be determined.
- ✓ Many nodes form a **neural layer**, **links** connect layers together, defining a NN model.
- ✓ Activation function ( $f$  or  $\sigma$ ), is generally a nonlinear data operator which facilitates identification of complex features.

# Step 1: Parametric Model : Convolutional Neural Network (CNN)

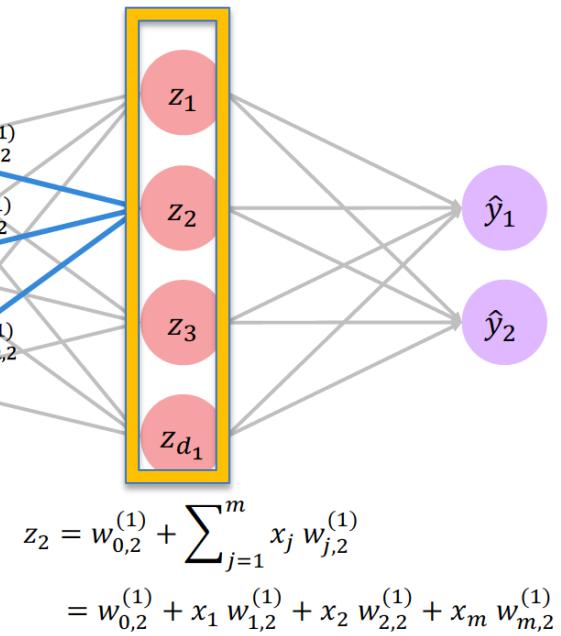
## Convolutional filter + Connected Neural Network

Input Image



Convolutional filtering

$$g(z)$$



Multilayer Perceptron

MLP : Neural Network

CNN : Convolutional NN

STEP 1 : Model Definition

**STEP 2 : Cost Function**

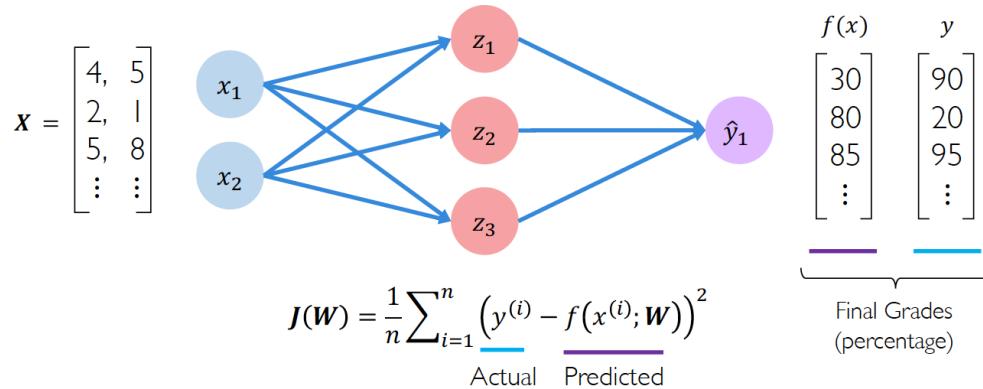
Step 3 : Optimization Scheme

Step 4 : Numerical Implementation

Step 5 : Evaluation

# NN Cost Function (regression) : Mean Squared Error Loss

Mean squared error loss can be used with regression models that output continuous real numbers

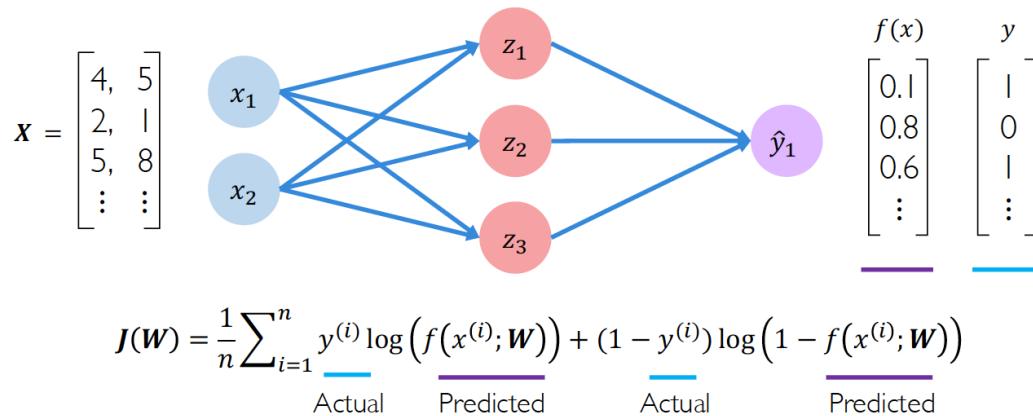


```
 loss = tf.reduce_mean( tf.square(tf.subtract(model.y, model.pred)) )
```

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \tilde{y}_i)^2$$

# NN Cost Function (Classification) : Cross Entropy Loss

Cross entropy loss can be used with models that output a probability between 0 and 1



```
 loss = tf.reduce_mean( tf.nn.softmax_cross_entropy_with_logits(model.y, model.pred) )
```

$$L_i = - \log \left( \frac{e^{s y_i}}{\sum_j e^{s_j}} \right)$$

MLP : Neural Network

CNN : Convolutional NN

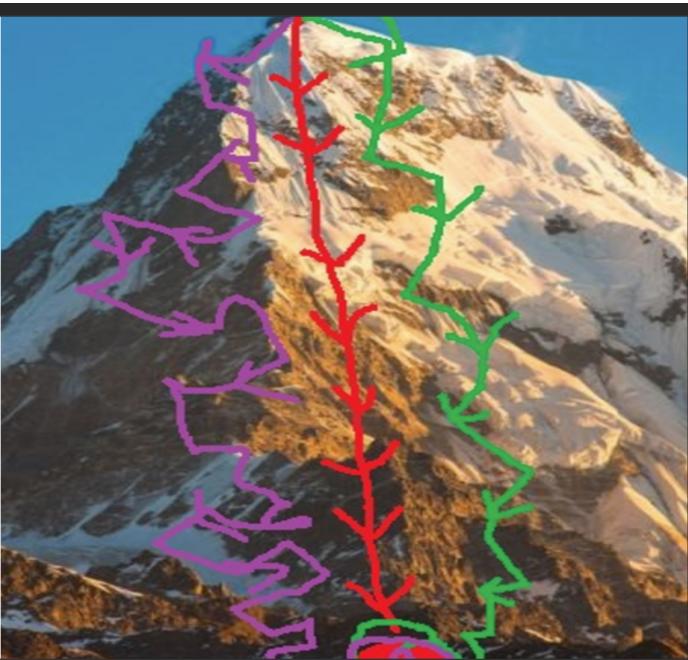
STEP 1 : Model Definition

STEP 2 : Cost Function

Step 3 : Optimization Scheme

Step 4 : Numerical Implementation

Step 5 : Evaluation



# Optimization Schemes: Gradient Descent : Simple Story

— Batch Gradient Descent  
 — Mini-batch Gradient Descent  
 — Stochastic Gradient Descent

Derivative is the slope of the tangent line

$$\frac{d}{dx}x^2 = 2x$$

$x = 1$

$x)$   
The slope of the tangent line is 2

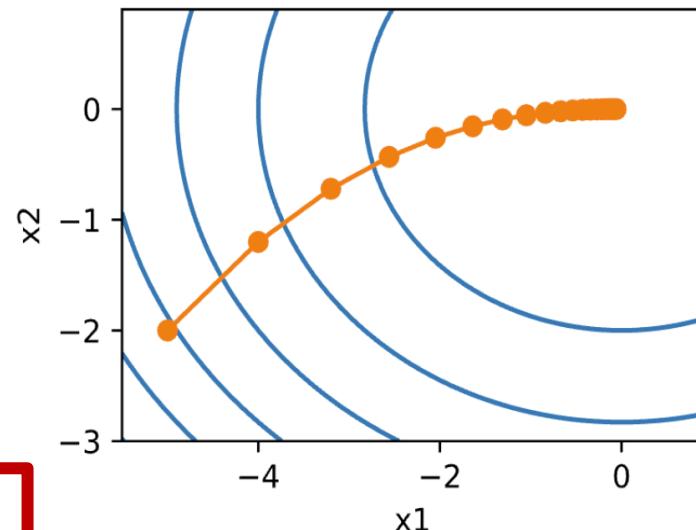
## Algorithm

- Choose initial  $\mathbf{x}_0$  ( $X = W$ , weight)
- At time  $t = 1, \dots, T$

$$\mathbf{x}_t = \mathbf{x}_{t-1} - \eta \nabla f(\mathbf{x}_{t-1})$$

- $\eta$  is called learning rate

$$W+ = W - (\text{learning rate}) * dJ/dw$$



# DNN Forward Backward Calculation : Weights and bias

## Algorithm

1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$



```
weights = tf.random_normal(shape, stddev=sigma)
```

2. Loop until convergence:

3. Compute gradient,  $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$



```
grads = tf.gradients(ys=loss, xs=weights)
```

4. Update weights,  $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$



```
weights_new = weights.assign(weights - lr * grads)
```

5. Return weights

## In one epoch

1. Pick a mini-batch
2. Feed it to Neural Network
3. Forward path calculation
4. Calculate the mean gradient of the mini-batch
5. Use the mean gradient calculated in step 4 to update the weights
6. Repeat steps 1–5 for the mini-batches we created

$$\frac{\sum_{j=1}^m \nabla C_{X_j}}{m} \approx \frac{\sum_x \nabla C_x}{n} = \nabla C,$$

$$w_k \rightarrow w'_k = w_k - \eta \frac{\partial C}{\partial w_k}$$
$$b_l \rightarrow b'_l = b_l - \eta \frac{\partial C}{\partial b_l}.$$

# Optimization: Batch SGD

Instead of computing a gradient across the entire dataset with size  $N$ , divides the dataset into a fixed-size subset (“minibatch”) with size  $m$ , and computes the gradient for each update on this smaller batch:

( $N$  is the dataset size,  $m$  is the minibatch size)

$$\begin{aligned}\mathbf{g} &= \frac{1}{m'} \sum_{i=1}^{m'} \nabla_{\theta} E_i(\mathbf{x}^{(n_i)}, y^{(n_i)}, \theta), \\ \theta &\leftarrow \theta - \eta \mathbf{g},\end{aligned}$$

Given: training set  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$

Initialize all  $\Theta^{(l)}$  randomly (NOT to 0!)

Loop // each iteration is called an epoch

Loop // each iteration is a mini-batch

Set  $\Delta_{i,j}^{(l)} = 0 \quad \forall l, i, j$

(Used to accumulate gradient)

Sample  $m$  training instances  $\mathcal{X} = \{(\mathbf{x}'_1, y'_1), \dots, (\mathbf{x}'_m, y'_m)\}$  without replacement

For each instance in  $\mathcal{X}$ ,  $(\mathbf{x}_k, y_k)$ :

Set  $\mathbf{a}^{(1)} = \mathbf{x}_k$

Compute  $\{\mathbf{a}^{(2)}, \dots, \mathbf{a}^{(L)}\}$  via forward propagation

Compute  $\delta^{(L)} = \mathbf{a}^{(L)} - y_k$

Compute errors  $\{\delta^{(L-1)}, \dots, \delta^{(2)}\}$

Compute gradients  $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

Compute mini-batch regularized gradient  $D_{ij}^{(l)} = \begin{cases} \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ \frac{1}{m} \Delta_{ij}^{(l)} & \text{otherwise} \end{cases}$

Update weights via gradient step  $\Theta_{ij}^{(l)} = \Theta_{ij}^{(l)} - \alpha D_{ij}^{(l)}$

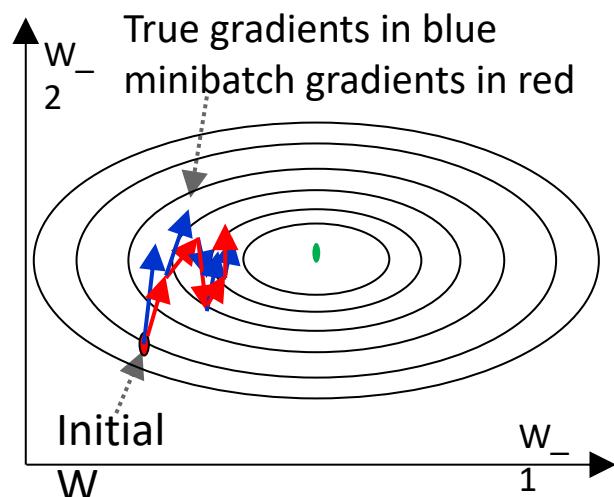
Until all training instances are seen

Until weights converge or max #epochs is reached

Epoch step



Batch iteration



MLP : Neural Network  
CNN : Convolutional NN

STEP 1 : Model Definition

STEP 2 : Cost Function

Step 3 : Optimization Scheme

Step 4 : Numerical Implementation

TensorFlow MNIST MLP & CNN Examples

# Training

“Training a neural network”

What does it mean?

What does the network train?

minimize error of the computed values against the  
labeled values (loss function)

What are the training parameters

$W$ ,  $b$  ,**and values of the filters ( $W_f$ )**

How does it work?

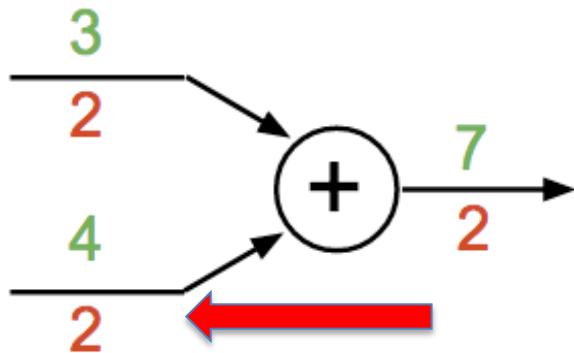
optimization based on gradient descent algorithm  
go back and forth in the NN model to adjust for  
 $w$  and  $b \Rightarrow$  need derivatives w.r.t. parameters

Go through forward calculation  
training with backpropagation

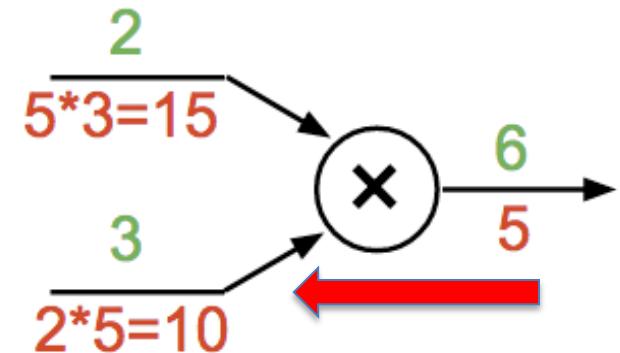
# Forward and Backpropagation Operators

## Forward path : backward path

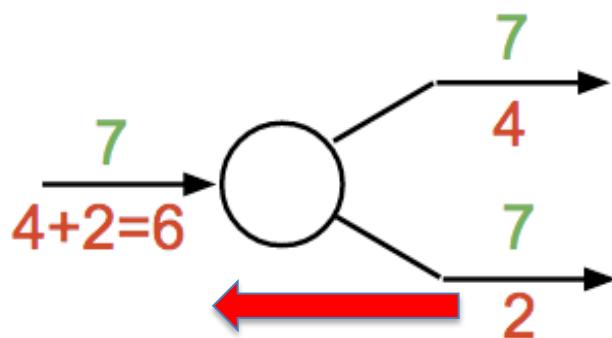
**add** gate: gradient distributor



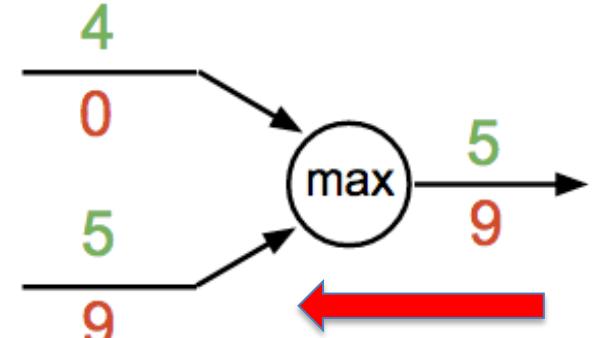
**mul** gate: “swap multiplier”



**copy** gate: gradient adder

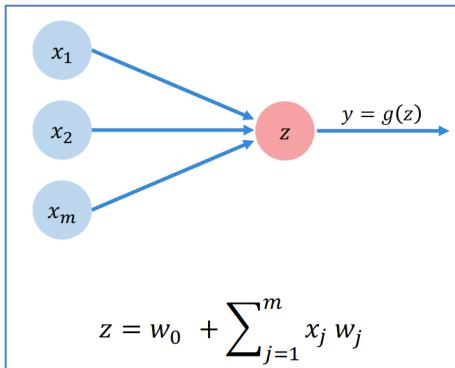
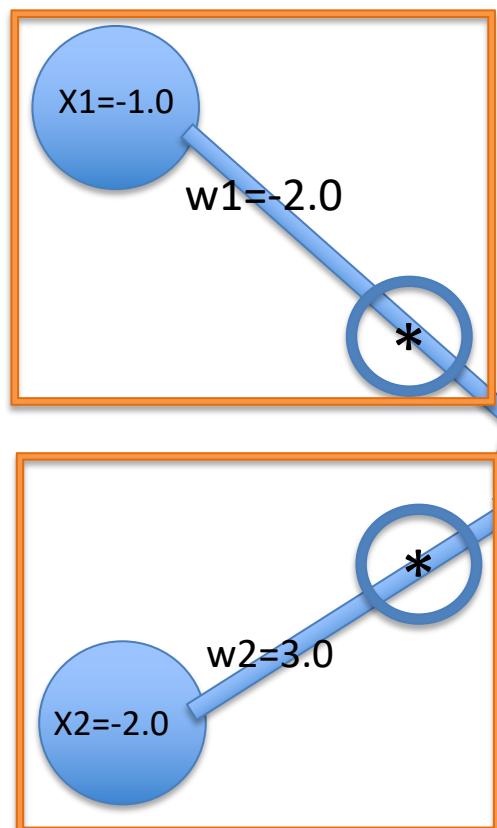


**max** gate: gradient router



# Forward and backward computation operators

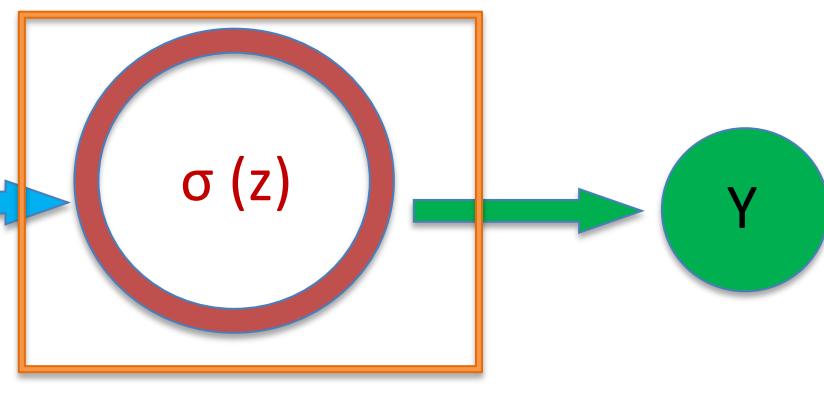
$$Z = (x_1 * w_1 + x_2 * w_2) + b = (-1.0 * -2.0 + -2.0 * 3.0) + -3.0 = 4.0 - 3.0 = 1.0$$



Sigmoid function

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$Y = \sigma(z) = 1 / (1 + e^{-z}) \\ = 1 / (1 + e^{-1}) = 0.731$$



multiple operator :  
downstream gradient =  
upstream gradient \* input value

add operator :  
Downstream gradient =  
upstream gradient

function operator :  
downstream gradient =  
Upstream gradient \* local function gradient

# Summary : Flow of NN

```
model.fit(x_train, y_train, epochs=3, batch_size=10)
```

Define Network

Forward Pass

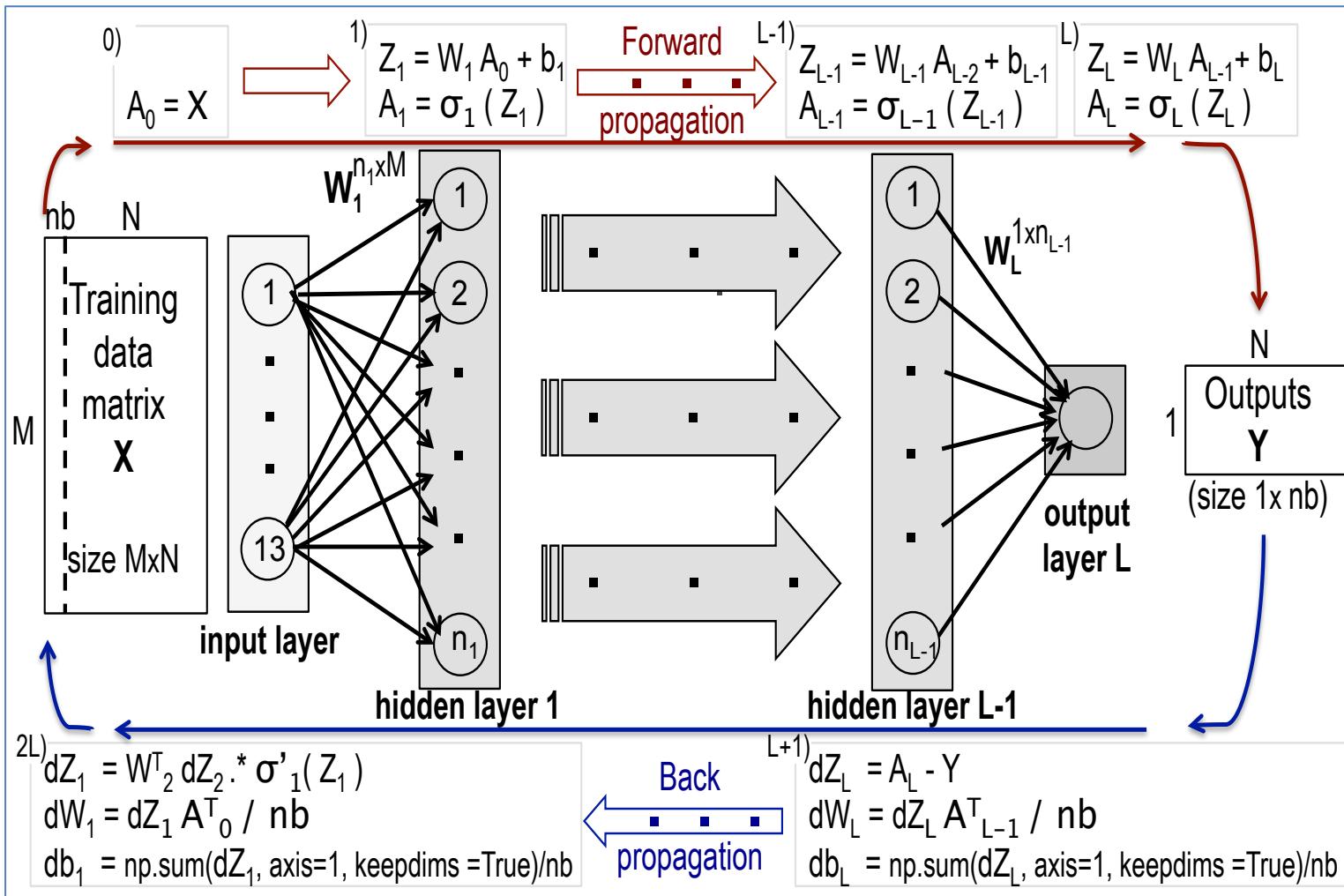
Calculate the analytical gradients

Update weights and bias

backpropagation

Gradient decent

Quantify loss



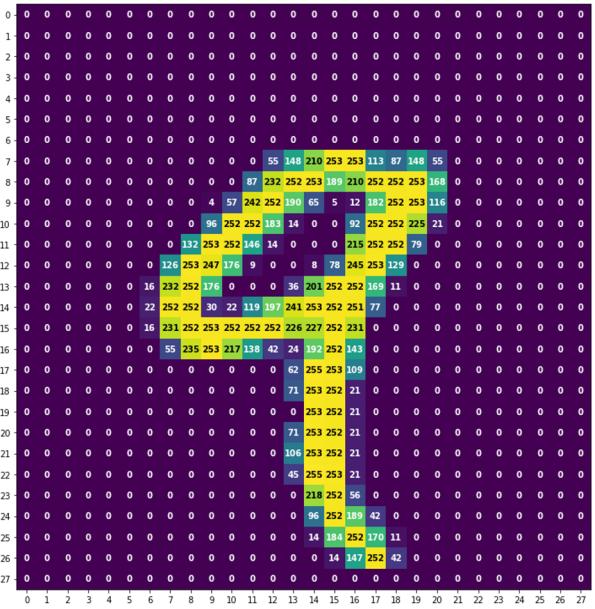
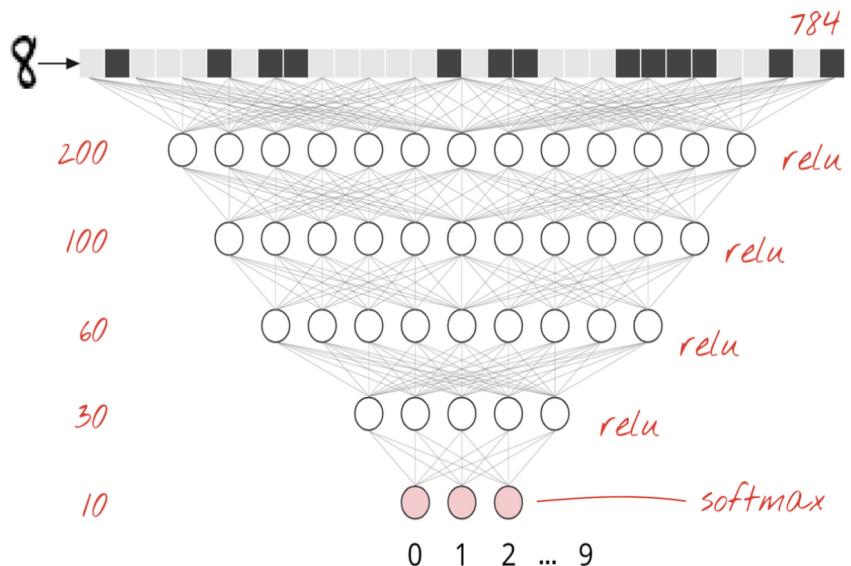
# MNIST Example (28x28 pixels)

Handwritten digits in the MNIST dataset are 28x28 pixel grayscale images. The simplest approach for classifying them is to use the  $28 \times 28 = 784$  pixels as inputs for a 1-layer neural network.

Image : input

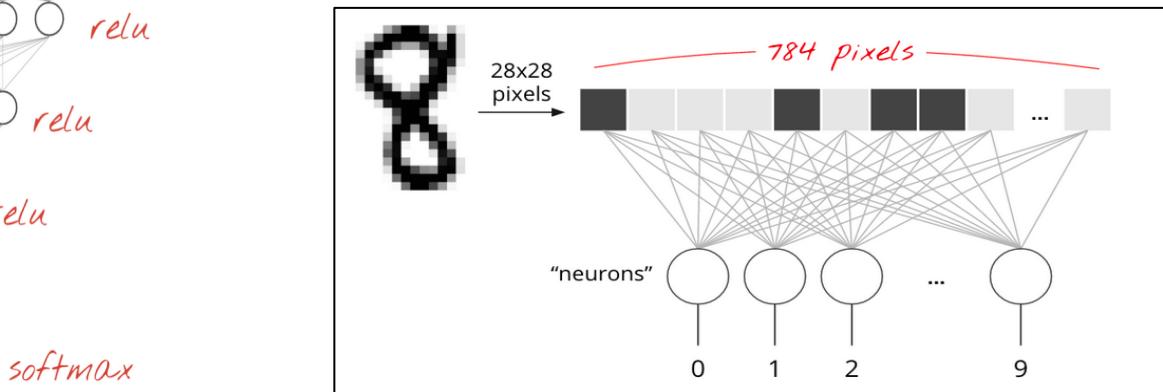
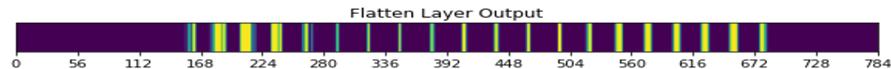
training digits and their labels	
5	9
9	0
1	5
5	0
4	2
4	5
7	5
0	6
7	5
0	0
0	0
8	8
8	0
6	9
9	1
validation digits and their labels	
7	2
1	0
4	1
9	5
9	0
6	9
0	1
5	9
7	3
4	9
9	6
6	5

Labels : target



Flatten the 28 x 28 matrix

to a vector of 28 x 28 elements 784 elements



TensorFlow, Keras and deep learning, without a PhD

<https://codelabs.developers.google.com/codelabs/cloud-tensorflow-mnist#0>

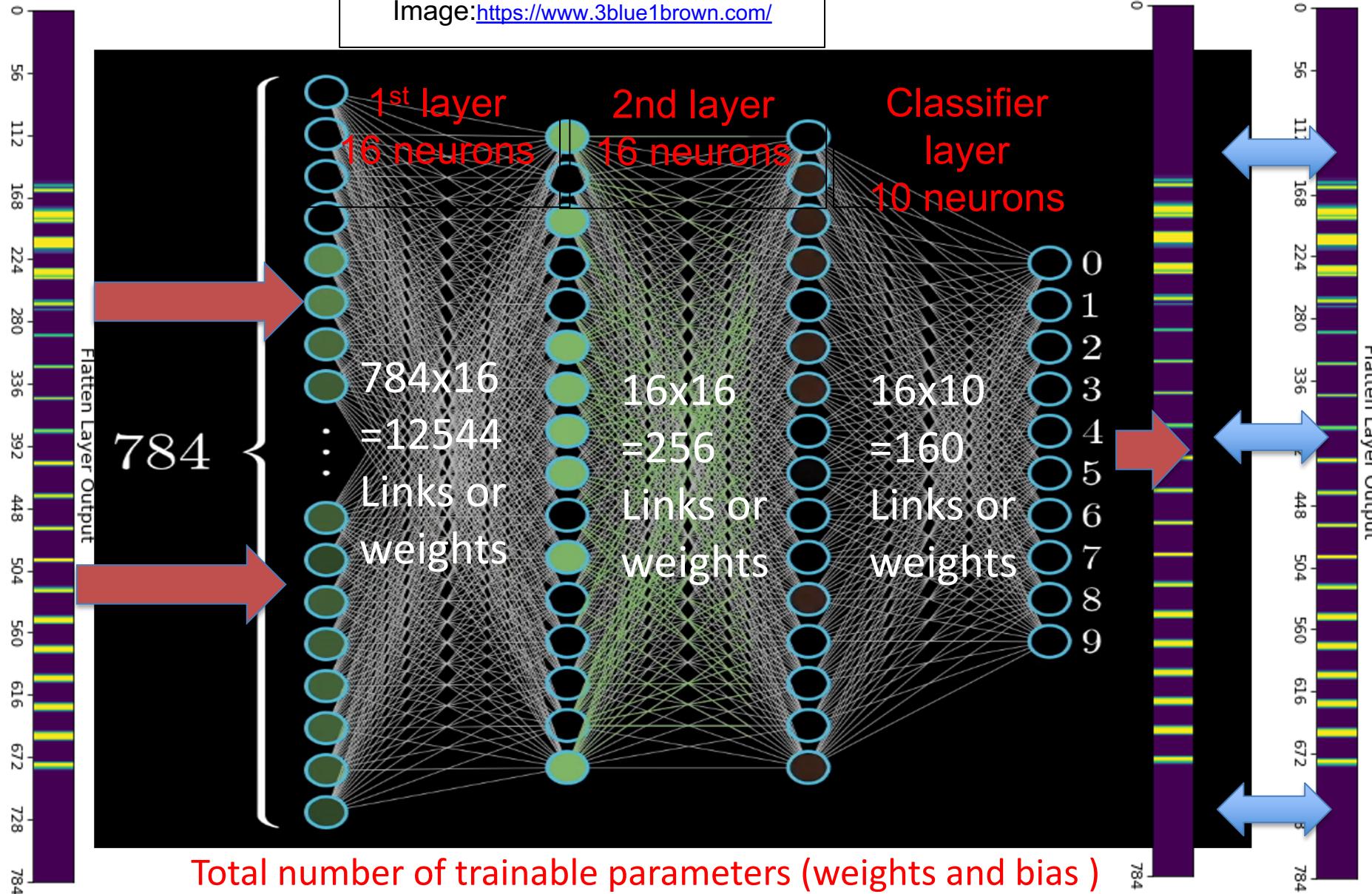
Input

# Simple MNIST MLP Neural Network

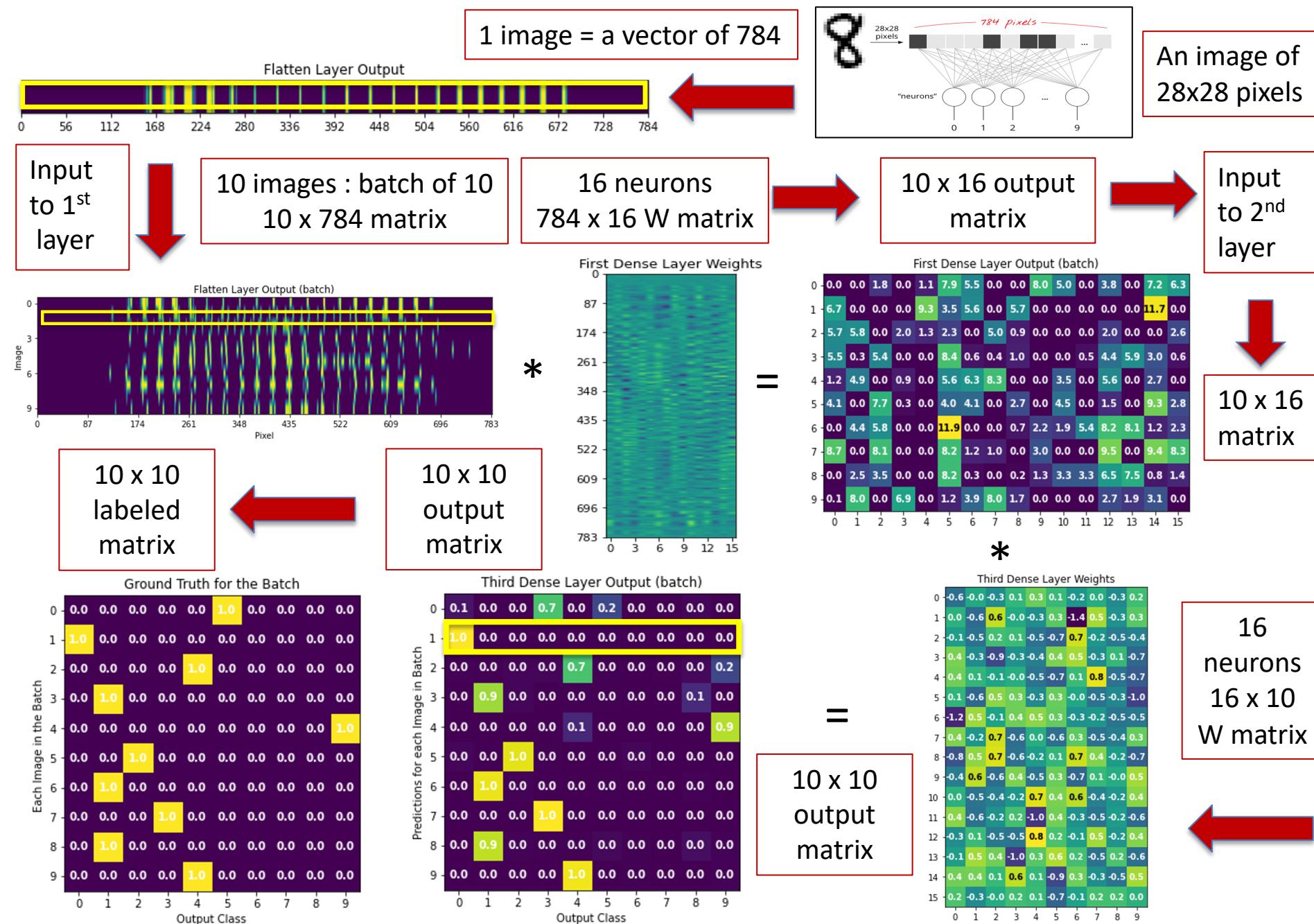
output

labels

Image:<https://www.3blue1brown.com/>



# Summary : Flow of MLP Dense layer NN



# Tensorflow Example

```
import tensorflow as tf

mnist = tf.keras.datasets.mnist

(x_train, y_train), (x_test, y_test) = mnist.load_data()
x_train, x_test = x_train / 255.0, x_test / 255.0

model = tf.keras.models.Sequential([
    tf.keras.layers.Flatten(input_shape=(28, 28)),
    tf.keras.layers.Dense(16, activation='relu'),
    tf.keras.layers.Dense(16, activation='relu'),
    tf.keras.layers.Dense(10, activation='softmax')
])
```

Tensorflow is imported

Sets the mnist dataset to variable “mnist”

Loads the mnist dataset

Builds the layers of the model  
4 layers in this model

```
model.compile(optimizer='sgd', loss='SparseCategoricalCrossentropy',
metrics=['accuracy'])

model.summary()

logdir = os.path.join("logs", datetime.datetime.now().strftime("%Y%m%d-%H%M%S"))
tensorboard_callback = tf.keras.callbacks.TensorBoard(logdir, histogram_freq=1)
```

Loss Function

Compiles the model with the SGD optimizer

Print summary

Use tensorborad

```

model.compile(optimizer='sgd',
              loss='SparseCategoricalCrossentropy',
              metrics=['accuracy'])

model.summary()

logdir = os.path.join("logs", datetime.datetime.now().strftime("%Y%m%d-%H%M%S"))
tensorboard_callback = tf.keras.callbacks.TensorBoard(logdir, histogram_freq=1)

model.fit(x_train, y_train, epochs=20, batch_size=50, validation_data=(x_test, y_test),
          callbacks=[tensorboard_callback])

model.evaluate(x_test, y_test, verbose=2)

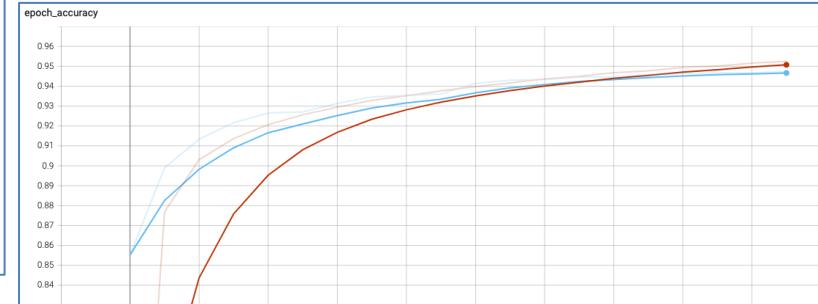
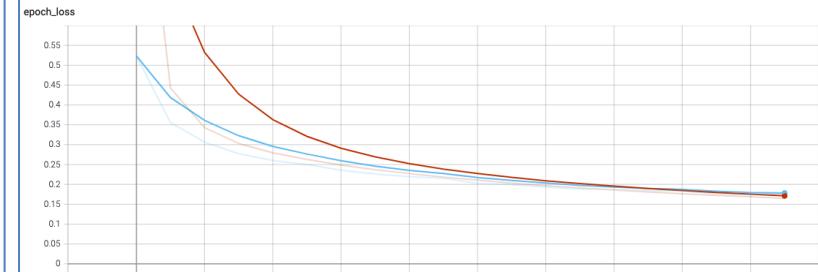
%tensorboard --logdir logs

```

↳ 2.4.1  
The tensorboard extension is already loaded. To reload it, use:  
%reload\_ext tensorboard  
Model: "sequential\_1"

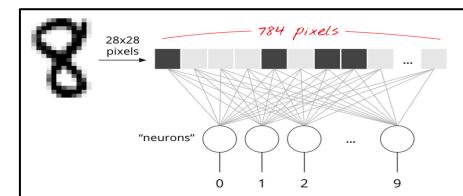
Layer (type)	Output Shape	Param #
flatten_1 (Flatten)	(None, 784)	0
dense_3 (Dense)	(None, 16)	12560
dense_4 (Dense)	(None, 16)	272
dense_5 (Dense)	(None, 10)	170

Total params: 13,002  
Trainable params: 13,002  
Non-trainable params: 0



# Summary : Flow of MLP Dense layer NN

1 image = a vector of 784

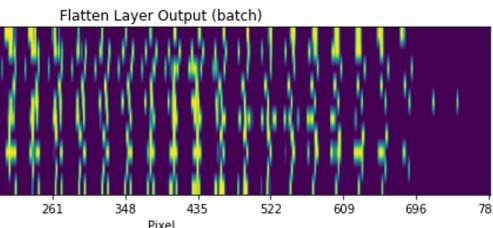


An image of 28x28 pixels

Input to 1<sup>st</sup> layer

10 images : batch of 10  
10 x 784 matrix

10 x 10 labeled matrix



10 x 10 output matrix

Ground Truth for the Batch

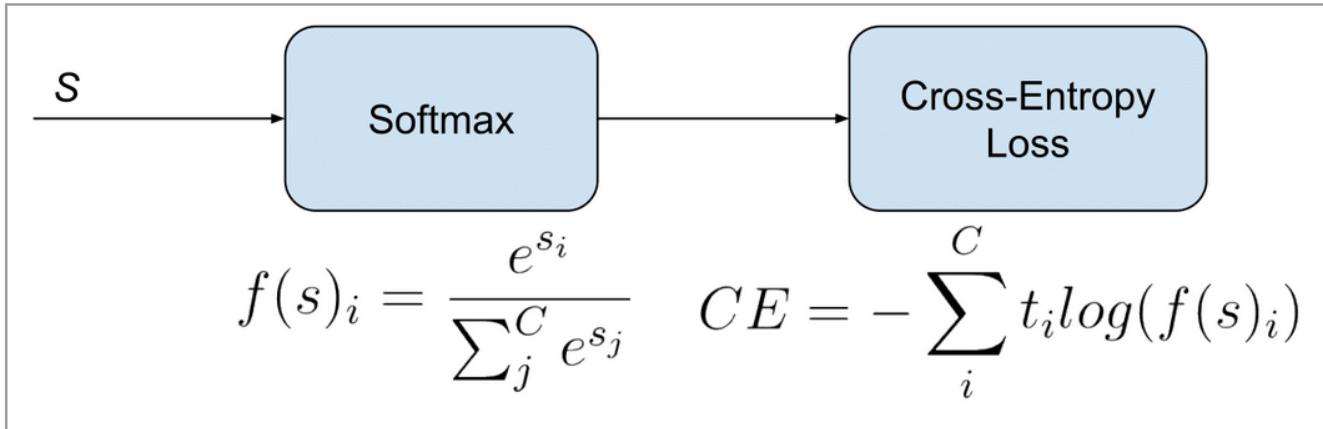
Each Image in the Batch	0	1	2	3	4	5	6	7	8	9	Output Class
0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	4
1	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1
2	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	3
3	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0
4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	9
5	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	2
6	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1
7	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	3
8	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0
9	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	4

Comparing  
Computed NN  
results  
To  
Labeled results  
(target)

Third Dense Layer Output (batch)

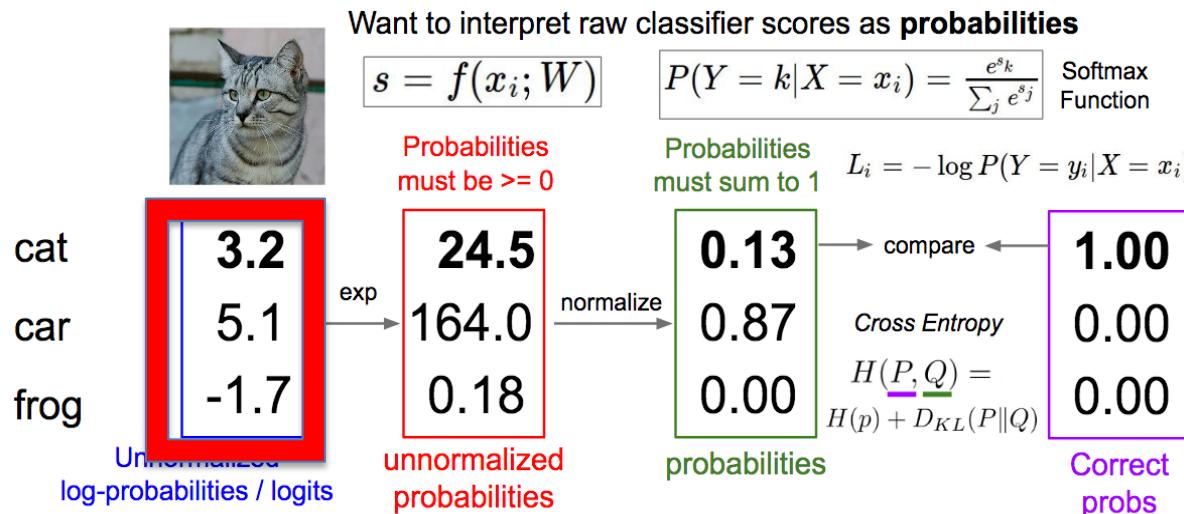
Predictions for each Image in Batch	0	1	2	3	4	5	6	7	8	9	Output Class
0	0.1	0.0	0.0	0.7	0.0	0.2	0.0	0.0	0.0	0.0	3
1	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1
2	0.0	0.0	0.0	0.0	0.0	0.7	0.0	0.0	0.0	0.2	2
3	0.0	0.9	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.1	0
4	0.0	0.0	0.0	0.0	0.1	0.0	0.0	0.0	0.0	0.9	4
5	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1
6	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0
7	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	3
8	0.0	0.9	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0
9	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	4

# Categorical (Softmax) Cross Entropy Loss (Statistical Learning)



$t_i$  and  $s_i$  are the groundtruth (label) and the computed score for each class  $i$  in  $C$ .

## Softmax Classifier (Multinomial Logistic Regression)



$$e^{3.2} = 24.5$$

$$e^{5.1} = 164.0$$

$$e^{-1.7} = 0.18$$

---

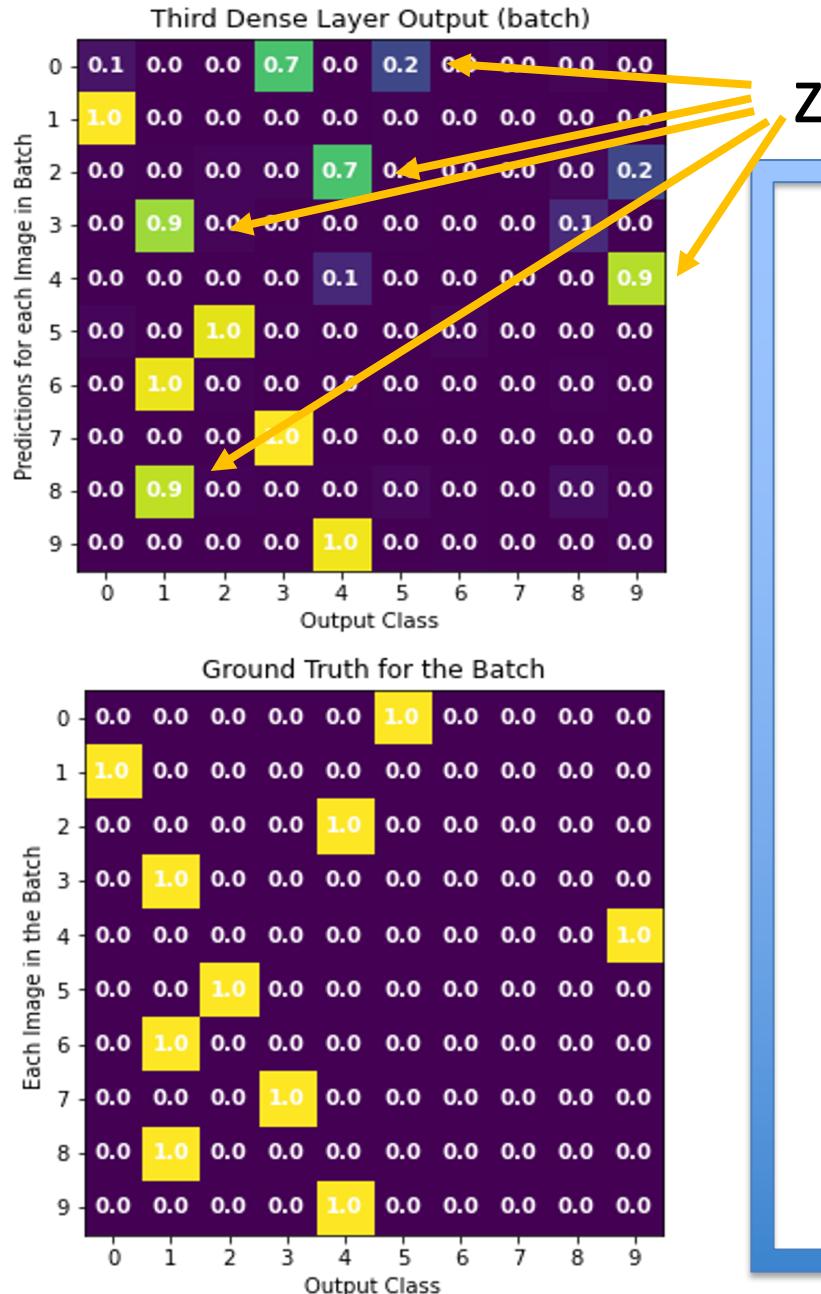
$$188.68$$

Software =  $24.5 / 188.68 = 0.13$

Software =  $164 / 188.68 = 0.87$

Software =  $-1.7 / 188.68 = 0.00$

# Evaluating the Error



# Categorical Cross Entropy Loss

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

$$L(i) = -\log 0.2 = 1.61$$

Image 0 = 5

$$-\log 1.0 = 0$$

Image 1 = 0

$$-\log 0.7 = 0.356$$

Image 2 = 4

$$-\log 0.9 = 0.105$$

Image 3 = 1

$$-\log 0.9 = 0.105$$

Image 4 = 9

$$-\log 1.0 = 0$$

Image 5 = 2

$$-\log 1.0 = 0$$

Image 6 = 1

$$-\log 1.0 = 0$$

Image 7 = 3

$$-\log 0.9 = 0.105$$

Image 8 = 2

$$-\log 1.0 = 0$$

Image 9 = 4

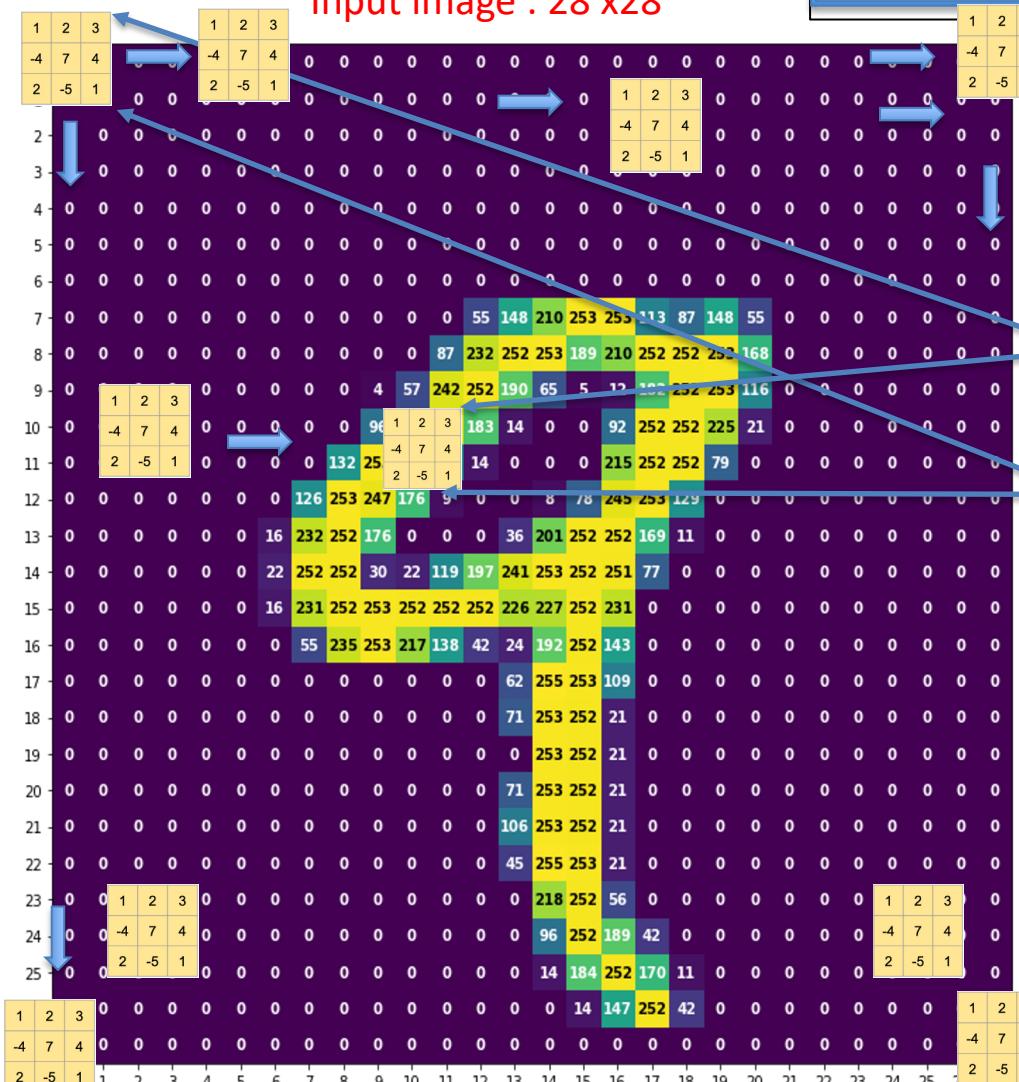
# MNIST Example (28x28 pixels image)

Handwritten digits in the MNIST dataset are 28x28 pixel grayscale images, use the original image (28x28 matrix).

Image : input

training digits and their labels
5 9 0 1 5 0 4 3 4 5 7 5 0 6 2 3 0 0 0 8 8 0 6 9
validation digits and their labels
7 2 1 0 4 1 4 9 5 9 0 6 9 0 1 5 9 7 3 4 9 6 4 5

Input image : 28 x28



Labels : target

12 filters, each one acts on an image until it is fully covered.

3 x 3 filter (kernel)

1	2	3
-4	7	4
2	-5	1

1	2	3
-4	7	4
2	-5	1

.....  
12 3x3 filters  
Same padding

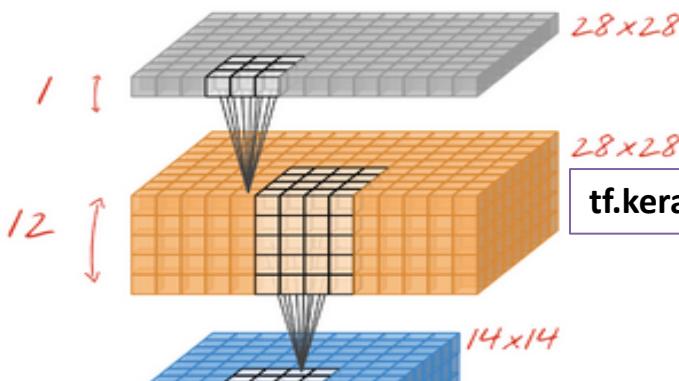
1	2	3
-4	7	4
2	-5	1



12 Output features : each 28 x28 matrix

# Convolution NN : MNIST 28x28 Pixels

[FW, FH, C, H] = [ filter size ? X ?, input channel (Depth), output channel (no. of filter) ]



`tf.keras.layers.Reshape(input_shape=(28*28,), target_shape=(28, 28, 1))`

*Convolutional 3x3 filters=12  
W<sub>1</sub>[3, 3, 1, 12]*

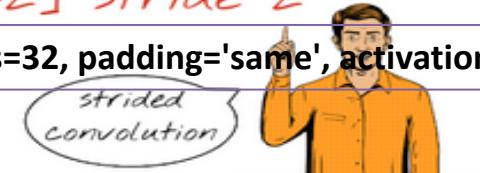
`tf.keras.layers.Conv2D(kernel_size=3, filters=12, padding='same', activation='relu')`

*Convolutional 6x6 filters=24  
W<sub>2</sub>[6, 6, 12, 24] stride 2*

`tf.keras.layers.Conv2D(kernel_size=6, filters=24, padding='same', activation='relu', strides=2)`

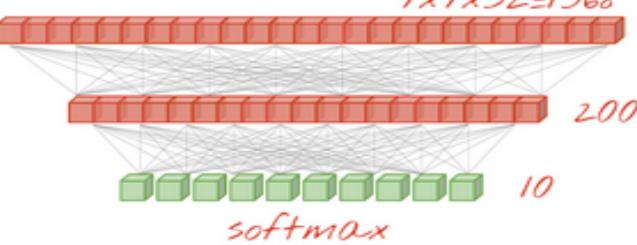
*Convolutional 6x6 filters=32  
W<sub>3</sub>[6, 6, 24, 32] stride 2*

`tf.keras.layers.Conv2D(kernel_size=6, filters=32, padding='same', activation='relu', strides=2)`



`tf.keras.layers.Flatten()`

<https://codelabs.developers.google.com/codelabs/cloud-tensorflow-mnist/#0>



*Dense layer*

*W<sub>4</sub>[1568, 200]*

`tf.keras.layers.Dense(200, activation='relu')`

*Softmax dense layer*

*W<sub>5</sub>[200, 10]*

`tf.keras.layers.Dense(10, activation='softmax')`

# CNN : example – MNIST

There are three main types of layers in a CNN: Convolutional Layer, Pooling Layer, and Fully-Connected Layer.

- ✓ INPUT [28x28x1] will hold the raw pixel values of the image, in this case an image of width 28, height 28 in grey scale (0-255)
- ✓ Apply twelve (12) 3x3 filters to one layer (channel) - W[3,3,Channel=1,number of filters=12]
- ✓ CONV layer will compute the output of neurons that are connected to local regions in the input, each computing a dot product between their weights and a small region they are connected to in the input volume. This may result in volume of [28x28x12] when we use 12 filters. Volume size = [ 28 width x height = 28 , 12 channels (depth) ]
- ✓ RELU layer will apply an elementwise activation function, such as the  $\max(0,x)$
- ✓ Apply 24 6x6 filters to 12 layers (channel) with stride 2 - W[6,6,12,24] stride 2
- ✓ This may result in volume of [14x14x24] when we use 24 filters, stride 2.
- ✓ RELU layer will apply an elementwise activation function, such as the  $\max(0,x)$
- ✓ Apply 32 6x6 filters to 24 layers (channel) with stride 2 - W[6,6,24,32] stride 2
- ✓ This may result in volume of [7x7x32] when we use 32 filters, stride 2.
- ✓ RELU layer will apply an elementwise activation function, such as the  $\max(0,x)$
- ✓ Flatten the volume (tensor) to be a vector of  $7 \times 6 \times 32 = 1568$  elements
- ✓ Use 200 neurons with the input vector, Fully connected layer, W [1558, 200], output=1 x 200
- ✓ FC (i.e. fully-connected) layer will compute the class scores, resulting in volume of size [1x1x10], where each of the 10 numbers correspond to a class score.
- ✓ In this way, ConvNets transform the original image layer by layer from the original pixel values to the final class scores.

```
import tensorflow as tf
print(tf.__version__)

import datetime, os

%load_ext tensorboard

import numpy as np
from tensorflow import keras
from tensorflow.keras import layers

# Model / data parameters
num_classes = 10
input_shape = (28, 28, 1)

# the data, split between train and test sets
(x_train, y_train), (x_test, y_test) = keras.datasets.mnist.load_data()

# Scale images to the [0, 1] range
x_train = x_train.astype("float32") / 255
x_test = x_test.astype("float32") / 255
print("x_train shape:", x_train.shape)
print(x_train.shape[0], "train samples")
print(x_test.shape[0], "test samples")

# Make sure images have shape (28, 28, 1)
x_train = x_train.reshape([-1,28,28,1])
x_test = x_test.reshape([-1,28,28,1])
print("x_train shape:", x_train.shape)
print(x_train.shape[0], "train samples")
print(x_test.shape[0], "test samples")

# convert class vectors to binary class matrices
y_train = keras.utils.to_categorical(y_train, num_classes)
y_test = keras.utils.to_categorical(y_test, num_classes)
```

[https://github.com/keras-team/keras-io/blob/master/examples/vision/mnist\\_convnet.py](https://github.com/keras-team/keras-io/blob/master/examples/vision/mnist_convnet.py)

[https://github.com/keras-team/keras-io/blob/master/examples/vision/mnist\\_convnet.py](https://github.com/keras-team/keras-io/blob/master/examples/vision/mnist_convnet.py)

```
## Build the model
model = keras.Sequential(
[
    keras.Input(shape=input_shape),
    layers.Conv2D(32, kernel_size=(3, 3), activation="relu"),
    layers.MaxPooling2D(pool_size=(2, 2)),
    layers.Conv2D(64, kernel_size=(3, 3), activation="relu"),
    layers.MaxPooling2D(pool_size=(2, 2)),
    layers.Flatten(),
    layers.Dropout(0.5),
    layers.Dense(num_classes, activation="softmax"),
]
)

model.summary()

## Train the model
batch_size = 128
epochs = 5

model.compile(loss="categorical_crossentropy", optimizer="adam", metrics=["accuracy"])

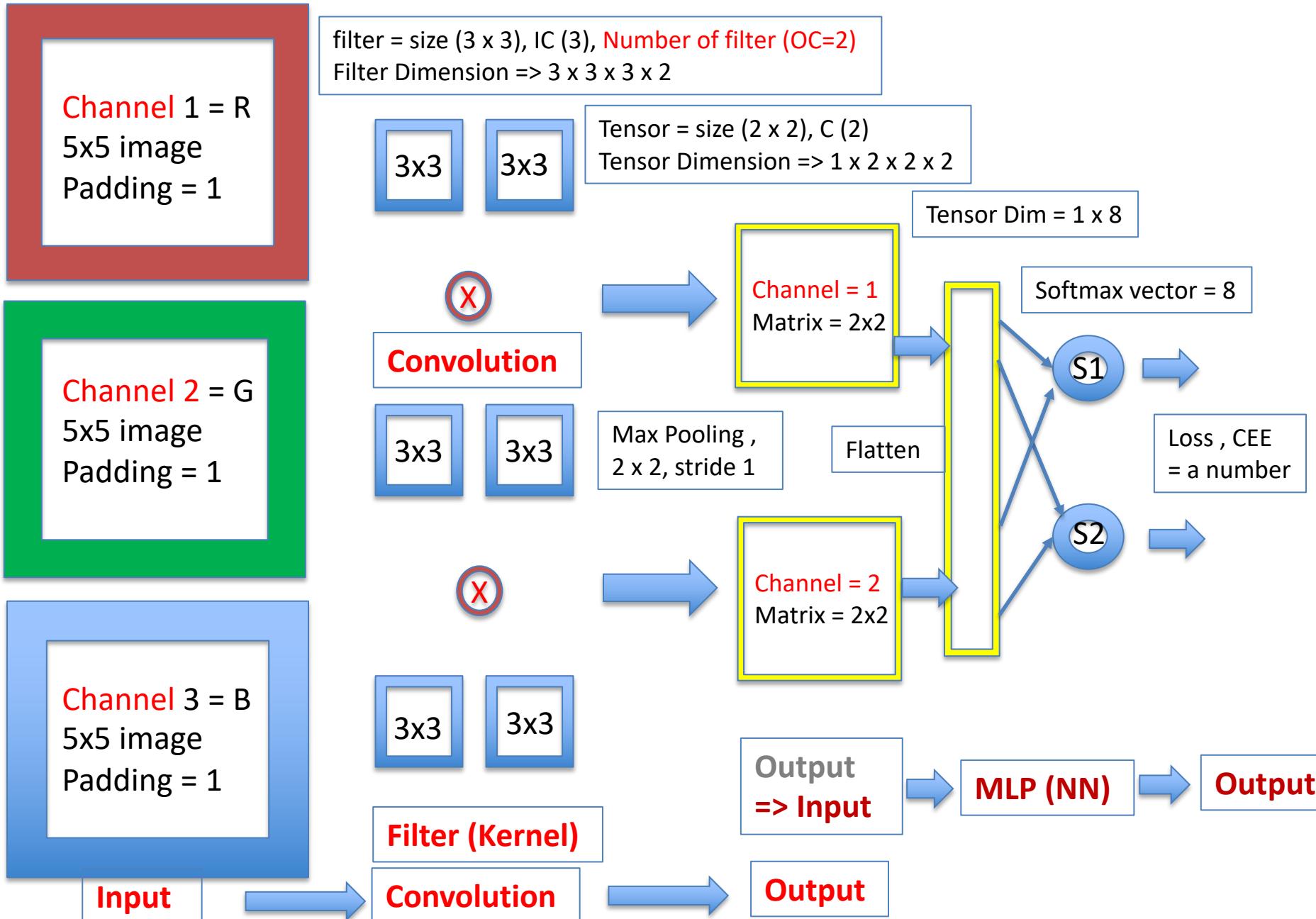
logdir = os.path.join("logs", datetime.datetime.now().strftime("%Y%m%d-%H%M%S"))
tensorboard_callback = tf.keras.callbacks.TensorBoard(logdir, histogram_freq=1)

model.fit(x_train, y_train, batch_size=batch_size, epochs=epochs, validation_split=0.1,
callbacks=[tensorboard_callback])

## Evaluate the trained model
%tensorboard --logdir logs
```

Input = size (5 x 5) Channel (IC = 3)+ padding (1), Stride 2  
Input dimension = 7 x 7 x 3 or 1 x 7 x 7 x 3

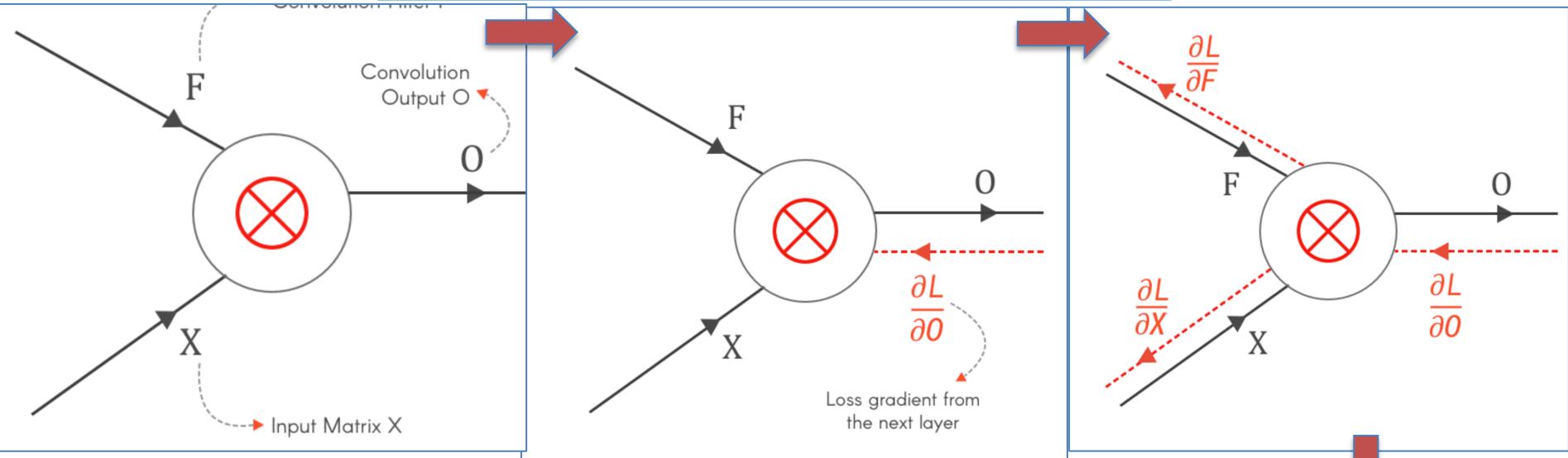
# CNN Model Summary



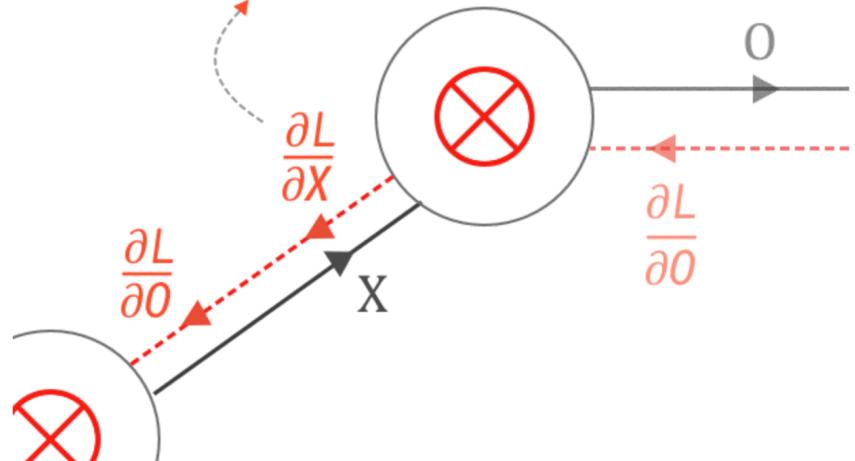
# CNN : Backward Path Calculation

## A Simple Example

# Forward path : backward path

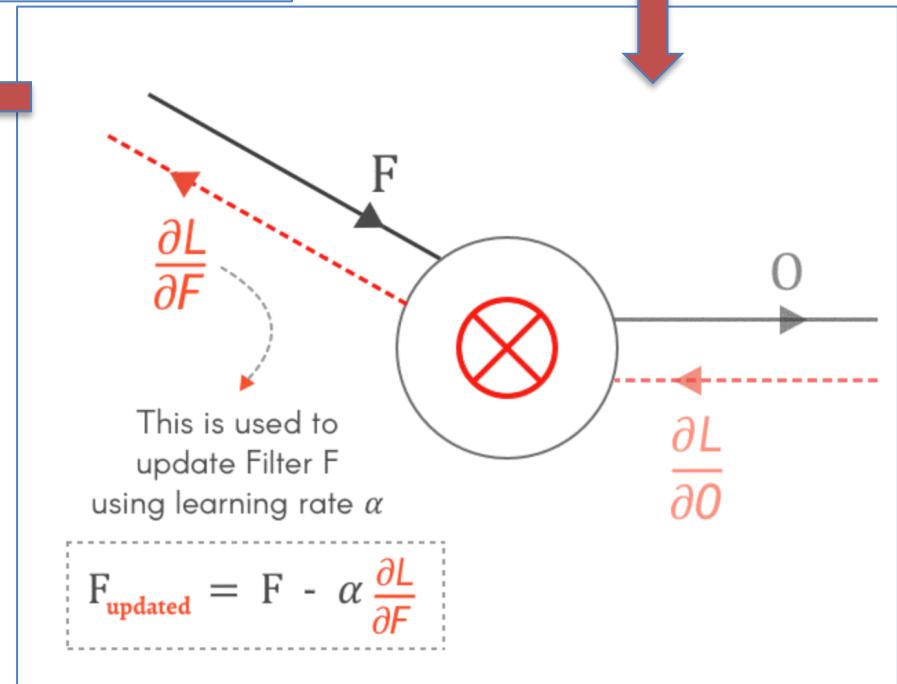


Since  $X$  is the output of the previous layer,  
 $\frac{\partial L}{\partial X}$  becomes the loss gradient for the previous layer



This is used to update Filter  $F$  using learning rate  $\alpha$

$$F_{\text{updated}} = F - \alpha \frac{\partial L}{\partial F}$$



# Example exercise

Compute the updated values of filter weights after backpropagation of the following problem -- 3x3 input matrix, 2x2 filter, stride 1, no padding, flatten it and use it as input connecting to a two-neuron layer, then pass the results to a softmax binary classification of two neurons.

3a) (0.5%) What is the values of the output at the end of the forward path

3b) (1.5%) What are the updated values of the weight and bias after backpropagation

1	0	1
0	1	0
1	0	1

w1	w2
w3	w4

$$b_1 = 0.1$$

$$\begin{aligned}w_1 &= 0.9 \\w_2 &= 0.1 \\w_3 &= 0.1 \\w_4 &= 0.9\end{aligned}$$



Apply ReLU activation function after the convolution step

Flatten the value after applying the ReLU activation function

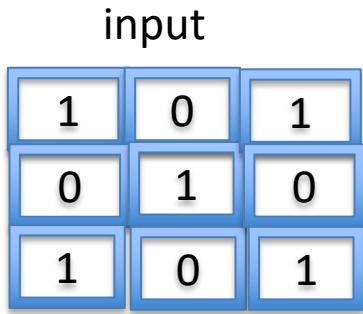
Assuming the weight of each link to this layer is  $w = 0.1$  to  $0.8$ ,  $b_2 = 0.2$

Fully connecting to a layer of two neurons

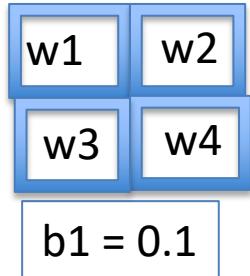
Assuming the weight of each link is  $w = 0.1$  to  $0.4$ ,  $b_3 = 0.3$

Fully Connected layer of 2 softmax neurons

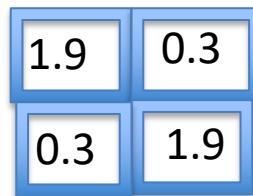
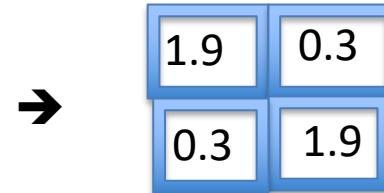
Labeled output of two classes are 1.0 and 0.00



Filter, stride 1



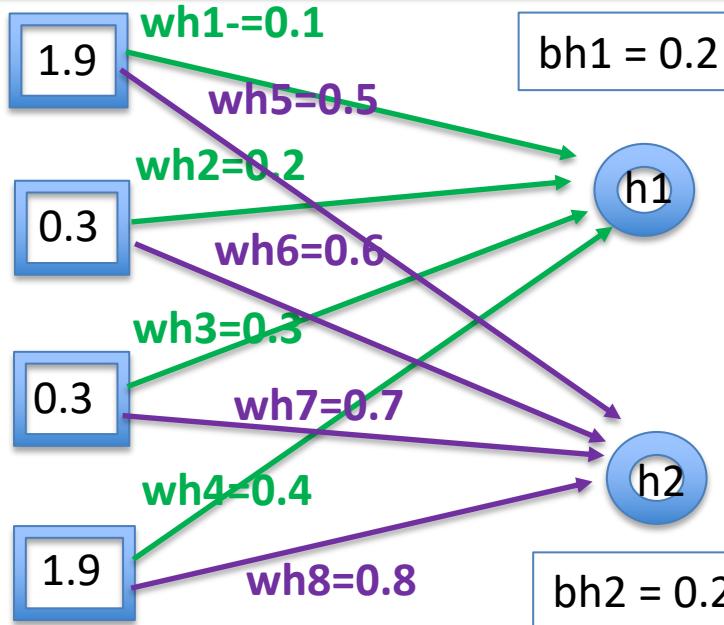
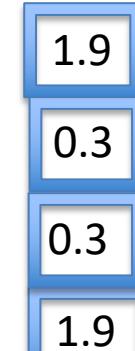
w1=0.9  
w2=0.1  
w3=0.1  
w4=0.9



ReLU

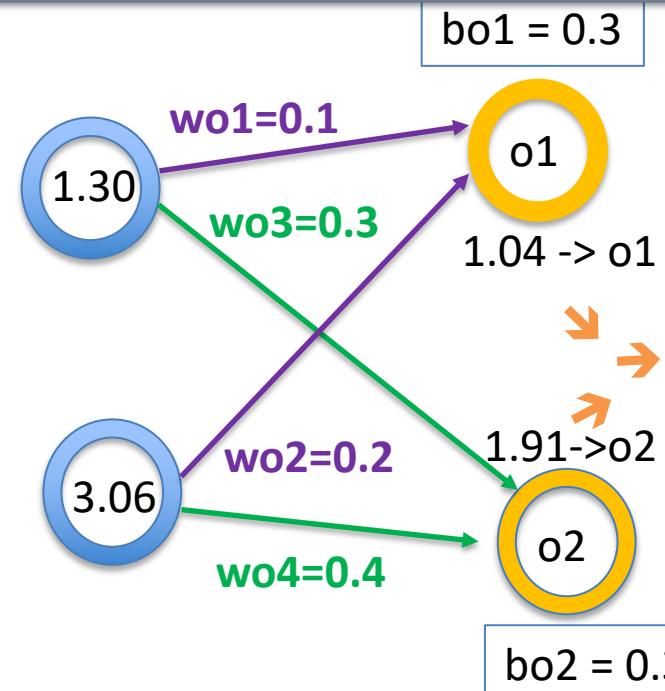


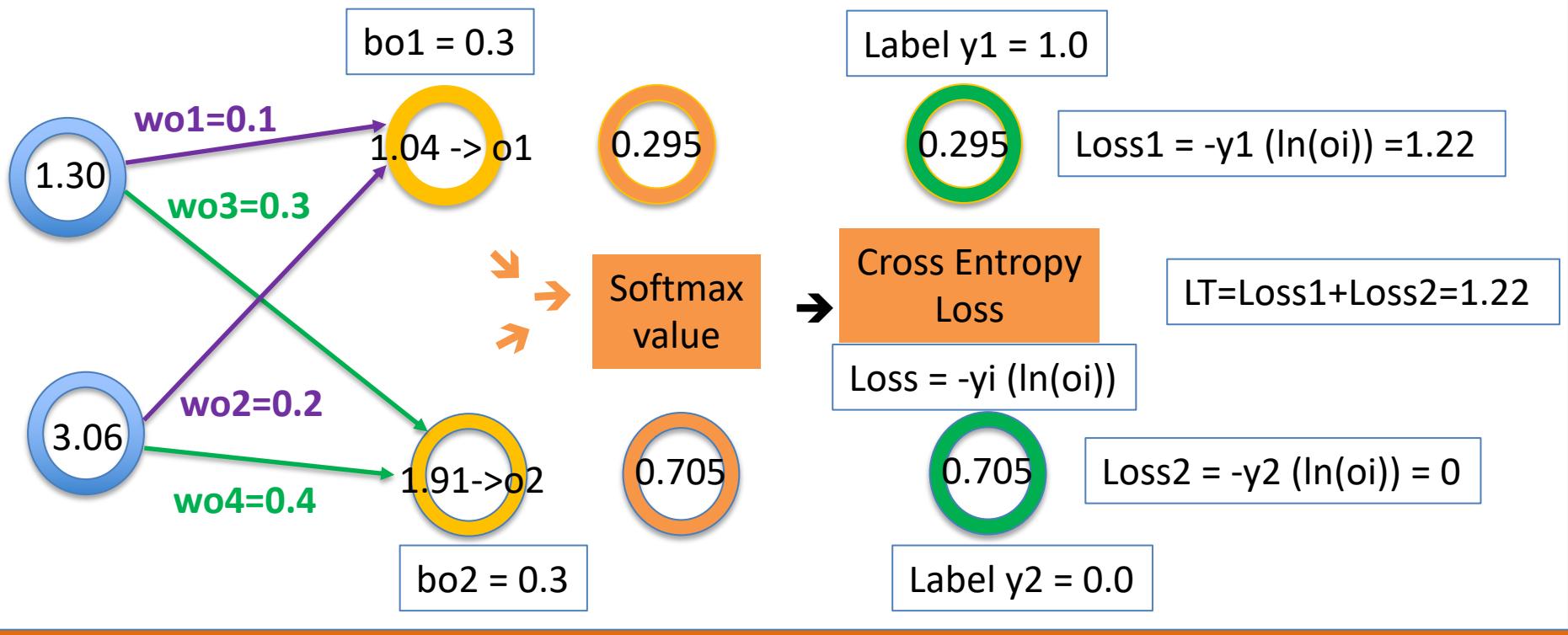
flatten



h1

h2





CNN : Backward Path – jumps over Cross Entropy and Softmax all together

Learning rate  
=  $n = 0.5$

$$\frac{\partial E}{\partial Z_k} = O_k - y_k$$

$$Z = Z(w, b)$$

$$S = S(z)$$

weights  
bias

$$w+ = w - n * (dE/dw)$$

$$dE/dw = (dE/dS) * (dS/dw)$$

Softmax  
value

$dL/dS$

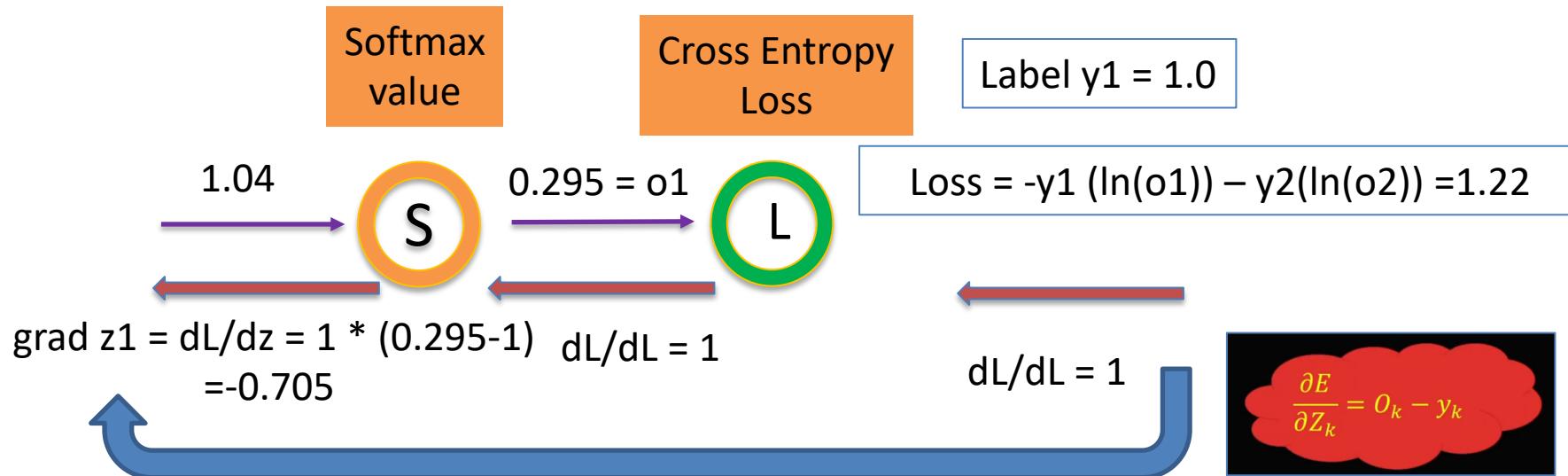
Cross Entropy  
Loss ( $L$ )

$dL/dL = 1$

# Flow Back from the Total Loss to Softmax to Z (S1)

$$\text{grad down} = (\text{grad up}) * (\text{local grad})$$

$$\text{grad } L = dL_i = O_i - Y_i$$



Jump from cross entropy and softmax (E , S) to (weight and bias) Z (o1), Y1 =1

$$E=L ; \frac{dE}{dz} = \frac{dL}{dz} = (\frac{dL}{ds} * \frac{ds}{dz})$$

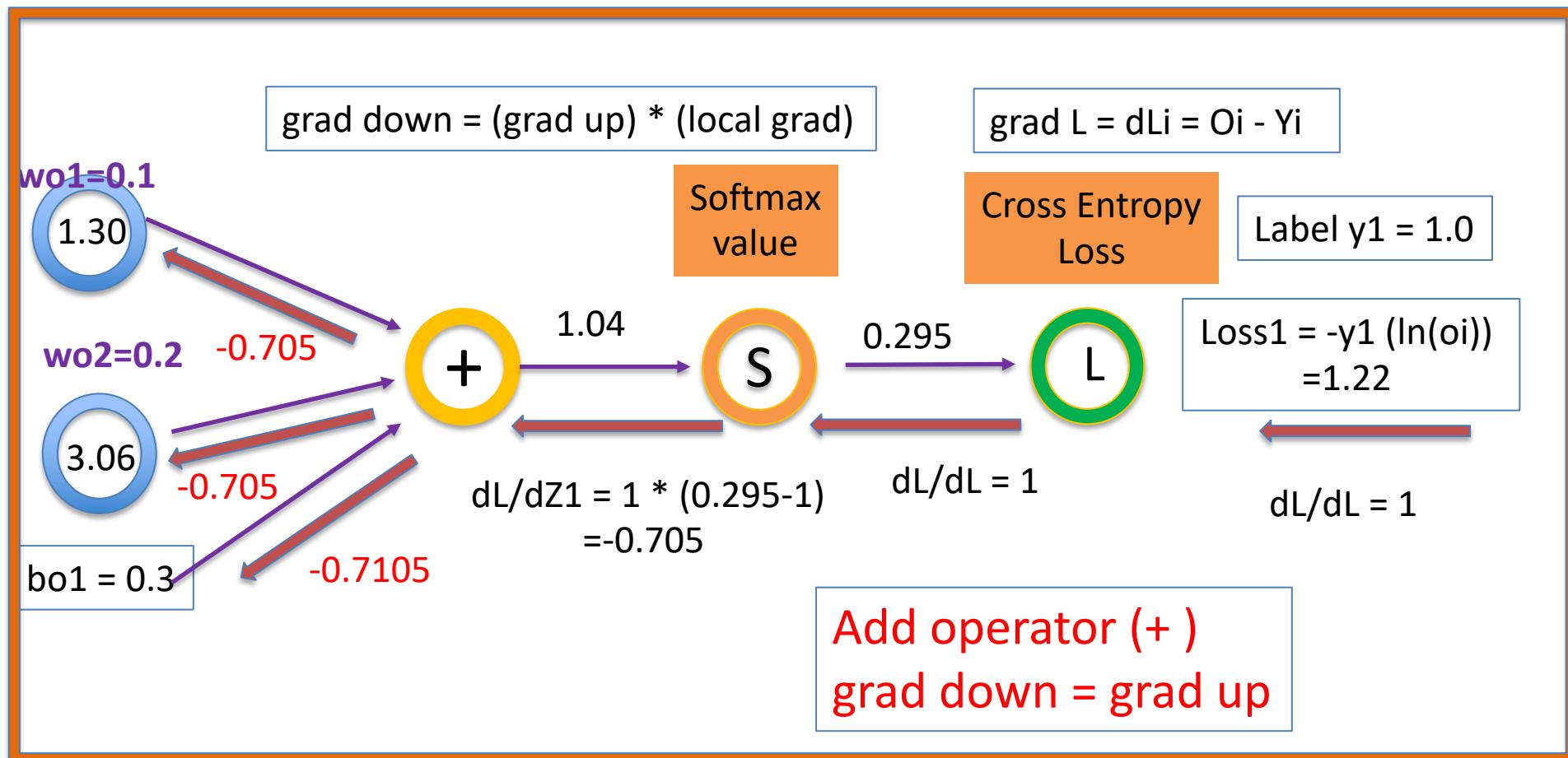
$$\frac{dE}{dZ} = \text{local grad} = O_1 - Y_1$$

$$\begin{aligned}\text{grad } Z_1 &= \text{grad down} = (\text{grad up} * \text{local grad}) \\ &= 1 * (0.295 - 1) = -0.705\end{aligned}$$

$$\frac{\partial E}{\partial Z_k} = O_k - y_k$$

# Flow Back from Cross Entropy and Softmax to W and b

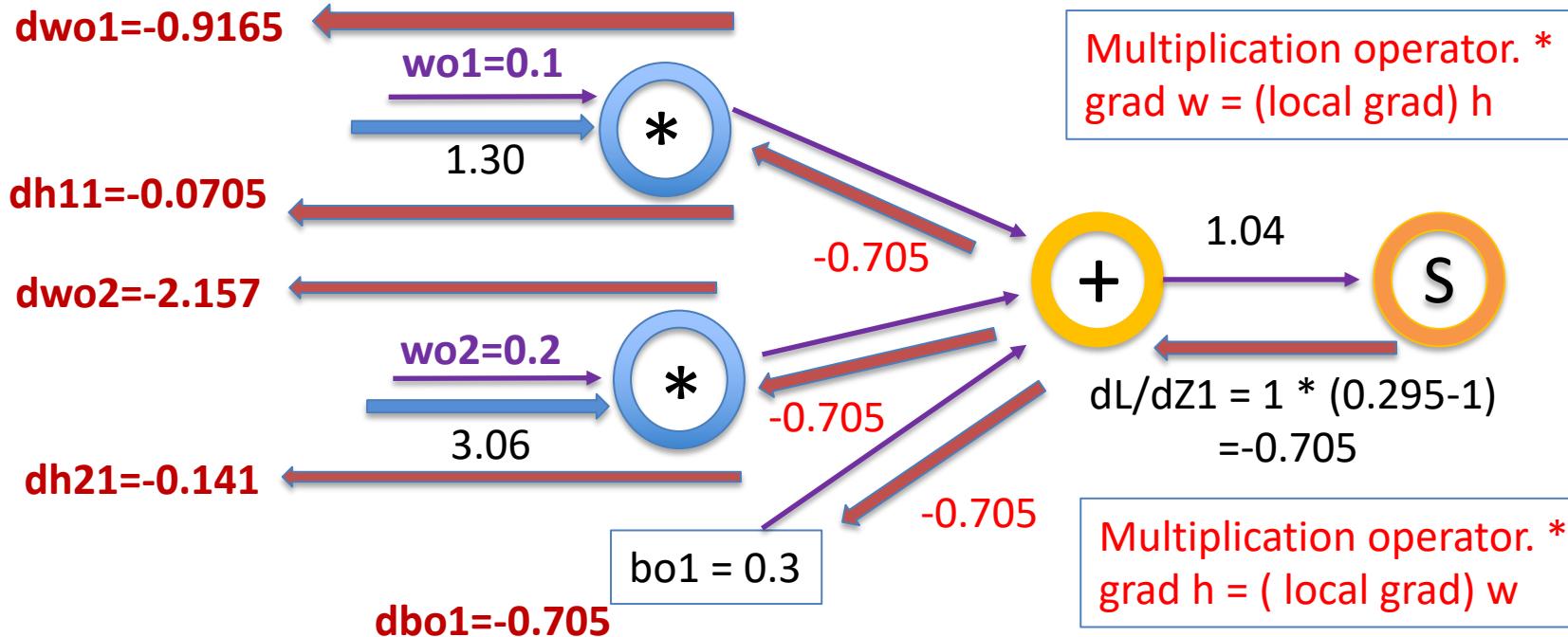
$L \rightarrow L1 \rightarrow S1(o_1) \rightarrow w_{o1}, w_{o2}, b_{o1}$



Multiplication operator. (\*)  
 $\text{grad } w = (\text{local grad}) h$

Multiplication operator. (\*)  
 $\text{grad } h = (\text{local grad}) w$

## Update values of wo1, wo2, bo1



In here, i = 1 ( o1),  
j = 1,2: h1, h2

$$\begin{aligned} d(L/dw_{o1}) &= dL/d(w_{11}) \\ &= -0.705 \cdot 1.30 = -0.9165 \end{aligned}$$

$$\begin{aligned} d(L/dw_{o2}) &= dL/d(w_{12}) \\ &= -0.705 \cdot 3.06 = -2.1573 \end{aligned}$$

$$\begin{aligned} d(L/db_{o1}) &= dL/d(b_1) \\ &= dL/dS = O_i - Y_i = -0.705 \end{aligned}$$

$$wo_1 += 0.1 - (0.5) * (-0.9165) = 0.558$$

$$wo_2 += 0.2 - (0.5) * (-2.1573) = 1.278$$

$$bo_1 += 0.3 - (0.5) * (-0.705) = 0.652$$

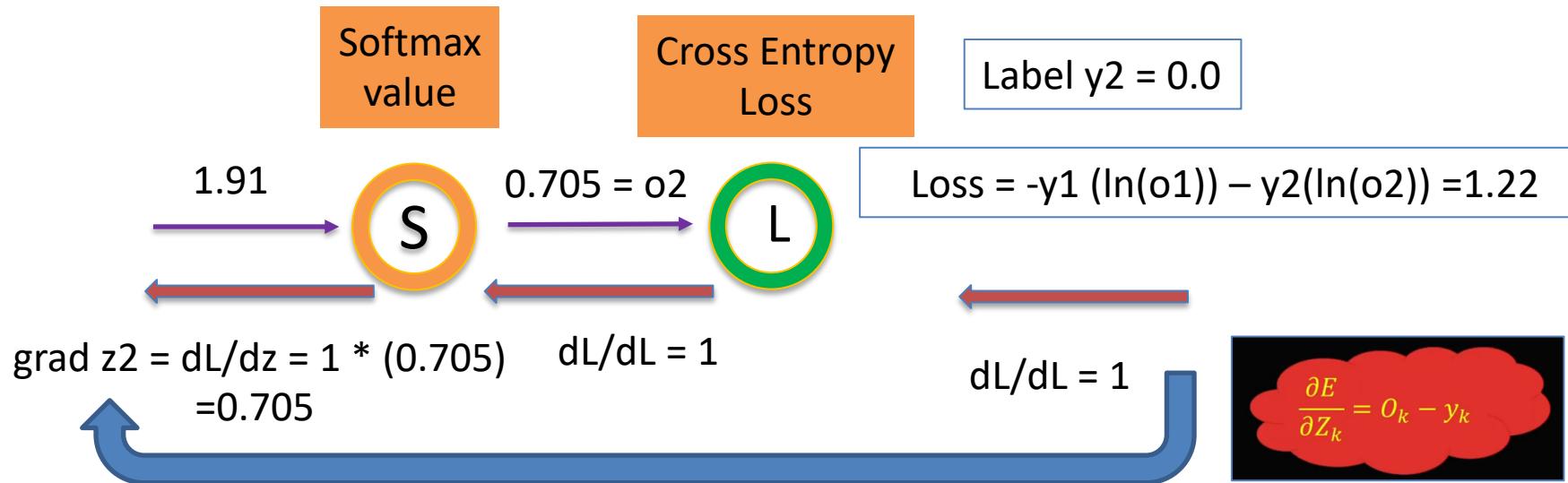
$$dh_{11} = (-0.0705)$$

$$dh_{21} = (-0.141)$$

# Flow Back from the CEE Loss to Softmax to Z : S2

grad down = (grad up) \* (local grad)

grad L = dLi = Oi - Yi



Jump from cross entropy and softmax (E, S) to weight and bias Z ( $o_2$ ),  $Y_2=0$

$$E=L ; dE/dz = dL/dz = (dL/ds * ds/dz)$$

$$dL/dz = \text{local grad} = O_2 - Y_2 = O_2$$

$$\begin{aligned} \text{grad } Z_2 &= \text{grad down} = (\text{grad up} * \text{local grad}) \\ &= 1 * (0.705) = 0.705 \end{aligned}$$

$$\frac{\partial E}{\partial Z_k} = o_k - y_k$$

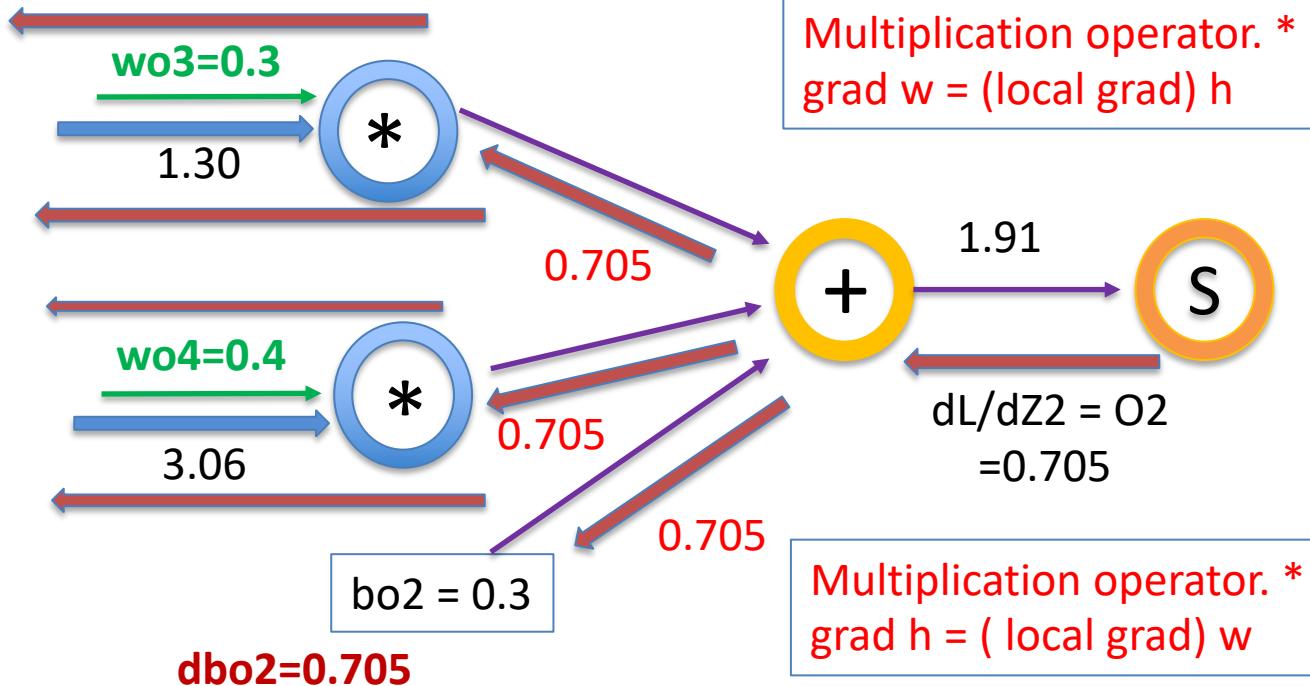
## Update values of wo3, wo4, bo2

$$dwo3=0.9165$$

$$dh12=0.211$$

$$dwo4=2.157$$

$$dh22=0.282$$



$$\begin{aligned} d(L/dwo3) &= dL/d(w21) \\ &= 0.705 \cdot 1.30 = 0.9165 \end{aligned}$$

$$\begin{aligned} d(L/dwo4) &= dL/d(w22) \\ &= 0.705 \cdot 3.06 = 2.157 \end{aligned}$$

$$\begin{aligned} d(L/dbo1) &= dL/d(b1) \\ &= dL/dS = Oi-Yi=0.705 \end{aligned}$$

$$wo3+= 0.3 - (0.5) \cdot (0.9165) = -0.158$$

$$wo4+= 0.4 - (0.5) \cdot (2.157) = -0.678$$

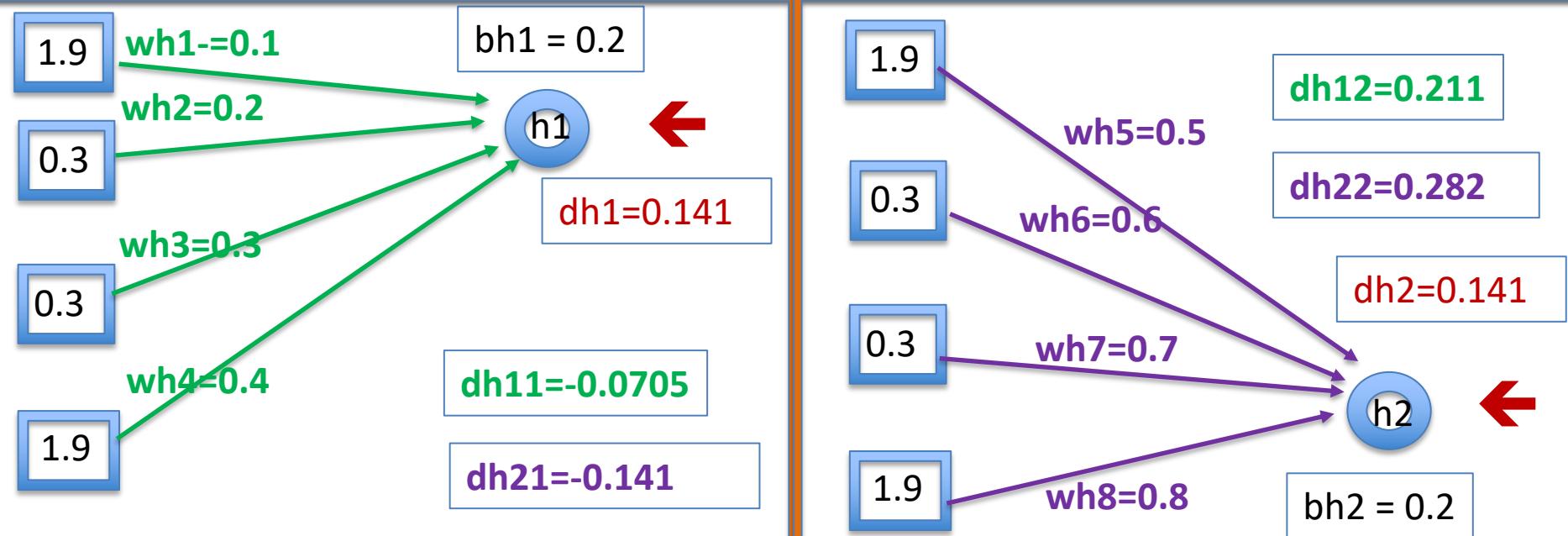
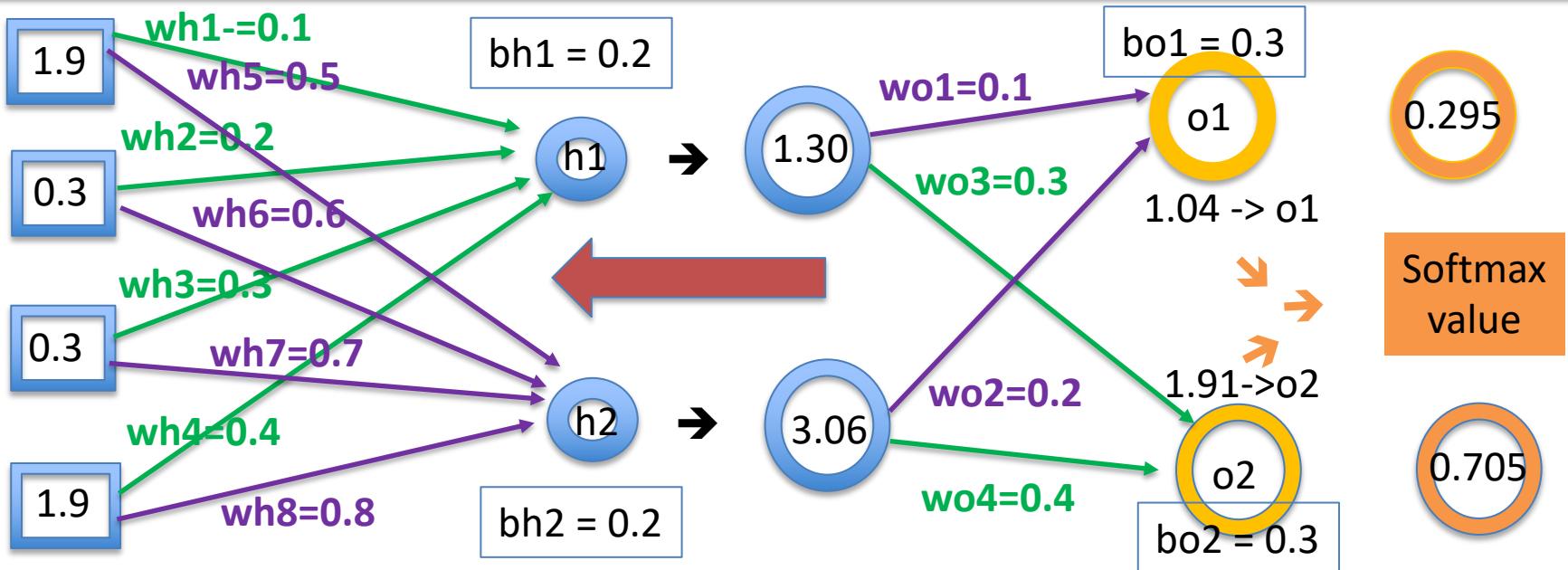
$$bo2+= 0.3 - (0.5) \cdot (0.705) = -0.352$$

$$dh12= (0.211)$$

$$dh22= (0.282)$$

In here , i =2 ( o2)  
j = 1,2: h1, h2

# Flow back from o1 to the hidden NN layer (h) , MLP calculation



# Gradients of C matrix (add grads of (h1 + h2))

$$dC1 = 0.0141 + 0.0705 = 0.0846$$

$$dC2 = 0.0282 + 0.0846 = 0.1128$$

$$dC3 = 0.0423 + 0.0987 = 0.141$$

$$dC4 = 0.0564 + 0.1128 = 0.1692$$

1.9

0.3

0.3

1.9

dC1 = 0.0846

dC2 = 0.1128

dC3 = 0.141

dC4 = 0.1692

$$dC1h1 = 0.1 \times 0.141 = 0.0141$$

$$dC2h1 = 0.2 \times 0.141 = 0.0282$$

$$dC3h1 = 0.3 \times 0.141 = 0.0423$$

$$dC4h1 = 0.4 \times 0.141 = 0.0564$$

$$dC1h2 = 0.5 \times 0.141 = 0.0705$$

$$dC2h2 = 0.6 \times 0.141 = 0.0846$$

$$dC3h2 = 0.7 \times 0.141 = 0.0987$$

$$dC4h2 = 0.8 \times 0.141 = 0.1128$$

1.9

0.3

0.3

1.9

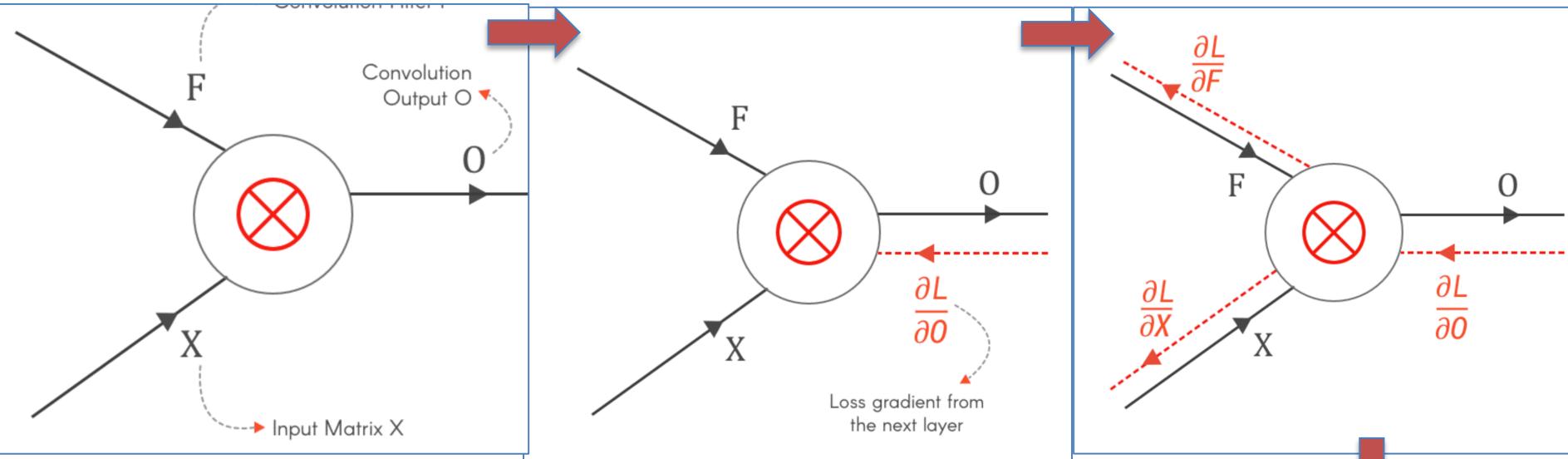
1.9

0.3

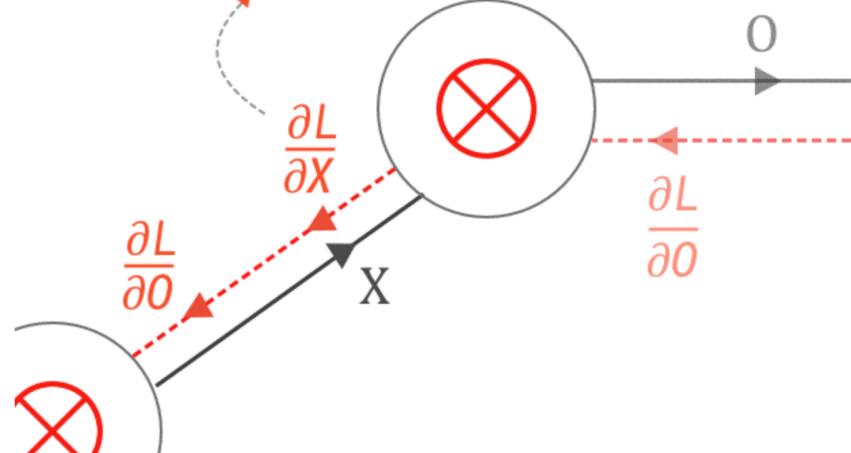
0.3

1.9

# Convolution : Forward path : backward path



Since  $X$  is the output of the previous layer,  
 $\frac{\partial L}{\partial X}$  becomes the loss gradient for the previous layer



This is used to update Filter  $F$  using learning rate  $\alpha$

$$F_{\text{updated}} = F - \alpha \frac{\partial L}{\partial F}$$

Local Gradients —— A

$$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22}$$

Finding derivatives with respect to  $F_{11}$ ,  $F_{12}$ ,  $F_{21}$  and  $F_{22}$

$$\frac{\partial O_{11}}{\partial F_{11}} = X_{11} \quad \frac{\partial O_{11}}{\partial F_{12}} = X_{12} \quad \frac{\partial O_{11}}{\partial F_{21}} = X_{21} \quad \frac{\partial O_{11}}{\partial F_{22}} = X_{22}$$

Similarly, we can find the local gradients for  $O_{12}$ ,  $O_{21}$  and  $O_{22}$

$$\frac{\partial L}{\partial F_{11}} = \frac{\partial L}{\partial O_{11}} * \frac{\partial O_{11}}{\partial F_{11}} + \frac{\partial L}{\partial O_{12}} * \frac{\partial O_{12}}{\partial F_{11}} + \frac{\partial L}{\partial O_{21}} * \frac{\partial O_{21}}{\partial F_{11}} + \frac{\partial L}{\partial O_{22}} * \frac{\partial O_{22}}{\partial F_{11}}$$

$$\frac{\partial L}{\partial F_{12}} = \frac{\partial L}{\partial O_{11}} * \frac{\partial O_{11}}{\partial F_{12}} + \frac{\partial L}{\partial O_{12}} * \frac{\partial O_{12}}{\partial F_{12}} + \frac{\partial L}{\partial O_{21}} * \frac{\partial O_{21}}{\partial F_{12}} + \frac{\partial L}{\partial O_{22}} * \frac{\partial O_{22}}{\partial F_{12}}$$

$$\frac{\partial L}{\partial F_{21}} = \frac{\partial L}{\partial O_{11}} * \frac{\partial O_{11}}{\partial F_{21}} + \frac{\partial L}{\partial O_{12}} * \frac{\partial O_{12}}{\partial F_{21}} + \frac{\partial L}{\partial O_{21}} * \frac{\partial O_{21}}{\partial F_{21}} + \frac{\partial L}{\partial O_{22}} * \frac{\partial O_{22}}{\partial F_{21}}$$

$$\frac{\partial L}{\partial F_{22}} = \frac{\partial L}{\partial O_{11}} * \frac{\partial O_{11}}{\partial F_{22}} + \frac{\partial L}{\partial O_{12}} * \frac{\partial O_{12}}{\partial F_{22}} + \frac{\partial L}{\partial O_{21}} * \frac{\partial O_{21}}{\partial F_{22}} + \frac{\partial L}{\partial O_{22}} * \frac{\partial O_{22}}{\partial F_{22}}$$

a convolution operation between input X and loss gradient  $\frac{\partial L}{\partial O}$

$$\frac{\partial L}{\partial F_{11}} = \frac{\partial L}{\partial O_{11}} * X_{11} + \frac{\partial L}{\partial O_{12}} * X_{12} + \frac{\partial L}{\partial O_{21}} * X_{21} + \frac{\partial L}{\partial O_{22}} * X_{22}$$

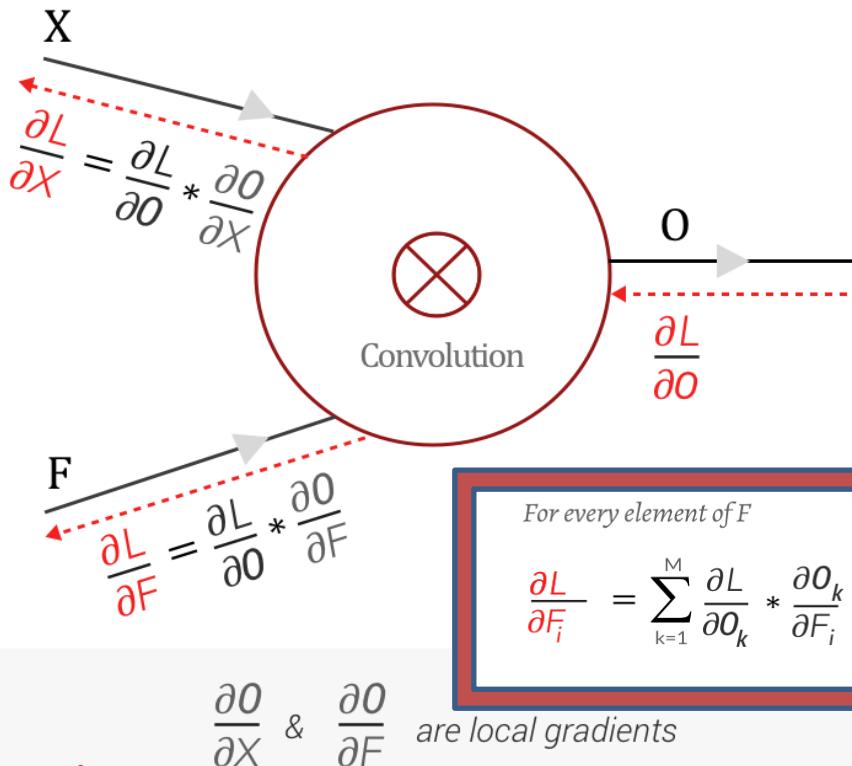
$$\frac{\partial L}{\partial F_{12}} = \frac{\partial L}{\partial O_{11}} * X_{12} + \frac{\partial L}{\partial O_{12}} * X_{13} + \frac{\partial L}{\partial O_{21}} * X_{22} + \frac{\partial L}{\partial O_{22}} * X_{23}$$

$$\frac{\partial L}{\partial F_{21}} = \frac{\partial L}{\partial O_{11}} * X_{21} + \frac{\partial L}{\partial O_{12}} * X_{22} + \frac{\partial L}{\partial O_{21}} * X_{31} + \frac{\partial L}{\partial O_{22}} * X_{32}$$

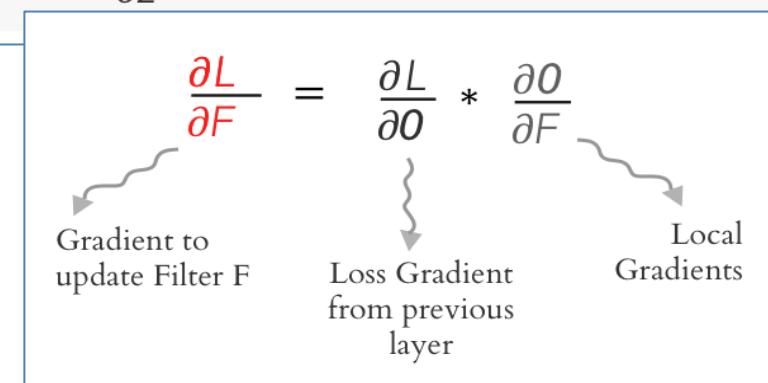
$$\frac{\partial L}{\partial F_{22}} = \frac{\partial L}{\partial O_{11}} * X_{22} + \frac{\partial L}{\partial O_{12}} * X_{23} + \frac{\partial L}{\partial O_{21}} * X_{32} + \frac{\partial L}{\partial O_{22}} * X_{33}$$

2x2 filter

<https://medium.com/@pavisj/convolutions-and-backpropagations-46026a8f5d2c>



$\frac{\partial L}{\partial Z}$  is the loss from the previous layer which has to be backpropagated to other layers



$$\begin{bmatrix} \frac{\partial L}{\partial F_{11}} & \frac{\partial L}{\partial F_{12}} \\ \frac{\partial L}{\partial F_{21}} & \frac{\partial L}{\partial F_{22}} \end{bmatrix}$$

= Convolution

$$\begin{bmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{bmatrix}$$

$$\begin{bmatrix} \frac{\partial L}{\partial O_{11}} & \frac{\partial L}{\partial O_{12}} \\ \frac{\partial L}{\partial O_{21}} & \frac{\partial L}{\partial O_{22}} \end{bmatrix}$$

**$\frac{\partial L}{\partial F}$  is nothing but the convolution between Input X and Loss Gradient from the next layer  $\frac{\partial L}{\partial O}$**

where

$$\begin{bmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{bmatrix}$$

= Input X

$$\begin{bmatrix} \frac{\partial L}{\partial O_{11}} & \frac{\partial L}{\partial O_{12}} \\ \frac{\partial L}{\partial O_{21}} & \frac{\partial L}{\partial O_{22}} \end{bmatrix}$$

=  $\frac{\partial L}{\partial O}$  Loss gradient from previous layer

$$\frac{\partial L}{\partial X_{11}} = \frac{\partial L}{\partial O_{11}} * F_{11}$$

$$\frac{\partial L}{\partial X_{12}} = \frac{\partial L}{\partial O_{11}} * F_{12} + \frac{\partial L}{\partial O_{12}} * F_{11}$$

$$\frac{\partial L}{\partial X_{13}} = \frac{\partial L}{\partial O_{12}} * F_{12}$$

$$\frac{\partial L}{\partial X_{21}} = \frac{\partial L}{\partial O_{11}} * F_{21} + \frac{\partial L}{\partial O_{21}} * F_{11}$$

$$\frac{\partial L}{\partial X_{22}} = \frac{\partial L}{\partial O_{11}} * F_{22} + \frac{\partial L}{\partial O_{12}} * F_{21} + \frac{\partial L}{\partial O_{21}} * F_{12} + \frac{\partial L}{\partial O_{22}} * F_{11}$$

Local Gradients:

(B)

2x2 filter

$$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22}$$

Differentiating with respect to  $X_{11}, X_{12}, X_{21}$  and  $X_{22}$

$$\frac{\partial O_{11}}{\partial X_{11}} = F_{11} \quad \frac{\partial O_{11}}{\partial X_{12}} = F_{12} \quad \frac{\partial O_{11}}{\partial X_{21}} = F_{21} \quad \frac{\partial O_{11}}{\partial X_{22}} = F_{22}$$

Similarly, we can find local gradients for  $O_{12}, O_{21}$  and  $O_{22}$

For every element of  $X_i$

$$\frac{\partial L}{\partial X_i} = \sum_{k=1}^M \frac{\partial L}{\partial O_k} * \frac{\partial O_k}{\partial X_i}$$

$$\frac{\partial L}{\partial X_{23}} = \frac{\partial L}{\partial O_{12}} * F_{22} + \frac{\partial L}{\partial O_{22}} * F_{12}$$

$$\frac{\partial L}{\partial X_{31}} = \frac{\partial L}{\partial O_{21}} * F_{21}$$

$$\frac{\partial L}{\partial X_{32}} = \frac{\partial L}{\partial O_{21}} * F_{22} + \frac{\partial L}{\partial O_{22}} * F_{21}$$

$$\frac{\partial L}{\partial X_{33}} = \frac{\partial L}{\partial O_{22}} * F_{22}$$

## Backpropagation in a Convolutional Layer of a CNN

Finding the gradients:

$$\frac{\partial L}{\partial F} = \text{Convolution} \left( \text{Input } X, \text{ Loss gradient } \frac{\partial L}{\partial O} \right)$$

$$\frac{\partial L}{\partial X} = \text{Full Convolution} \left( 180^\circ \text{rotated Filter } F, \text{ Gradient } \frac{\partial L}{\partial O} \right)$$

Loss gradient =  $dL/dC$

$dC1=0.0846$

$dC2=0.1128$

$dC3=0.141$

$dC4=0.1692$

1	0	1
0	1	0
1	0	1



$dC1=0.0846$

$dC2=0.1128$

$dC3=0.141$

$dC4=0.1692$

$dF1=0.2538$

$dF2=0.2538$

$dF3=0.2538$

$dF4=0.238$

$$= dL/dF$$

$$dF1 = 1 \times dC1 + 0 \times dC2 + 0 \times dC3 + 1 \times dC4 = 0.5238$$

$$dF2 = 0 \times dC1 + 1 \times dC2 + 1 \times dC3 + 0 \times dC4 = 0.2538$$

$$dF3 = 0 \times dC1 + 1 \times dC2 + 1 \times dC3 + 0 \times dC4 = 0.2538$$

$$dF4 = 1 \times dC1 + 0 \times dC2 + 0 \times dC3 + 1 \times dC4 = 0.2538$$

$F1+=0.7731$

$F2+=-0.027$

$F3+=-0.027$

$F4+=0.7731$

=

0.9	0.1
0.1	0.9

$$- (0.5) *$$

$dF1=0.2538$

$dF2=0.2538$

$dF3=0.2538$

$dF4=0.2538$

# Update values of b1 – bias of convolution

Local Gradients → A

$$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22} + b$$

Finding derivatives with respect to  $F_{11}$ ,  $F_{12}$ ,  $F_{21}$  and  $F_{22}$

$$\frac{\partial O_{11}}{\partial F_{11}} = X_{11} \quad \frac{\partial O_{11}}{\partial F_{12}} = X_{12} \quad \frac{\partial O_{11}}{\partial F_{21}} = X_{21} \quad \frac{\partial O_{11}}{\partial F_{22}} = X_{22}$$

Similarly, we can find the local gradients for  $O_{12}$ ,  $O_{21}$  and  $O_{22}$

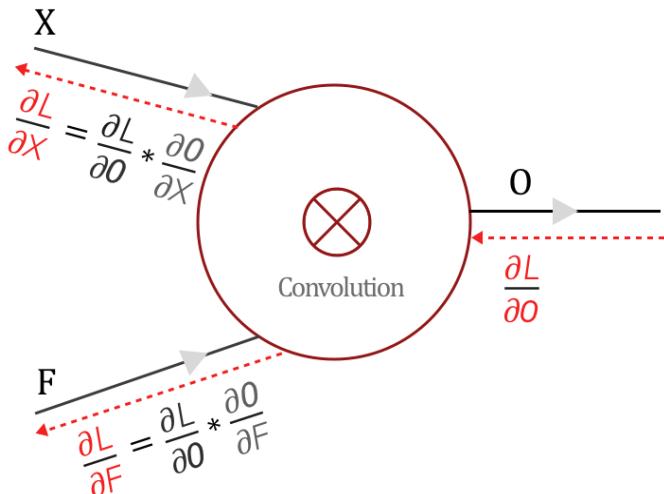
Loss gradient =  $dL/dC$

$$dC1=0.0846$$

$$dC2=0.1128$$

$$dC3=0.141$$

$$dC4=0.1692$$



$\frac{\partial O}{\partial X}$  &  $\frac{\partial O}{\partial F}$  are local gradients

$\frac{\partial L}{\partial Z}$  is the loss from the previous layer which has to be backpropagated to other layers

$$dL/db = (dL/dO) * (dO/db)$$

$$O_{11} = X_{11}F_{11} + X_{12}F_{12} + X_{21}F_{21} + X_{22}F_{22} + b_1$$

So  $dO/db = 1$ ;  $O = C$  in this case

$$dL/db_i = (dL/dC) * (dO/db) = (dL/dC)_i$$

Add each contribution together

$$\text{grad } b = dC1 + dC2 + dC3 + dC4$$

$$\text{grad } b = 0.0846 + 0.1128 + 0.141 + 0.1692 = 0.5076$$

$$b1+ = b - 0.5 * 0.5076 = 0.1 - 0.2538 = -0.1538$$

1	0	1
0	1	0
1	0	1

X

w1	w2
w3	w4

I1 = new data

I2 = new data

I3 = new data

I4 = new data

I5 = new data

I6 = new data

I7 = new data

I8 = new data

I9 = new data

F1+=0.7731

F2+=-0.027

F3+=-0.027

F4+=0.7731

$$w1+=F1+=0.7731$$

$$w2+=F2+=-0.027$$

$$w3+=F3+=-0.027$$

$$w4+=F4+=0.7731$$

$$\begin{aligned} b1+ &= b1 - 0.5 * 0.5076 \\ &= 0.1 - 0.5 * 0.5076 \\ &= -0.1538 \end{aligned}$$

$$wh1+=0.1-0.5 \times 0.268=-0.034$$

$$wh2+=0.2-0.5 \times 0.042=0.179$$

$$wh3+=0.3-0.5 \times 0.042=0.279$$

$$wh4+=0.4-0.5 \times 0.268=0.266$$

$$wh5+=0.5-0.5 \times 0.268=0.364$$

$$wh6+=0.6-0.5 \times 0.042=0.579$$

$$wh7+=0.7-0.5 \times 0.042=0.679$$

$$wh8+=0.8-0.5 \times 0=268-0.666$$

$$bh1+=0.2-0.5 \times 0.1428=0.1295$$

$$bh2+=0.2-0.5 \times 0.1428=0.1295$$

$$wo1+= 0.1 - (0.5) * (-0.9996) = 0.558$$

$$wo2+= 0.2 - (0.5) * (-2.256) = 1.278$$

$$wo3+= 0.3 - (0.5) * (0.9995) = -0.158$$

$$wo4+= 0.4 - (0.5) * (2.256) = -0.7678$$

$$bo1+= 0.3 - (0.5) * (-0.705) = 0.625$$

$$bo2+= 0.3 - (0.5) * (0.705) = -0.0525$$

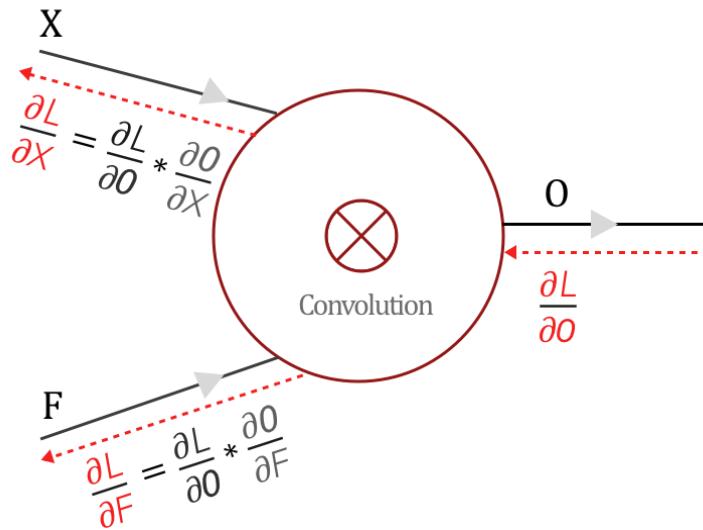
$$\begin{aligned} bo+ &= (bo1+bo2)/2 = \\ &= (-0.057 + 0.657) / .2 = 0.3 !! \end{aligned}$$

# Exercise

## CNN Backward Path Calculation

Given the following input, filter, and  $dL/dO$  matrices.

- Compute the matrix  $dL/dF$
- Compute the matrix  $dL/dX$



$\frac{\partial O}{\partial X}$  &  $\frac{\partial O}{\partial F}$  are local gradients

$\frac{\partial L}{\partial z}$  is the loss from the previous layer which has to be backpropagated to other layers

Input 3 x 3

1	0	-1
1	0	1
0	1	1

Filter = 2 x 2

-1	2
3	1

$dL/dO =$

-0.2	0.2
0.1	-0.1

## CNN Backward Path Calculation - Answers

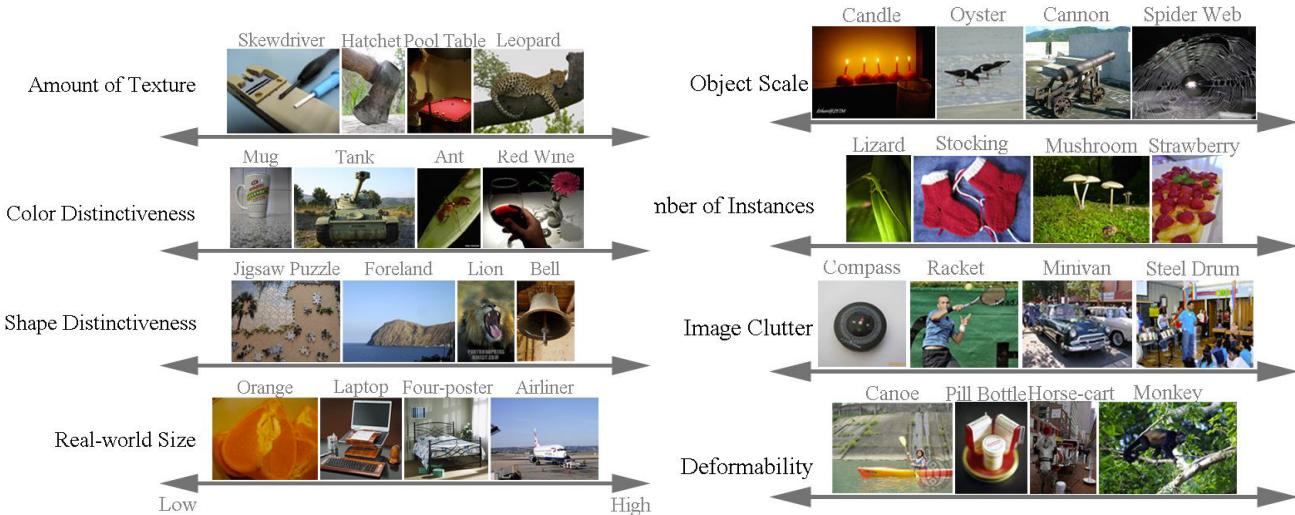
$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & 1 \\ \hline 0 & 1 & 1 \\ \hline \end{array} \quad \text{X} \quad \begin{array}{|c|c|} \hline -0.2 & 0.2 \\ \hline 0.1 & -0.1 \\ \hline \end{array} = \begin{array}{|c|c|} \hline -0.1 & -0.3 \\ \hline -0.3 & 0.2 \\ \hline \end{array} = dL/dF$$

dL/dX

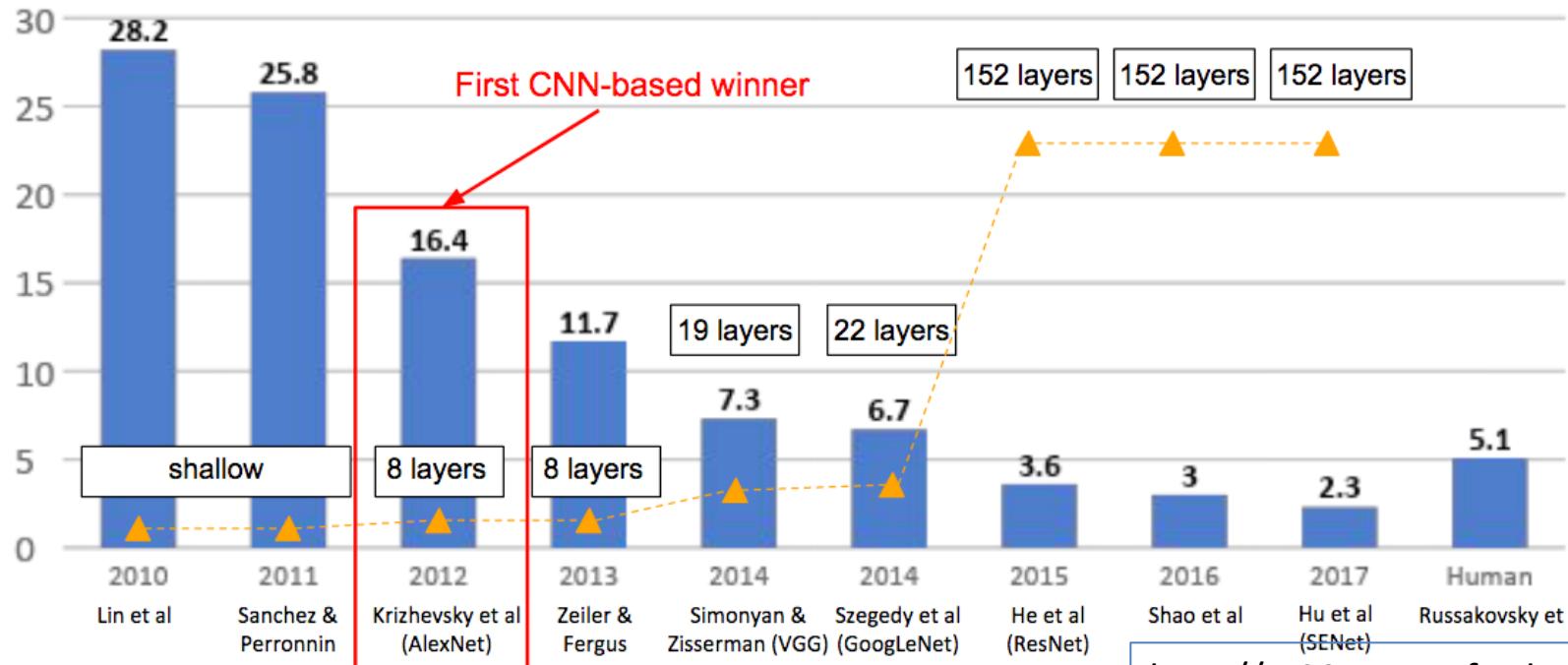
$$\begin{array}{|c|c|} \hline 1 & 3 \\ \hline 2 & -1 \\ \hline \end{array} \quad \text{X} \quad \begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 0 & -0.2 & 0.2 \\ \hline 0 & 0.1 & -0.1 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 0.2 & -0.6 & 0.4 \\ \hline -0.7 & 0.7 & 0 \\ \hline 0.3 & -0.2 & -0.1 \\ \hline \end{array}$$

rotated 180 w

- ✓ 1000 classes
- ✓ RGB images
- ✓ Cluttered Images (one category)
- ✓ Full resolution images
- ✓ 1.2 million training images
- ✓ 100 thousand test images

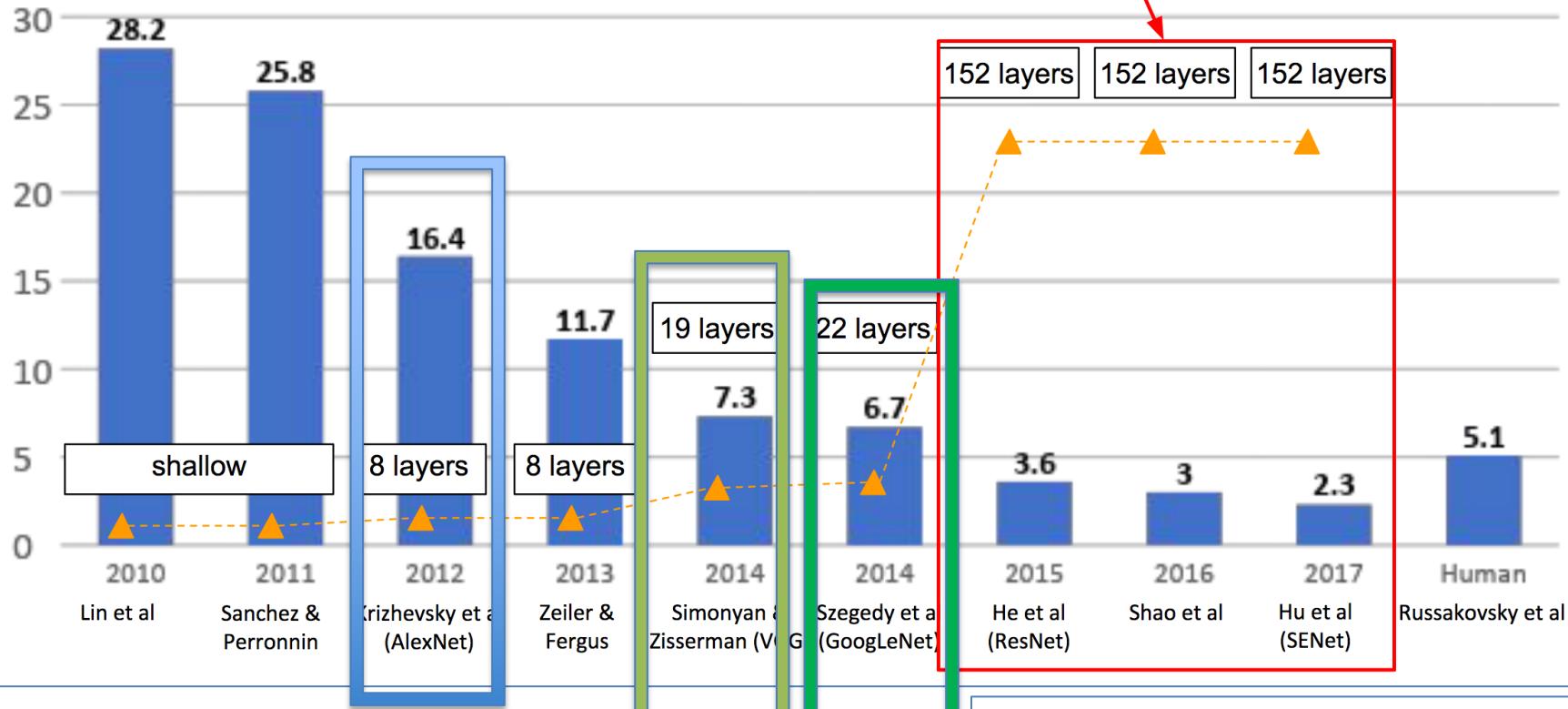


## ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners

“Revolution of Depth”



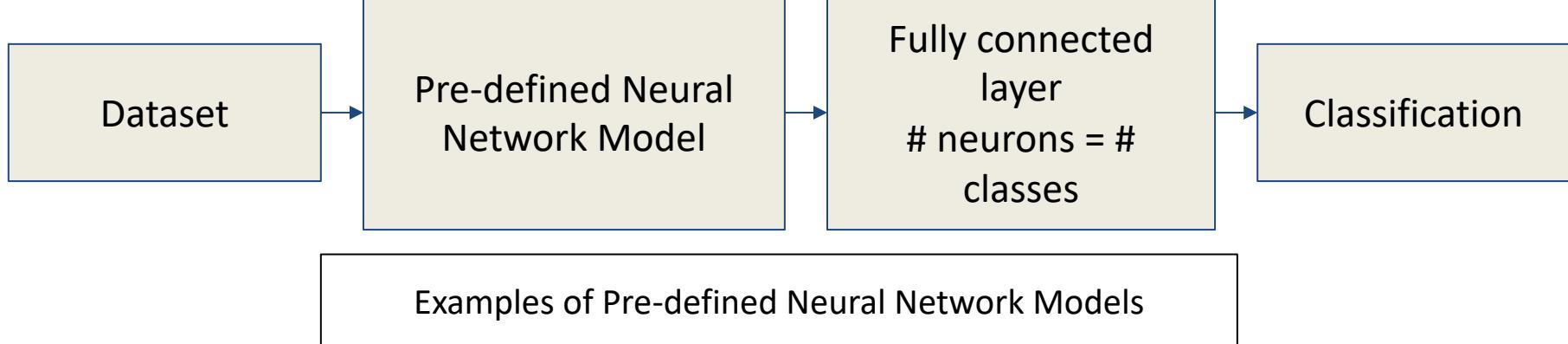
AlexNet – DNN  
Multiple layers  
~Ad hoc design

VGG – DNN  
Structure 3x3  
Deep NN,  
Lots of  
parameters  
Accuracy

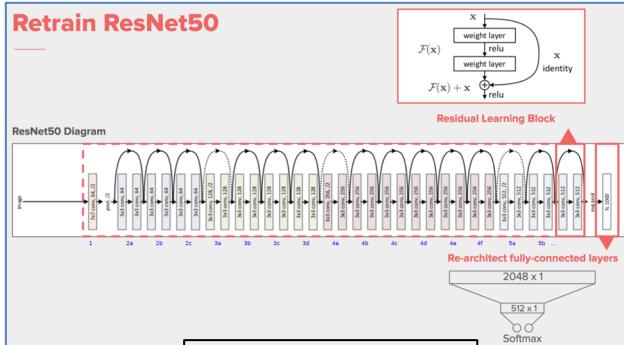
ResNet – DNN, Stem  
Residual block to preserve gradient flow  
Allow very deep network  
Accurate, efficiency, 18 t- 152

GoogleNet – DNN  
Stem to reduce feature size  
Block structure for efficiency  
Group pooling to reduce parameter count

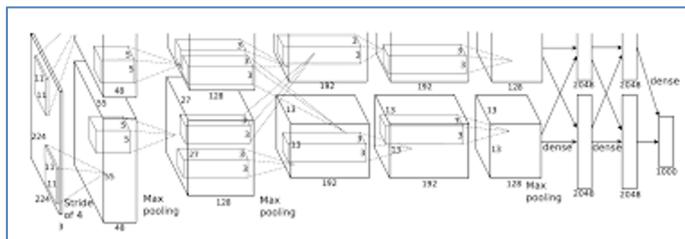
# Neural Network Models



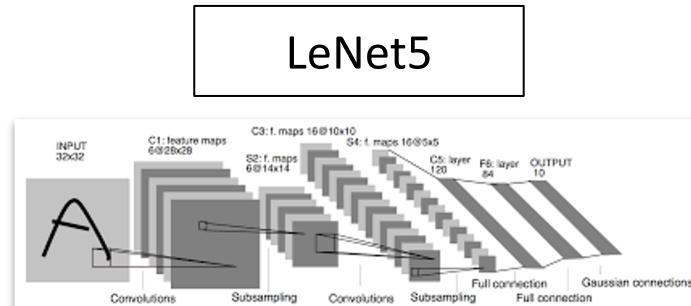
## ResNet50



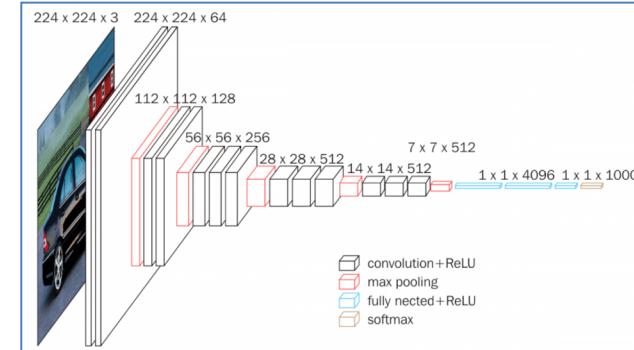
## AlexNet



## LeNet5

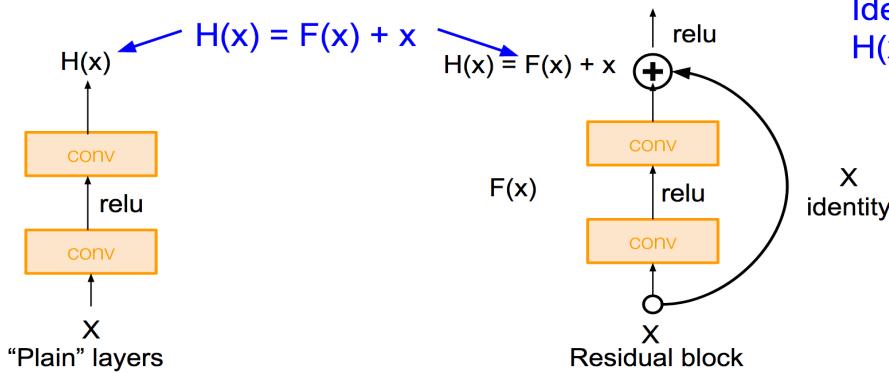


## VGG16



# ResNet

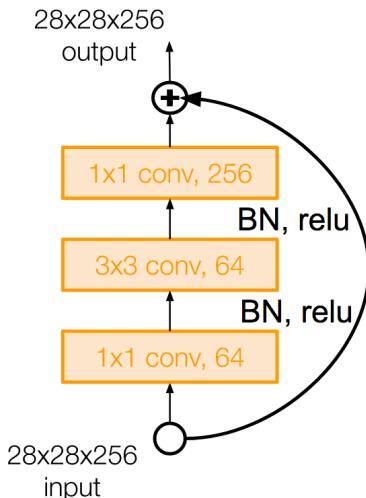
Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



For deeper networks (ResNet-50+), use “bottleneck” layer to improve efficiency (similar to GoogLeNet)

- ✓ Full ResNet architecture:-
- ✓ Stack residual blocks
- ✓ Every residual block has two 3x3 conv layers
- ✓ Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
- ✓ Additional conv layer at the beginning (stem)
- ✓ No FC layers at the end (only FC 1000 to output classes)

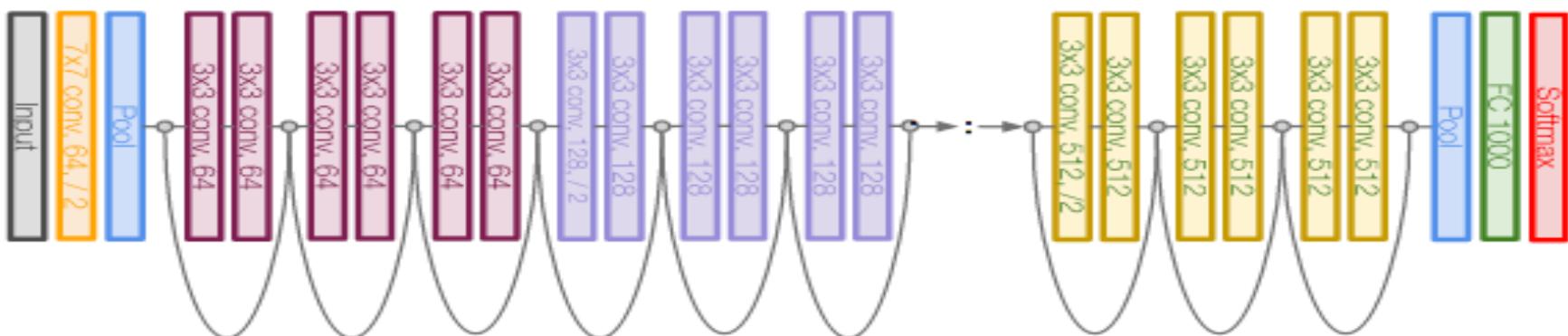
Total depths of 18, 34, 50, 101, or 152 layers for ImageNet



1x1 conv, 256 filters projects back to 256 feature maps (28x28x256)

3x3 conv operates over only 64 feature maps

1x1 conv, 64 filters to project to 28x28x64

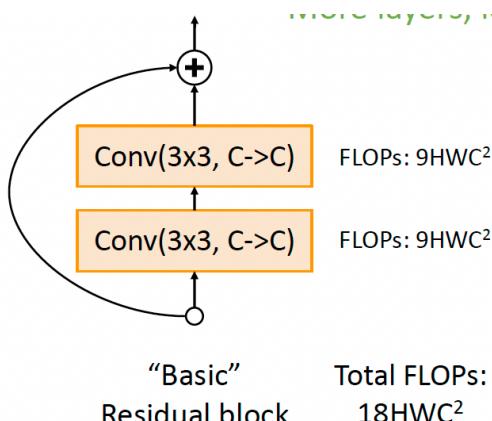


A residual network is a stack of many residual blocks

Regular design, like VGG: each residual block has two 3x3 conv

Network is divided into **stages**: the first block of each stage halves the resolution (with stride-2 conv) and doubles the number of channels

## Basic Block ResNet18, 34

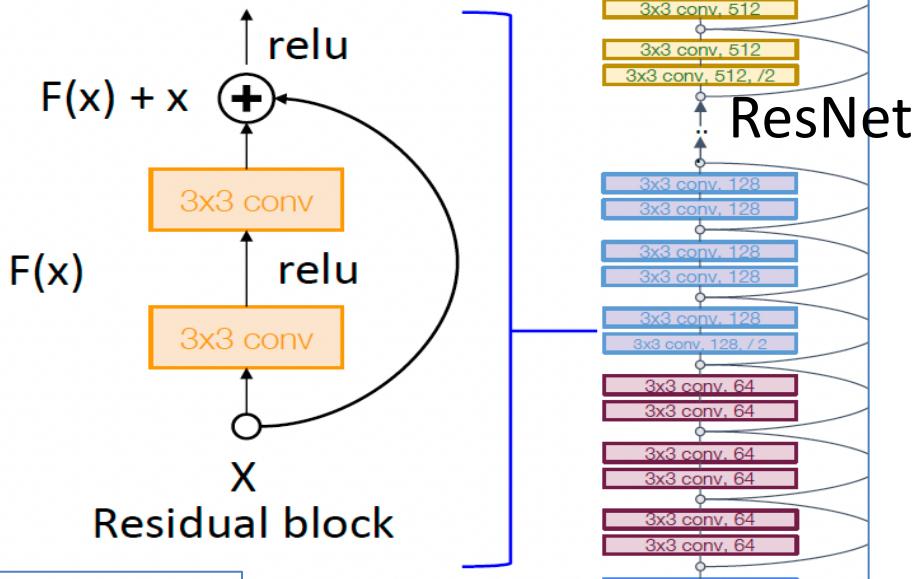


**ResNet-18:**  
Stem: 1 conv layer  
Stage 1 (C=64): 2 res. block = 4 conv  
Stage 2 (C=128): 2 res. block = 4 conv  
Stage 3 (C=256): 2 res. block = 4 conv  
Stage 4 (C=512): 2 res. block = 4 conv  
Linear

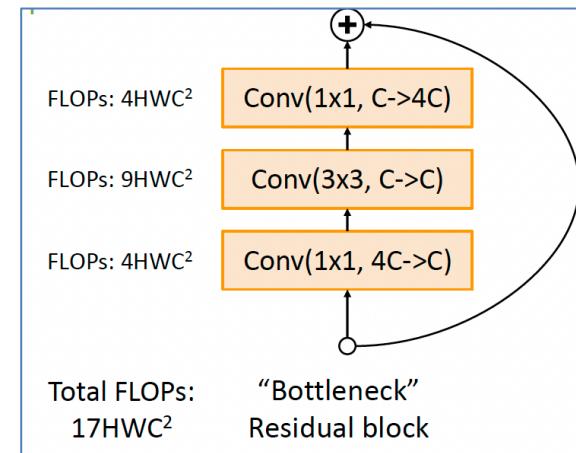
ImageNet top-5 error: 10.92  
GFLOP: 1.8

**ResNet-34:**  
Stem: 1 conv layer  
Stage 1: 3 res. block = 6 conv  
Stage 2: 4 res. block = 8 conv  
Stage 3: 6 res. block = 12 conv  
Stage 4: 3 res. block = 6 conv  
Linear

ImageNet top-5 error: 8.58  
GFLOP: 3.6



## Bottleneck Block ResNet50 ...



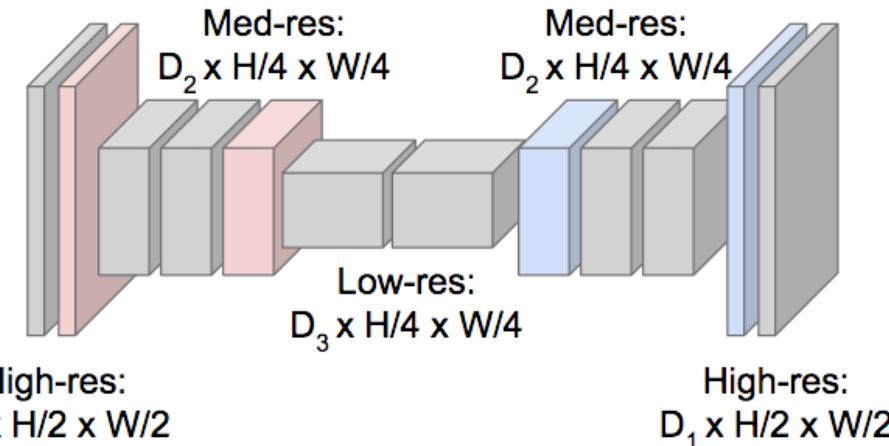
## Semantic Segmentation Idea: Fully Convolutional

**Downsampling:**  
Pooling, strided convolution



Input:  
 $3 \times H \times W$

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



**Upsampling:**  
Unpooling or strided transpose convolution



Predictions:  
 $H \times W$

### Max Pooling

Remember which element was max!

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



Input:  $4 \times 4$

Output:  $2 \times 2$

### Max Unpooling

Use positions from  
pooling layer

1	2
3	4



Input:  $2 \times 2$

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

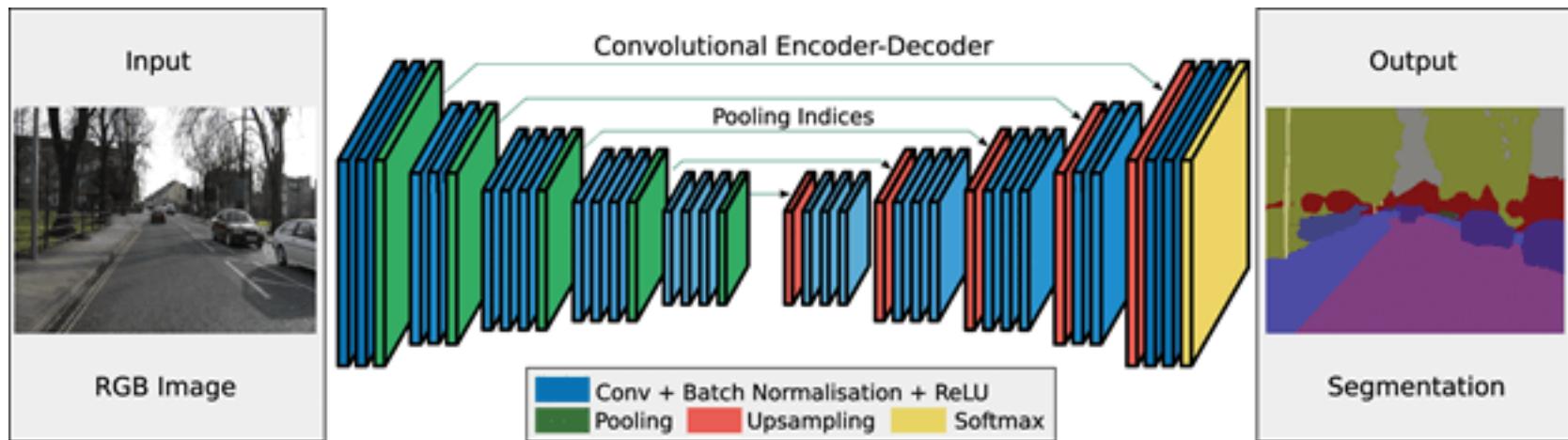
Output:  $4 \times 4$

# Image segmentation

Divide an image into multiple segments, every pixel in an image is associated with an object, instance segmentation and semantic segmentation :

In semantic segmentation, all objects of the same type are marked using one class label while in instance segmentation similar objects get their own separate labels.

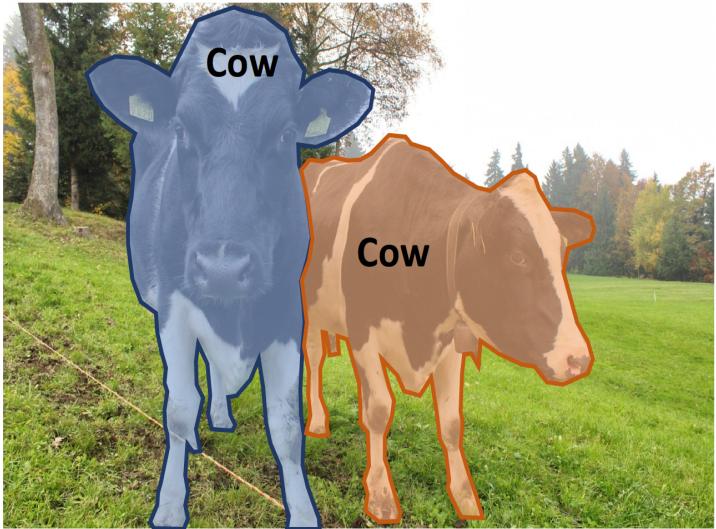
The basic architecture in image segmentation consists of an encoder and a decoder.



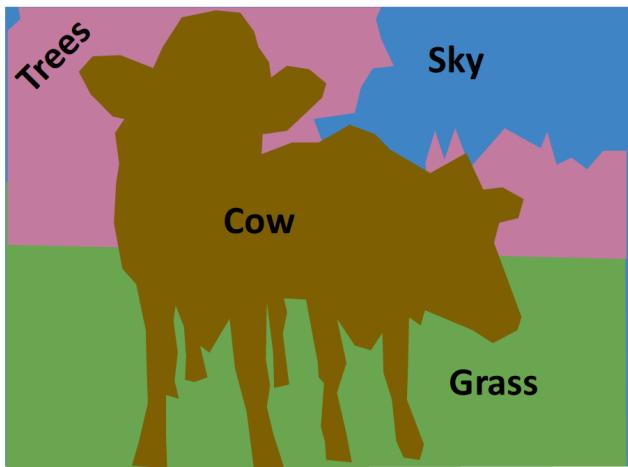
The encoder extracts features from the image through filters. The decoder is responsible for generating the final output which is usually a segmentation mask containing the outline of the object. Most of the architectures have this architecture or a variant of it.

**Instance Segmentation:**  
Detect all objects in the image, and identify the pixels that belong to each object (Only things!)

**Approach:** Perform object detection, then predict a segmentation mask for each object!



**Semantic Segmentation:** Identify both things and stuff, but doesn't separate instances



## Beyond Instance Segmentation: Human Keypoints

Represent the pose of a human by locating a set of **keypoints**

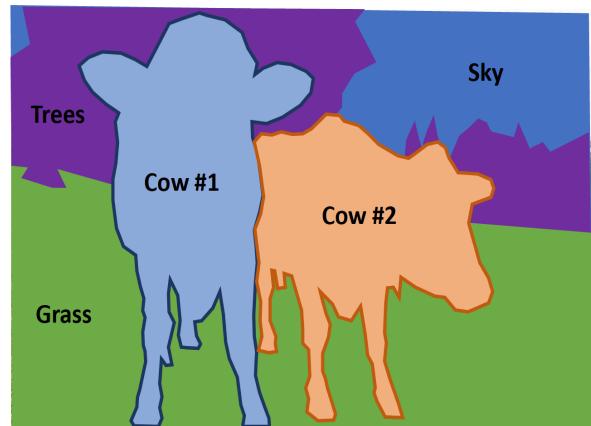
- e.g. 17 keypoints:
- Nose
  - Left / Right eye
  - Left / Right ear
  - Left / Right shoulder
  - Left / Right elbow
  - Left / Right wrist
  - Left / Right hip
  - Left / Right knee
  - Left / Right ankle



## Beyond Instance Segmentation: Panoptic Segmentation

Label all pixels in the image (both things and stuff)

For “thing” categories also separate into instances



Kirillov et al., "Panoptic Segmentation", CVPR 2019

Kirillov et al., "Panoptic Feature Pyramid Networks", CVPR 2019

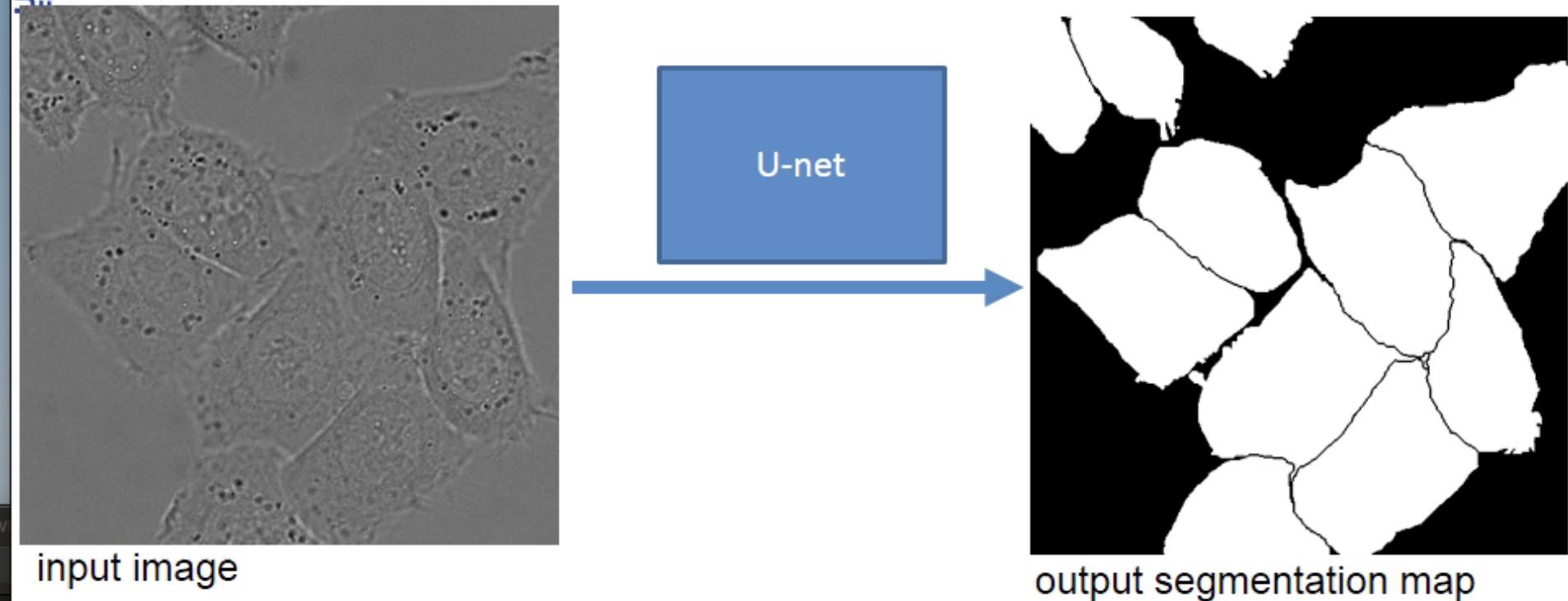
# U-net: Convolutional Networks for Biomedical Image Segmentation

---

Olaf Ronneberger, Philipp Fischer, and Thomas Brox

Computer Science Department and  
BIOSS Centre for Biological Signalling Studies  
University of Freiburg, Germany

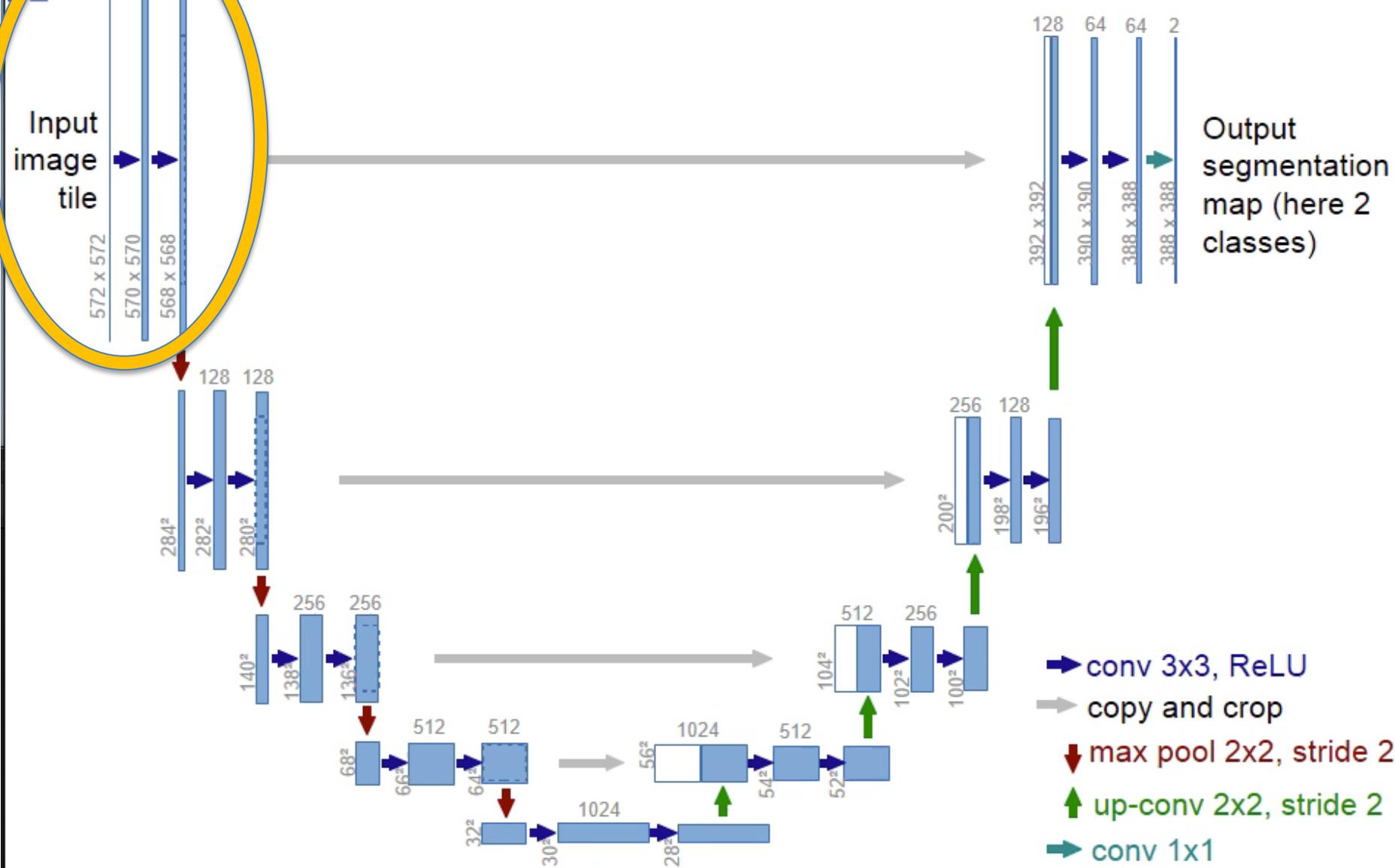
# Biomedical Image Segmentation with U-net

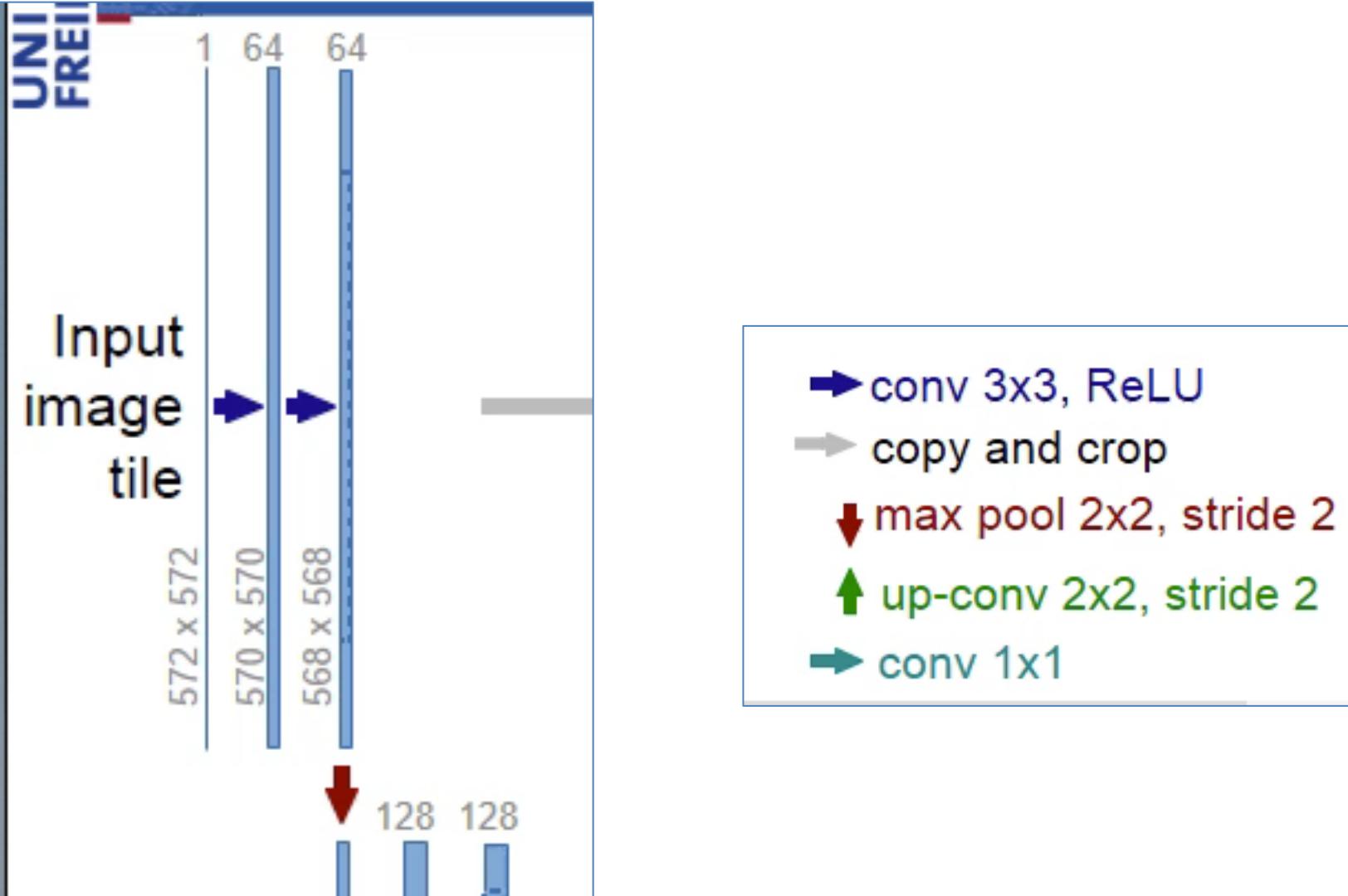


- U-net **learns segmentation** in an end-to-end setting
- **Very few annotated images** (approx. 30 per application)
- **Touching objects** of the same class

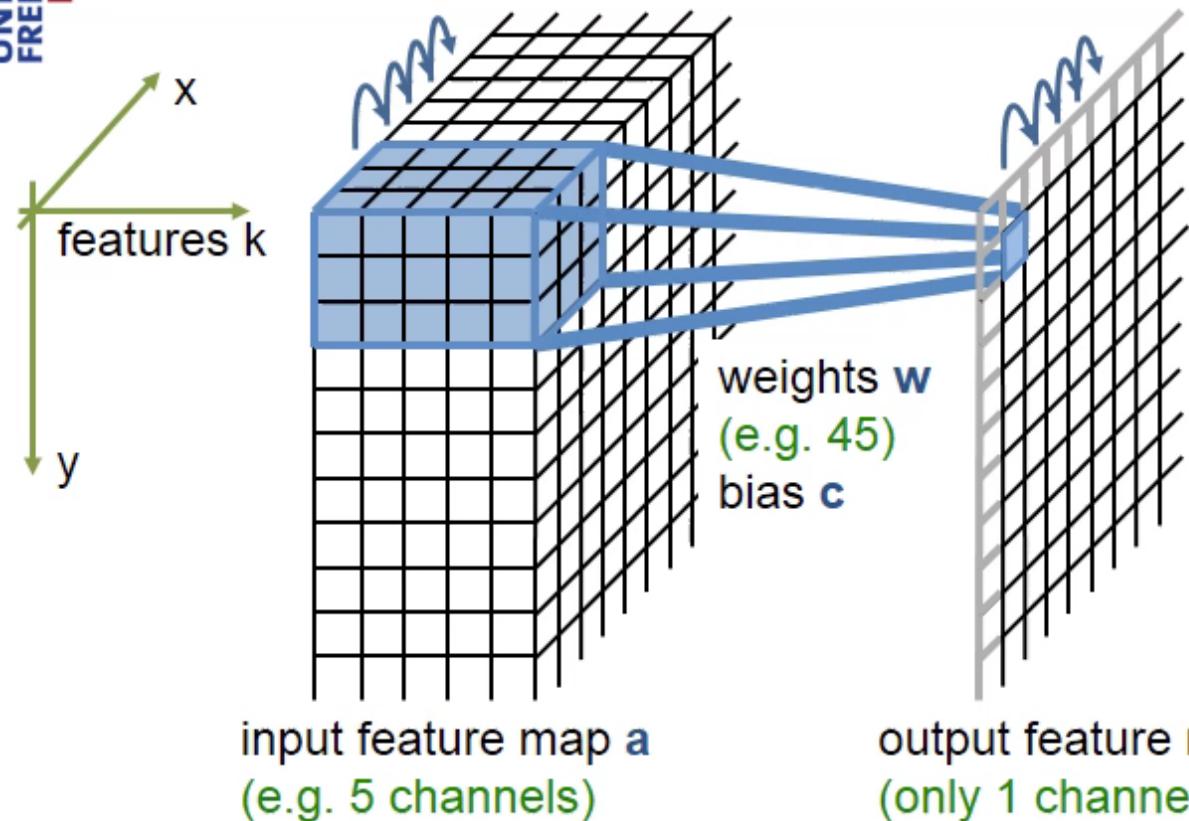
[Data provided by Dr. Gert van Cappellen, Erasmus Medical Center, Rotterdam, The Netherlands]

# U-net Architecture

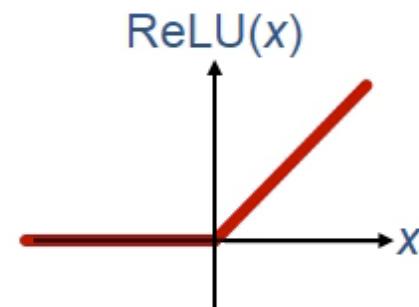




# 3x3 convolution + ReLU

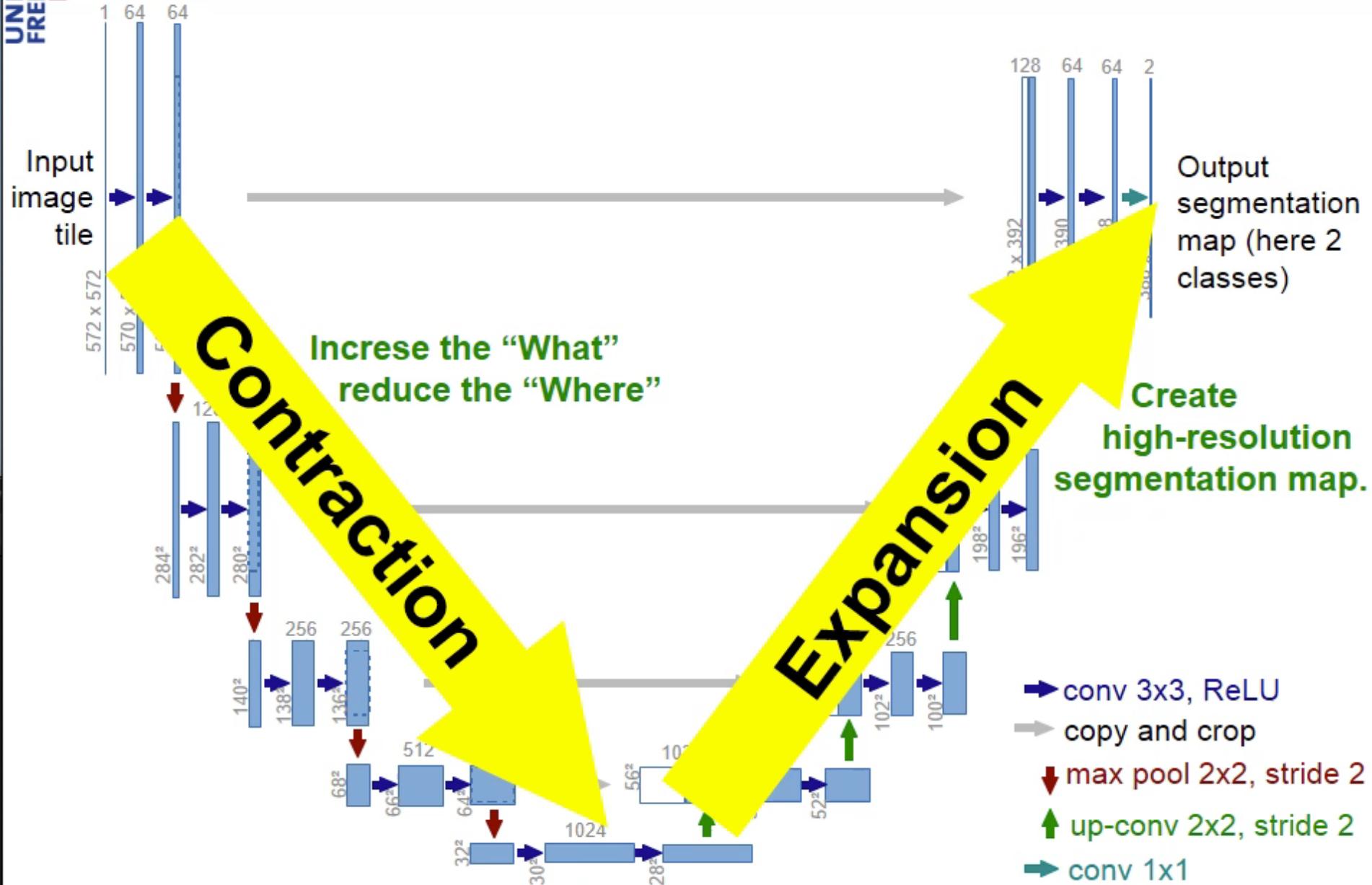


- Only valid part of convolution is used.
- For 3x3 convolutions a 1-pixel border is lost

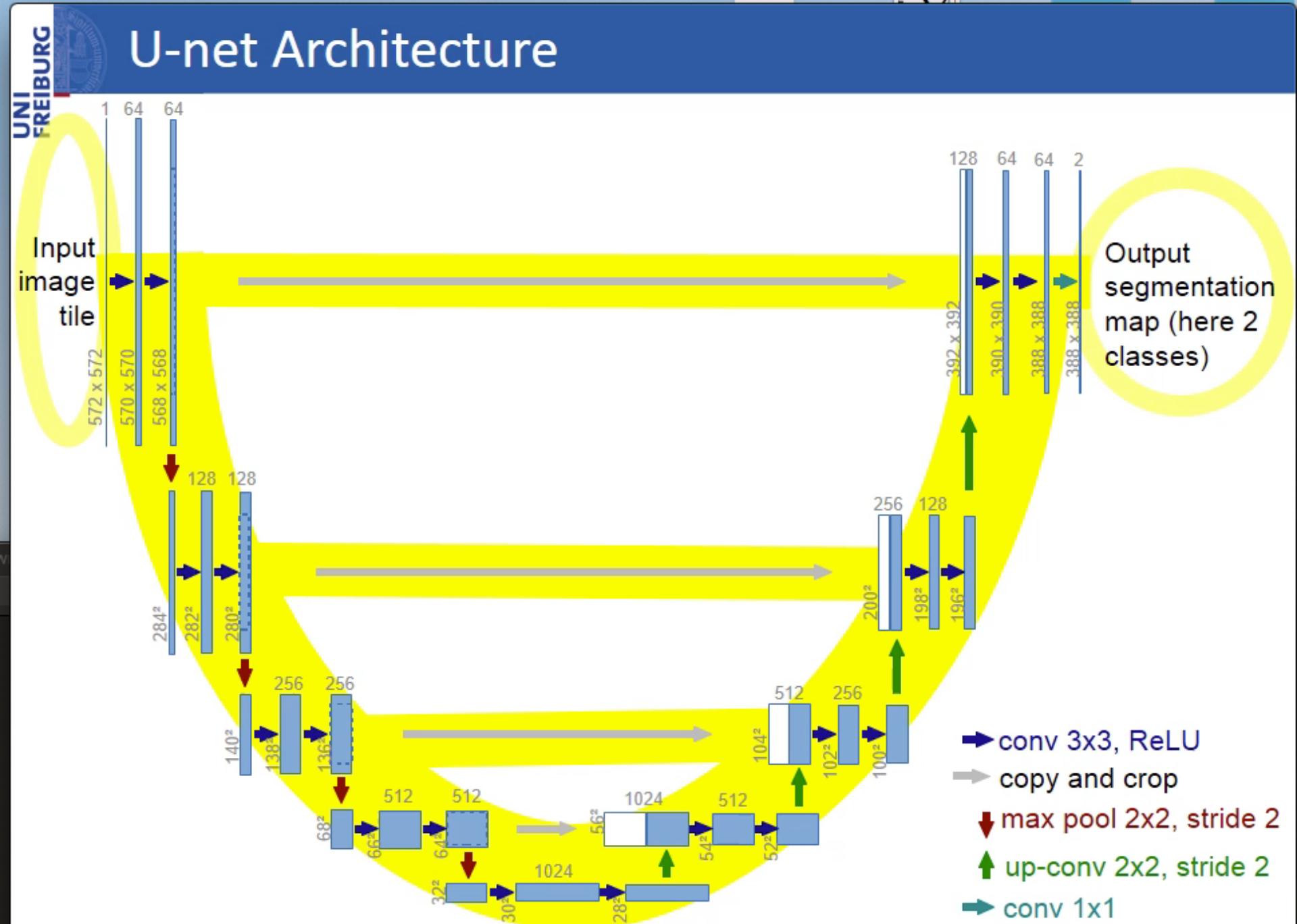


$$b_{x,y,l} = \text{ReLU}\left(\sum_{\substack{i \in \{-1,0,1\} \\ j \in \{-1,0,1\} \\ k \in \{1, \dots, K\}}} w_{i,j,k,l} \cdot a_{x+i,y+j,k} + c_l\right)$$

# U-Net Architecture

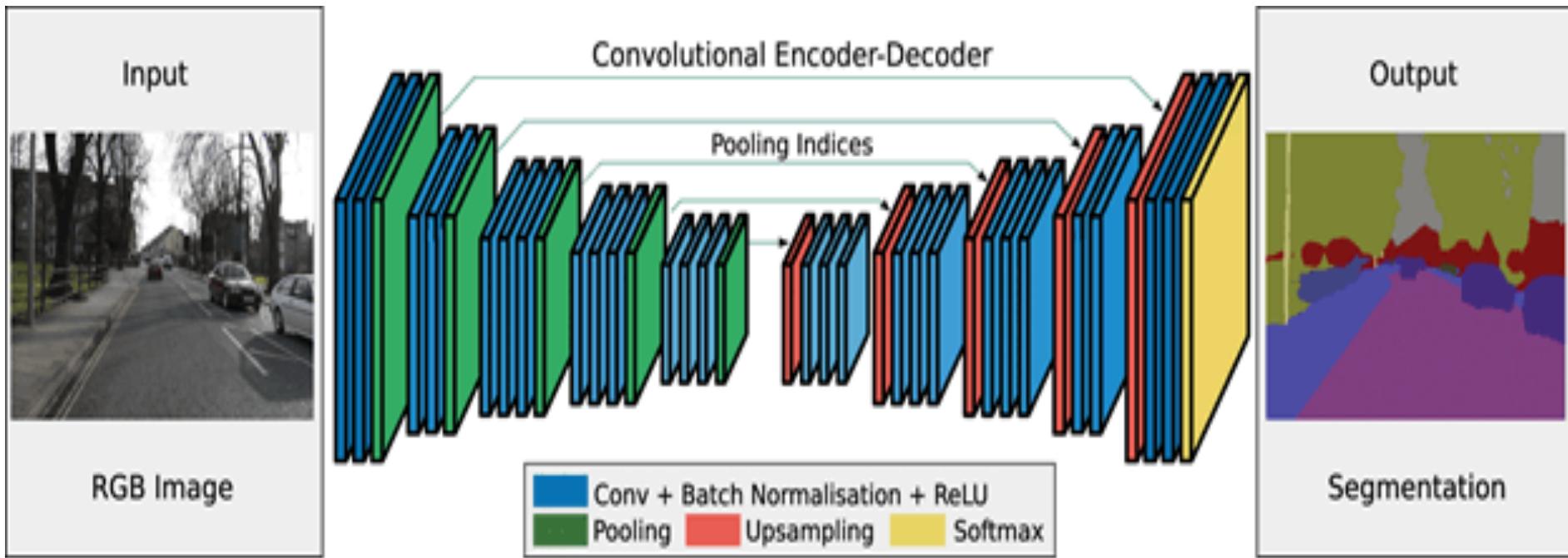


# U-net Architecture



# Image segmentation

The basic architecture in image segmentation consists of an encoder and a decoder.



The encoder extracts features from the image through filters. The decoder is responsible for generating the final output which is usually a segmentation mask containing the outline of the object. Most of the architectures have this architecture or a variant of it.

## CODE explanation

[https://keras.io/examples/vision/oxford\\_pets\\_image\\_segmentation/](https://keras.io/examples/vision/oxford_pets_image_segmentation/)



```
def get_model(img_size, num_classes):
    inputs = keras.Input(shape=img_size + (3,))

    ### [First half of the network: downsampling inputs] ###

    # Entry block
    x = layers.Conv2D(32, 3, strides=2, padding="same")(inputs)
    x = layers.BatchNormalization()(x)
    x = layers.Activation("relu")(x)

    previous_block_activation = x # Set aside residual

    # Blocks 1, 2, 3 are identical apart from the feature depth.
    for filters in [64, 128, 256]:
        x = layers.Activation("relu")(x)
        x = layers.SeparableConv2D(filters, 3, padding="same")(x)
        x = layers.BatchNormalization()(x)

        x = layers.Activation("relu")(x)
        x = layers.SeparableConv2D(filters, 3, padding="same")(x)
        x = layers.BatchNormalization()(x)

        x = layers.MaxPooling2D(3, strides=2, padding="same")(x)

        # Project residual
        residual = layers.Conv2D(filters, 1, strides=2, padding="same")(
            previous_block_activation
        )
        x = layers.add([x, residual]) # Add back residual
        previous_block_activation = x # Set aside next residual

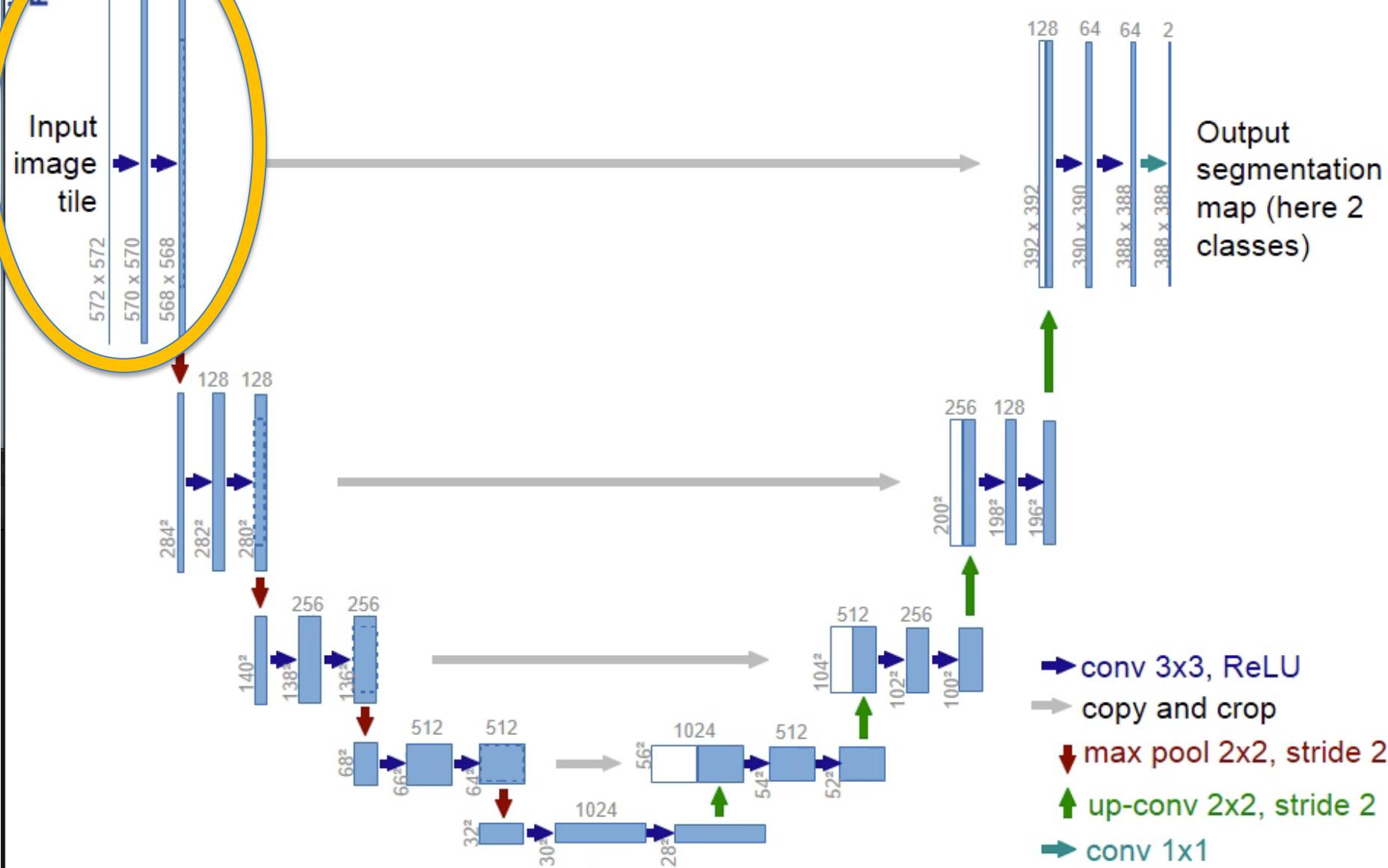
    # Final layers
    x = layers.SeparableConv2D(128, 3, padding="same")(x)
    x = layers.BatchNormalization()(x)
    x = layers.Activation("relu")(x)

    x = layers.GlobalAveragePooling2D()(x)
    x = layers.Dense(100, activation="softmax")(x)
    x = layers.Dropout(0.5)(x)

    return keras.Model(inputs, x)
```

```
input_dir = "images/"
target_dir = "annotat"
img_size = (160, 160)
num_classes = 3
batch_size = 32
```

# U-net Architecture



Img\_size = (160, 160, 3)

Num\_class=3

Batch\_size = 32

X = layers.Conv2D(32, 3, strides=2, padding='same') (input)

X = layers.BatchNormalization()(x)

X = layers.Activation("relu")(x)

→ conv 3x3, ReLU

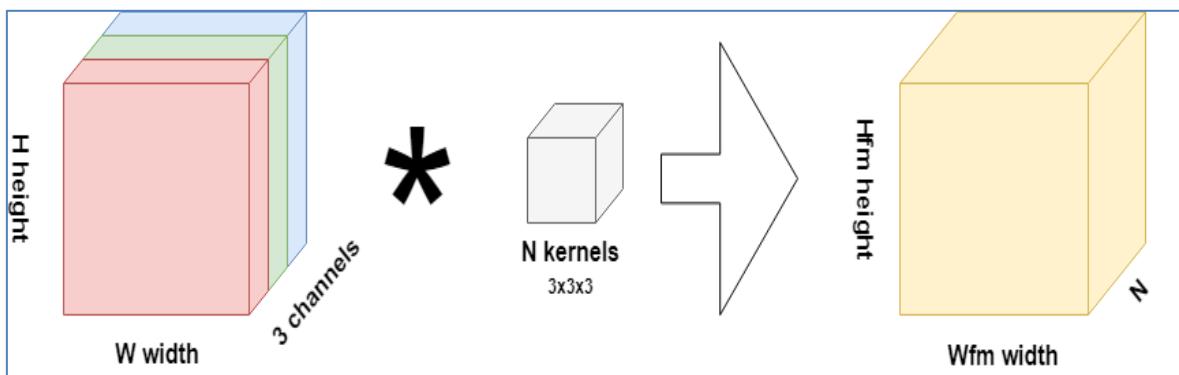
→ copy and crop

↓ max pool 2x2, stride 2

↑ up-conv 2x2, stride 2

→ conv 1x1

Height = 160, width = 160, depth (channel) = 3

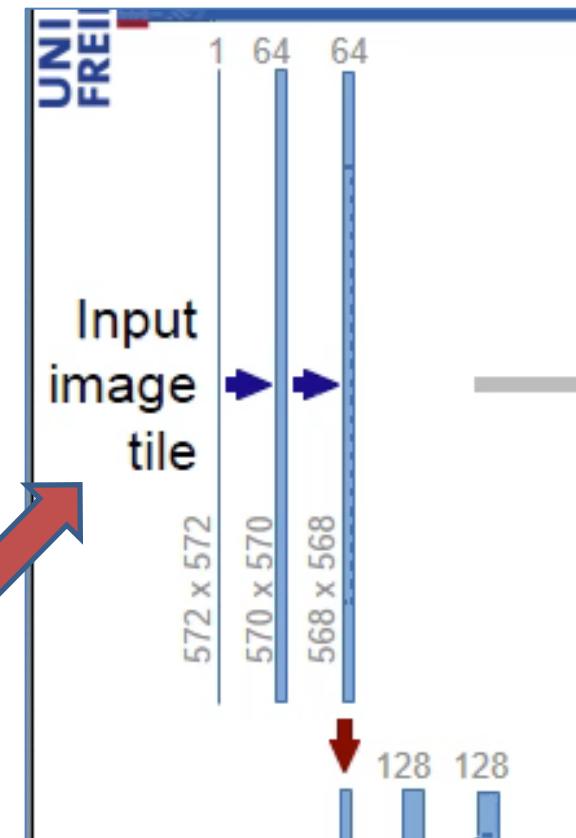


Input matrix =  $160 \times 160 \times 3$

$N = 32$

Stride = 2, same padding

output matrix =  $80 \times 80 \times 32$



Model: "functional\_1"

Layer (type)	Output Shape	Param #	Connected to
input_1 (InputLayer)	[(None, 160, 160, 3)]	0	(3x3x3+1) x32 = 896
conv2d (Conv2D)	(None, 80, 80, 32)	896	input_1[0][0]
batch_normalization (BatchNorma	(None, 80, 80, 32)	128	conv2d[0][0]
activation (Activation)	(None, 80, 80, 32)	0	batch_normalization[0]
activation_1 (Activation)	(None, 80, 80, 32)	0	activation[0][0]
separable_conv2d (SeparableConv	(None, 80, 80, 64)	2400	activation_1[0][0]
batch_normalization_1 (BatchNor	(None, 80, 80, 64)	256	separable_conv2d[0][0]
activation_2 (Activation)	(None, 80, 80, 64)	0	batch_normalization_1
separable_conv2d_1 (SeparableCo	(None, 80, 80, 64)	4736	activation_2[0][0]
batch_normalization_2 (BatchNor	(None, 80, 80, 64)	256	separable_conv2d_1[0]
max_pooling2d (MaxPooling2D)	(None, 40, 40, 64)	0	batch_normalization_2
conv2d_1 (Conv2D)	(None, 40, 40, 64)	2112	activation[0][0]
add (Add)	(None, 40, 40, 64)	0	max_pooling2d[0][0] conv2d_1[0][0]

```
def get_model(img_size, num_classes):
    inputs = keras.Input(shape=img_size + (3,))

    ### [First half of the network: downsampling inputs] ####

    # Entry block
    x = layers.Conv2D(32, 3, strides=2, padding="same")(inputs)
    x = layers.BatchNormalization()(x)
    x = layers.Activation("relu")(x)

    previous_block_activation = x # Set aside residual

    # Blocks 1, 2, 3 are identical apart from the feature depth.
    for filters in [64, 128, 256]:
        x = layers.Activation("relu")(x)
        x = layers.SeparableConv2D(filters, 3, padding="same")(x)
        x = layers.BatchNormalization()(x)

        x = layers.Activation("relu")(x)
        x = layers.SeparableConv2D(filters, 3, padding="same")(x)
        x = layers.BatchNormalization()(x)

        x = layers.MaxPooling2D(3, strides=2, padding="same")(x)

    # Project residual
    residual = layers.Conv2D(filters, 1, strides=2, padding="same")(
        previous_block_activation
    )
    x = layers.add([x, residual]) # Add back residual
    previous_block_activation = x # Set aside next residual
```

```
input_dir = "images/"
target_dir = "annotat
img_size = (160, 160)
num_classes = 3
batch_size = 32
```

No of Filters = 64

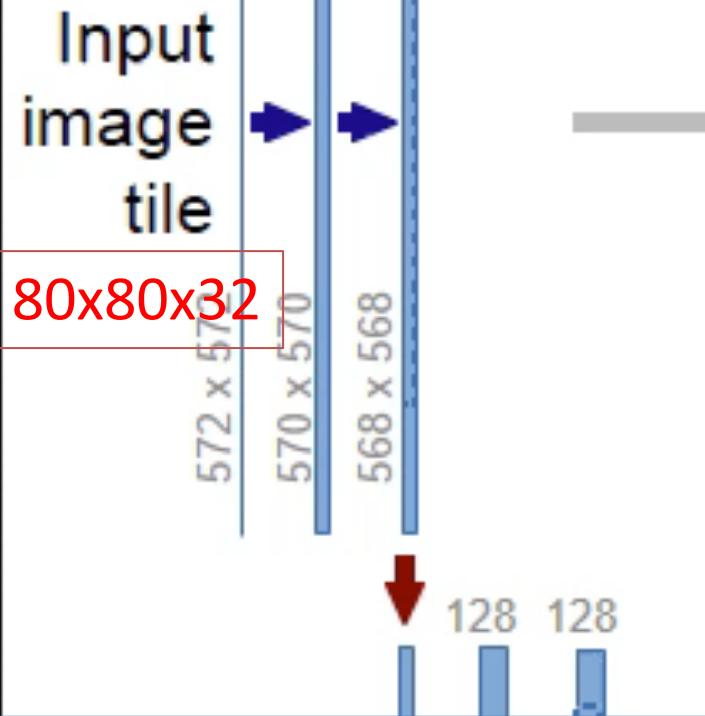
X=Input matrix =  $80 \times 80 \times 32$

N = 64

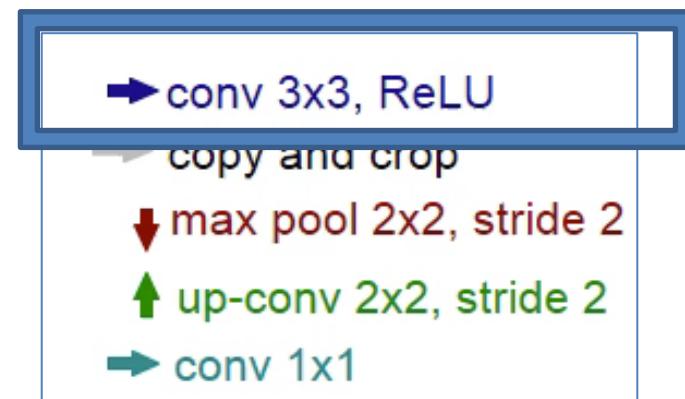
Filter size =  $3 \times 3$ , stride = 1, same padding

X = layers.Activation("relu")(x)  
X = layers.SeparableConv2D(64, 3, padding='same')(x)  
X = layers.BatchNormalization()(x)

output matrix =  $80 \times 80 \times 64$

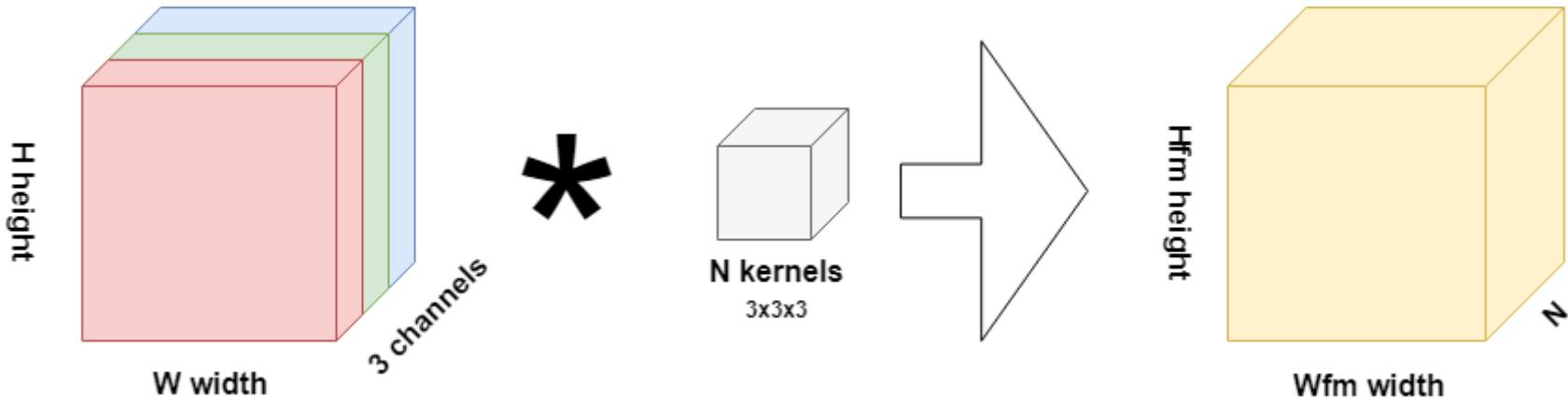


1<sup>st</sup> step : M = 32,  $(3 \times 3 \times 1) \times 32 = 288$   
2<sup>nd</sup> step : weights =  $(1 \times 1 \times 32) \times 64 = 2048$   
Bias = 64  
Total parameters = 2400



## Keras function : SeparableConv2D (....)

A little bit different from standard CONV2D but gets the same results  
Separate it to two stages

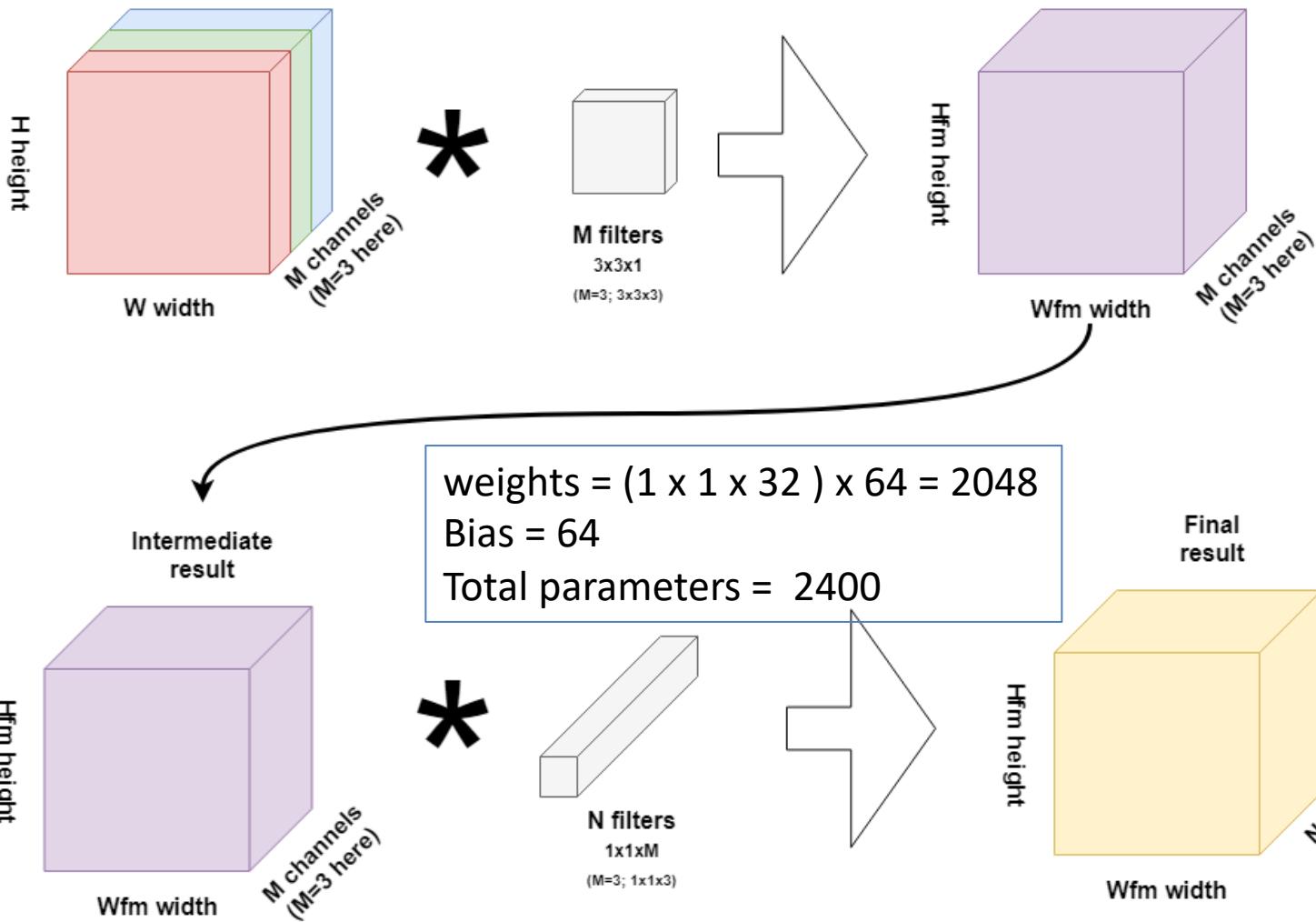


If your  $15 \times 15$  pixels image is RGB, and by consequence has 3 channels, you'll need  $(15-3+1) \times (15-3+1) \times 3 \times 3 \times N = 4563N$  multiplications to complete the full interpretation of *one image*. If you're working with ten kernels, so  $N = 10$ , you'll need over 45000 multiplications. Today, 20 to 50 kernels are not uncommon, datasets often span thousands of images and neural networks often compose multiple convolutional layers in their architecture.

<https://www.machinecurve.com/index.php/2019/09/24/creating-depthwise-separable-convolutions-in-keras/>

$$M \text{ (input channel)} = 32; \text{ weights} = (3 \times 3 \times 1) \times 32 = 288$$

Intermediate result



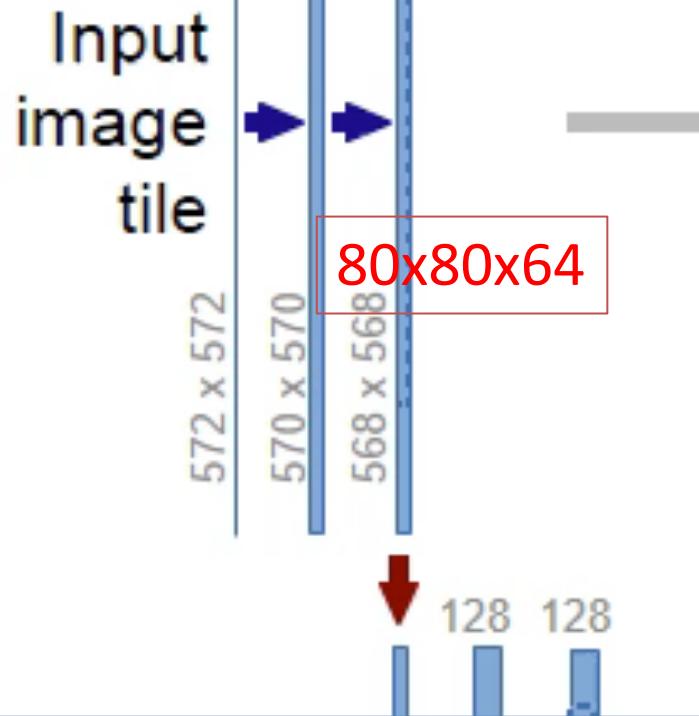
With those, you essentially split your  $N$  traditional kernels into *depthwise convolutions* and *pointwise convolutions*. In the first subprocess, you convolve with  $M$  filters on a layer basis, adding the kernels ‘pointwise’ in the second subprocess. While achieving the same result, you’ll need only **9633 convolutions**

No of Filters = 64

X= Input matrix =  $80 \times 80 \times 64$  N = 64

Filter size = 3 x3, stride = 1, same padding

X = layers.Activation("relu")(x)  
X = layers.SeparableConv2D(64, 3, padding='same')(x)  
X = layers.BatchNormalization()(x)



output matrix =  $80 \times 80 \times 64$

1<sup>st</sup> step : M = 64,  $(3 \times 3 \times 1) \times 64 = 576$

2<sup>nd</sup> step : weights =  $(1 \times 1 \times 64) \times 64 = 4096$

Bias = 64

Total parameters = 4736

→ conv 3x3, ReLU

copy and crop

↓ max pool 2x2, stride 2

↑ up-conv 2x2, stride 2

→ conv 1x1

**X = layers.BatchNormalization()(x)**

Batch Normalization has 4 parameters per set of filters (feature)  
Number of parameters = 4 x number of filters (bias) = 128, 256

As you can read there, in order to make the batch normalization work during training, they need to keep track of the distributions of each normalized dimensions. To do so, since you are in mode=0 by default, they compute 4 parameters per feature on the previous layer. Those parameters are making sure that you properly propagate and backpropagate the information. These parameters are in fact [gamma weights, beta weights, moving\_mean(non-trainable), moving\_variance(non-trainable)], of the size of the input layer.

[https://en.wikipedia.org/wiki/Batch\\_normalization](https://en.wikipedia.org/wiki/Batch_normalization)

input matrix = 80 x 80 x 64



**X = layers.MaxPooling2D(3, stride= 2, padding='same')(x)**  
**Reduce size by half because of stride, paddling**



output matrix = 40 x 40 x 64

Model: "functional\_1"

Layer (type)	Output Shape	Param #	Connected to
input_1 (InputLayer)	[None, 160, 160, 3]	0	(3x3x3+1) x32 = 896
conv2d (Conv2D)	(None, 80, 80, 32)	896	input_1[0][0]
batch_normalization (BatchNorma	(None, 80, 80, 32)	128	conv2d[0][0]
activation (Activation)	(None, 80, 80, 32)	0	batch_normalization[0]
activation_1 (Activation)	(None, 80, 80, 32)	0	activation[0][0]
separable_conv2d (SeparableConv	(None, 80, 80, 64)	2400	activation_1[0][0]
batch_normalization_1 (BatchNor	(None, 80, 80, 64)	256	separable_conv2d[0][0]
activation_2 (Activation)	(None, 80, 80, 64)	0	batch_normalization_1
separable_conv2d_1 (SeparableCo	(None, 80, 80, 64)	4736	activation_2[0][0]
batch_normalization_2 (BatchNor	(None, 80, 80, 64)	256	separable_conv2d_1[0][0]
max_pooling2d (MaxPooling2D)	(None, 40, 40, 64)	0	batch_normalization_2
conv2d_1 (Conv2D)	(None, 40, 40, 64)	2112	activation[0][0]
add (Add)	(None, 40, 40, 64)	0	max_pooling2d[0][0] conv2d_1[0][0]

```
def get_model(img_size, num_classes):
    inputs = keras.Input(shape=img_size + (3,))

    ### [First half of the network: downsampling inputs] ####

    # Entry block
    x = layers.Conv2D(32, 3, strides=2, padding="same")(inputs)
    x = layers.BatchNormalization()(x)
    x = layers.Activation("relu")(x)

    previous_block_activation = x # Set aside residual

    # Blocks 1, 2, 3 are identical apart from the feature depth.
    for filters in [64, 128, 256]:
        x = layers.Activation("relu")(x)
        x = layers.SeparableConv2D(filters, 3, padding="same")(x)
        x = layers.BatchNormalization()(x)

        x = layers.Activation("relu")(x)
        x = layers.SeparableConv2D(filters, 3, padding="same")(x)
        x = layers.BatchNormalization()(x)

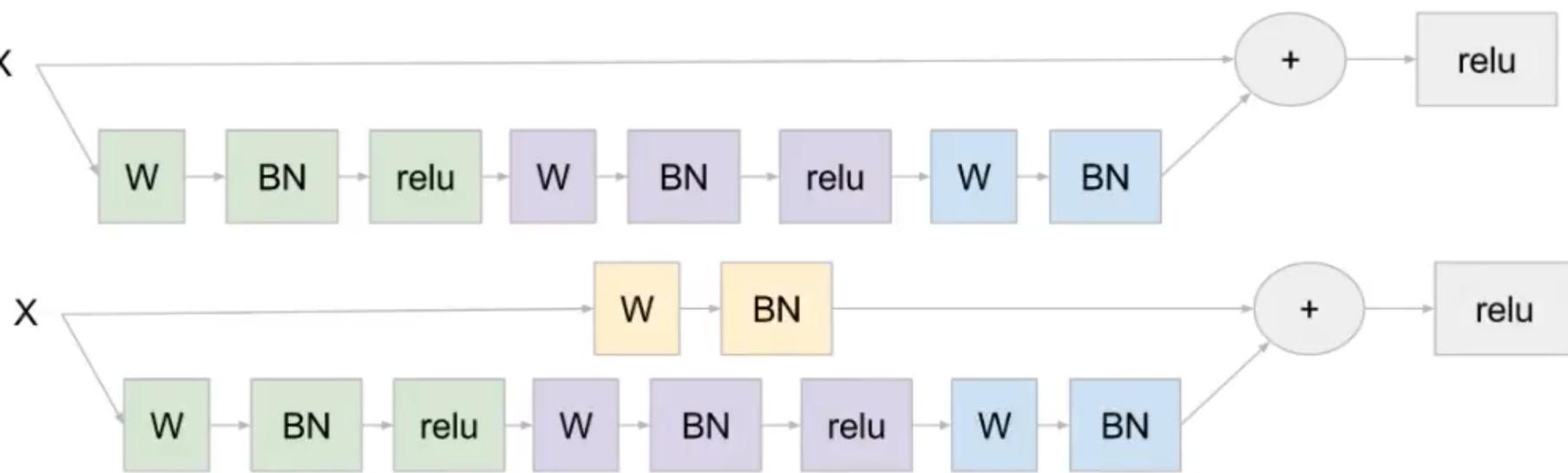
        x = layers.MaxPooling2D(3, strides=2, padding="same")(x)

    # Project residual
    residual = layers.Conv2D(filters, 1, strides=2, padding="same")(
        previous_block_activation
    )
    x = layers.add([x, residual]) # Add back residual
    previous_block_activation = x # Set aside next residual
```

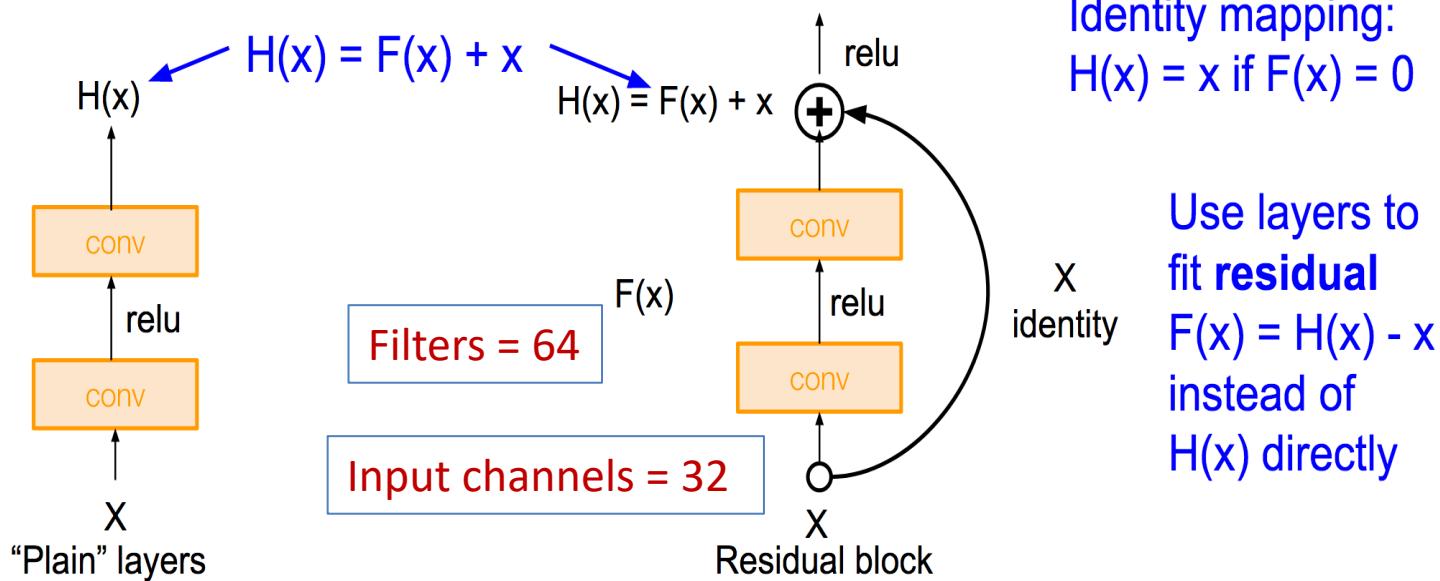
```
input_dir = "images/"
target_dir = "annotat
img_size = (160, 160)
num_classes = 3
batch_size = 32
```

Model: "functional\_1"

Layer (type)	Output Shape	Param #	Connected to
input_1 (InputLayer)	[None, 160, 160, 3]	0	(3x3x3+1) x32 = 896
conv2d (Conv2D)	(None, 80, 80, 32)	896	input_1[0][0]
batch_normalization (BatchNorma	(None, 80, 80, 32)	128	conv2d[0][0]
activation (Activation)	(None, 80, 80, 32)	0	batch_normalization[0]
activation_1 (Activation)	(None, 80, 80, 32)	0	activation[0][0]
separable_conv2d (SeparableConv	(None, 80, 80, 64)	2400	activation_1[0][0]
batch_normalization_1 (BatchNor	(None, 80, 80, 64)	256	separable_conv2d[0][0]
activation_2 (Activation)	(None, 80, 80, 64)	0	batch_normalization_1
separable_conv2d_1 (SeparableCo	(None, 80, 80, 64)	4736	activation_2[0][0]
batch_normalization_2 (BatchNor	(None, 80, 80, 64)	256	separable_conv2d_1[0][0]
max_pooling2d (MaxPooling2D)	(None, 40, 40, 64)	0	batch_normalization_2
conv2d_1 (Conv2D)	(None, 40, 40, 64)	2112	activation[0][0]
add (Add)	(None, 40, 40, 64)	0	max_pooling2d[0][0] conv2d_1[0][0]



Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



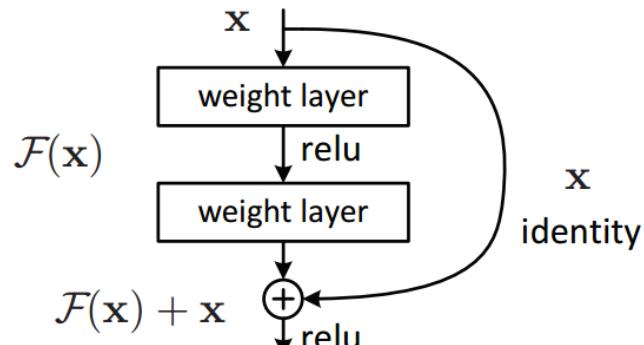
No of Filters = 64

Input matrix = 80 x 80 x 32

N = 64

1 64 64

80x80x32



572

570

568

40x40x64

128 128

Channel = 32, Filter size = 1 x1, stride = 2, same padding

residual = layers.Conv2D(64, 1, stride=2,  
padding='same')(previous\_block\_activation)  
 $X = \text{layers.add}(x, \text{residual})$

output matrix = 40 x 40 x 64

parameter,  $(1 \times 1 \times 32) \times 64 = 2048$   
Bias = 64

Total parameters = 2112

Add to the current X  
Size = 40x40x64

- conv 3x3, ReLU
- copy and crop
- ↓ max pool 2x2, stride 2
- ↑ up-conv 2x2, stride 2

→ conv 1x1

Model: "functional\_1"

Layer (type)	Output Shape	Param #	Connected to
input_1 (InputLayer)	[None, 160, 160, 3]	0	(3x3x3+1) x32 = 896
conv2d (Conv2D)	(None, 80, 80, 32)	896	input_1[0][0]
batch_normalization (BatchNorma	(None, 80, 80, 32)	128	conv2d[0][0]
activation (Activation)	(None, 80, 80, 32)	0	batch_normalization[0]
activation_1 (Activation)	(None, 80, 80, 32)	0	activation[0][0]
separable_conv2d (SeparableConv	(None, 80, 80, 64)	2400	activation_1[0][0]
batch_normalization_1 (BatchNor	(None, 80, 80, 64)	256	separable_conv2d[0][0]
activation_2 (Activation)	(None, 80, 80, 64)	0	batch_normalization_1
separable_conv2d_1 (SeparableCo	(None, 80, 80, 64)	4736	activation_2[0][0]
batch_normalization_2 (BatchNor	(None, 80, 80, 64)	256	separable_conv2d_1[0][0]
max_pooling2d (MaxPooling2D)	(None, 40, 40, 64)	0	batch_normalization_2
conv2d_1 (Conv2D)	(None, 40, 40, 64)	2112	activation[0][0]
add (Add)	(None, 40, 40, 64)	0	max_pooling2d[0][0] conv2d_1[0][0]

```
def get_model(img_size, num_classes):
    inputs = keras.Input(shape=img_size + (3,))

    ### [First half of the network: downsampling inputs] ###

    # Entry block
    x = layers.Conv2D(32, 3, strides=2, padding="same")(inputs)
    x = layers.BatchNormalization()(x)
    x = layers.Activation("relu")(x)

    previous_block_activation = x # Set aside residual

    # Blocks 1, 2, 3 are identical apart from the feature depth.
    for filters in [64, 128, 256]:
        x = layers.Activation("relu")(x)
        x = layers.SeparableConv2D(filters, 3, padding="same")(x)
        x = layers.BatchNormalization()(x)

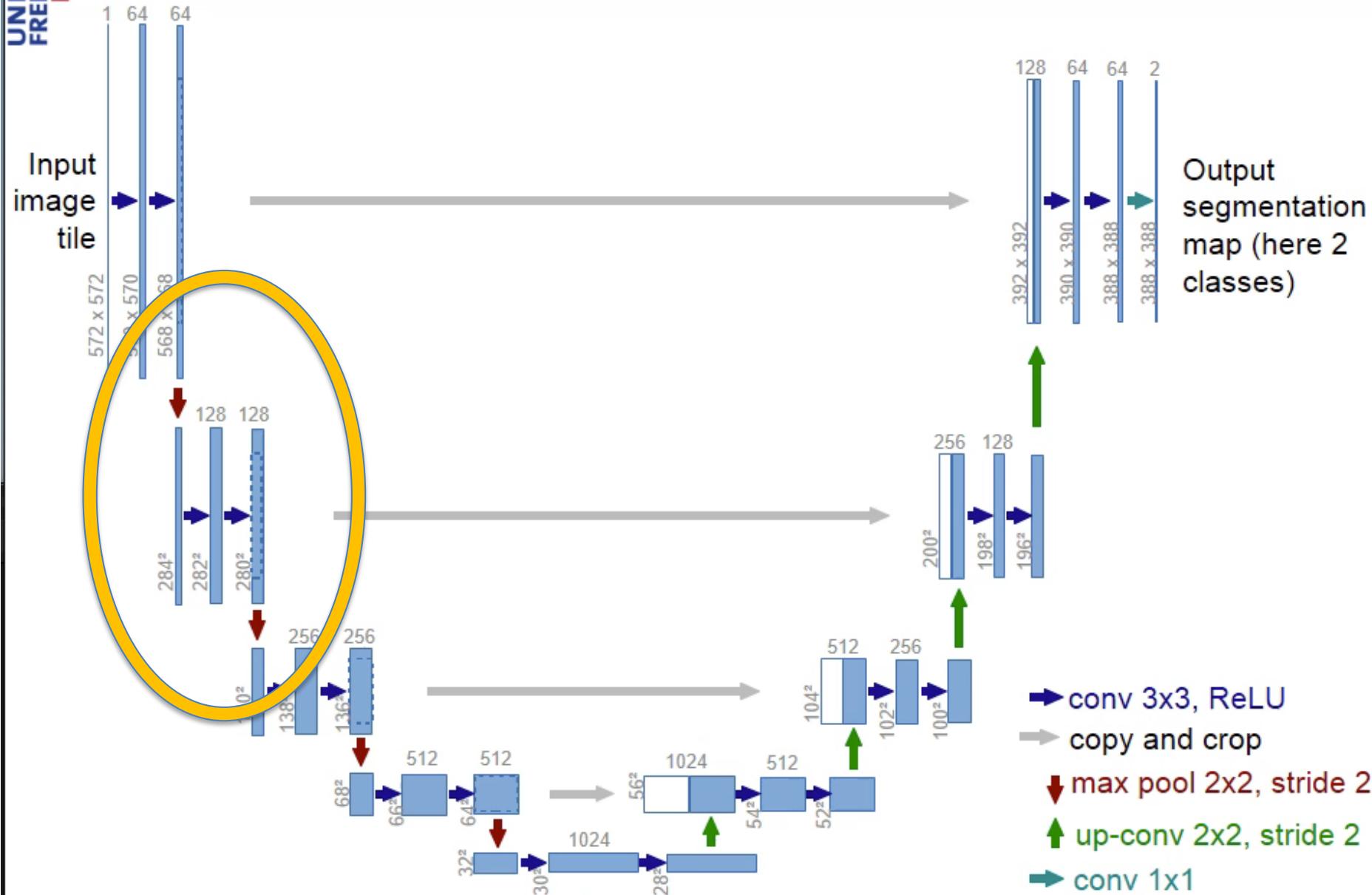
        x = layers.Activation("relu")(x)
        x = layers.SeparableConv2D(filters, 3, padding="same")(x)
        x = layers.BatchNormalization()(x)

        x = layers.MaxPooling2D(3, strides=2, padding="same")(x)

    # Project residual
    residual = layers.Conv2D(filters, 1, strides=2, padding="same")(
        previous_block_activation
    )
    x = layers.add([x, residual]) # Add back residual
    previous_block_activation = x # Set aside next residual
```

```
input_dir = "images/"
target_dir = "annotat"
img_size = (160, 160)
num_classes = 3
batch_size = 32
```

# U-net Architecture



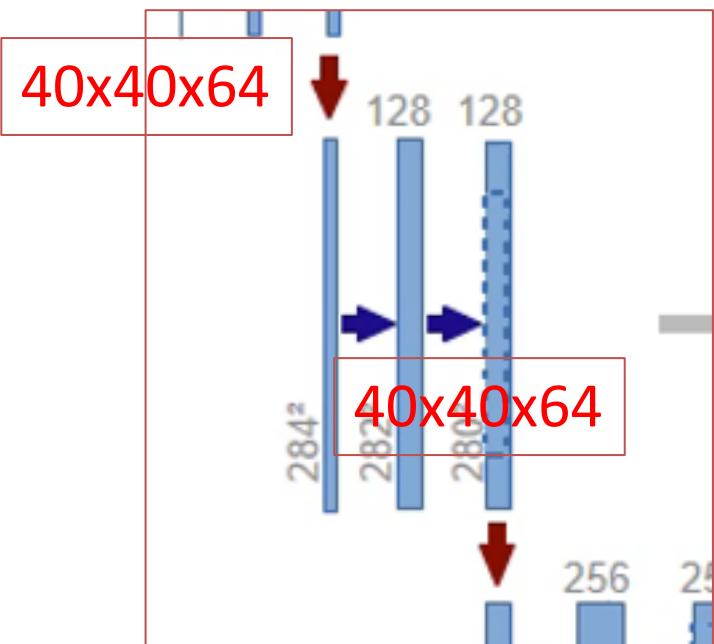
No of Filters = 128

Input matrix =  $40 \times 40 \times 64$

N = 128

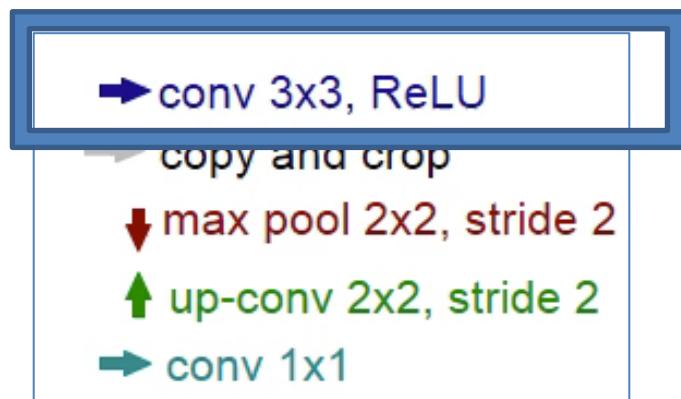
Filter size =  $3 \times 3$ , stride = 1, same padding

X = layers.Activation("relu")(x)  
X = layers.SeparableConv2D(128, 3, padding='same')(x)  
X = layers.BatchNormalization()(x)



output matrix =  $40 \times 40 \times 128$

1<sup>st</sup> step : M = 64,  $(3 \times 3 \times 1) \times 64 = 576$   
2<sup>nd</sup> step : weights =  $(1 \times 1 \times 64) \times 128 = 8192$   
Bias = 128  
Total parameters = 8896



No of Filters = 128

Input matrix =  $40 \times 40 \times 128$

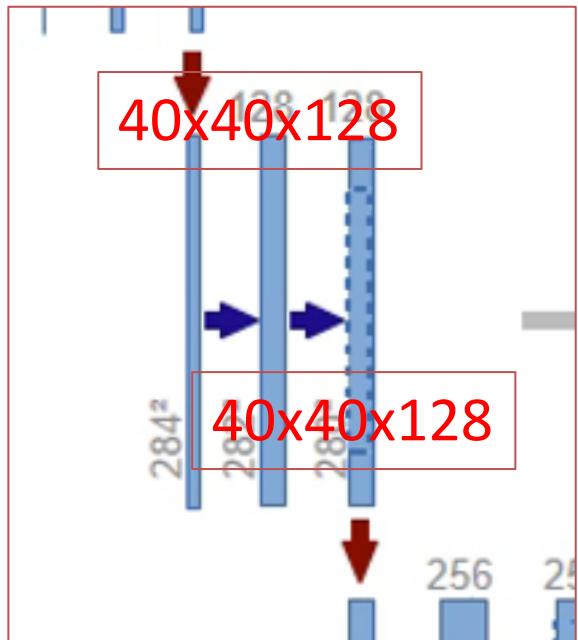
$N = 128$

Filter size =  $3 \times 3$ , stride = 1, same padding

$X = \text{layers.Activation}(\text{"relu"})(x)$

$X = \text{layers.SeparableConv2D}(128, 3, \text{padding}=\text{'same'})(x)$

$X = \text{layers.BatchNormalization}()(x)$



output matrix =  $40 \times 40 \times 128$

1<sup>st</sup> step :  $M = 128, (3 \times 3 \times 1) \times 128 = 1152$

2<sup>nd</sup> step : weights =  $(1 \times 1 \times 128) \times 128 = 16384$

Bias = 128

Total parameters = 17664

→ conv 3x3, ReLU

copy and crop

↓ max pool 2x2, stride 2

↑ up-conv 2x2, stride 2

→ conv 1x1

**X = layers.BatchNormalization()(x)**

Batch Normalization has 4 parameters per set of filters (feature)  
Number of parameters = 4 x number of filters (bias) = 256

As you can read there, in order to make the batch normalization work during training, they need to keep track of the distributions of each normalized dimensions. To do so, since you are in mode=0 by default, they compute 4 parameters per feature on the previous layer. Those parameters are making sure that you properly propagate and backpropagate the information. These parameters are in fact [gamma weights, beta weights, moving\_mean(non-trainable), moving\_variance(non-trainable)], of the size of the input layer.

[https://en.wikipedia.org/wiki/Batch\\_normalization](https://en.wikipedia.org/wiki/Batch_normalization)

input matrix = 40 x 40 x 128



**X = layers.MaxPooling2D(3, stride= 2, padding='same')(x)**  
**Reduce size by half because of stride, paddling**

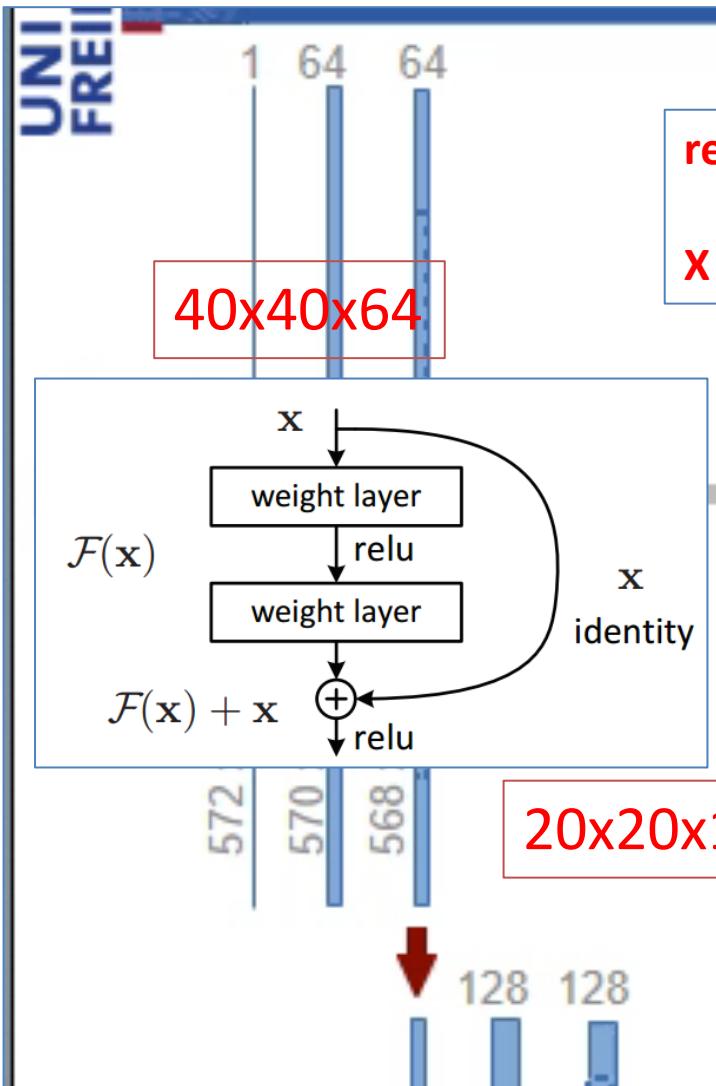


output matrix = 20 x 20 x 128

No of Filters = 128

Input matrix = 40 x 40 x 64

N = 128



Channel = 64, Filter size = 1 x1, stride = 2, same padding

```
residual = layers.Conv2D(128, 1, stride=2,  
padding='same')(previous_block_activation)  
X = layers.add(x, residual)
```

output matrix = 20 x 20 x 128

parameter,  $(1 \times 1 \times 64) \times 128 = 2048$   
Bias = 128  
Total parameters = 8320

Add to the current X  
Size = 20x20x128

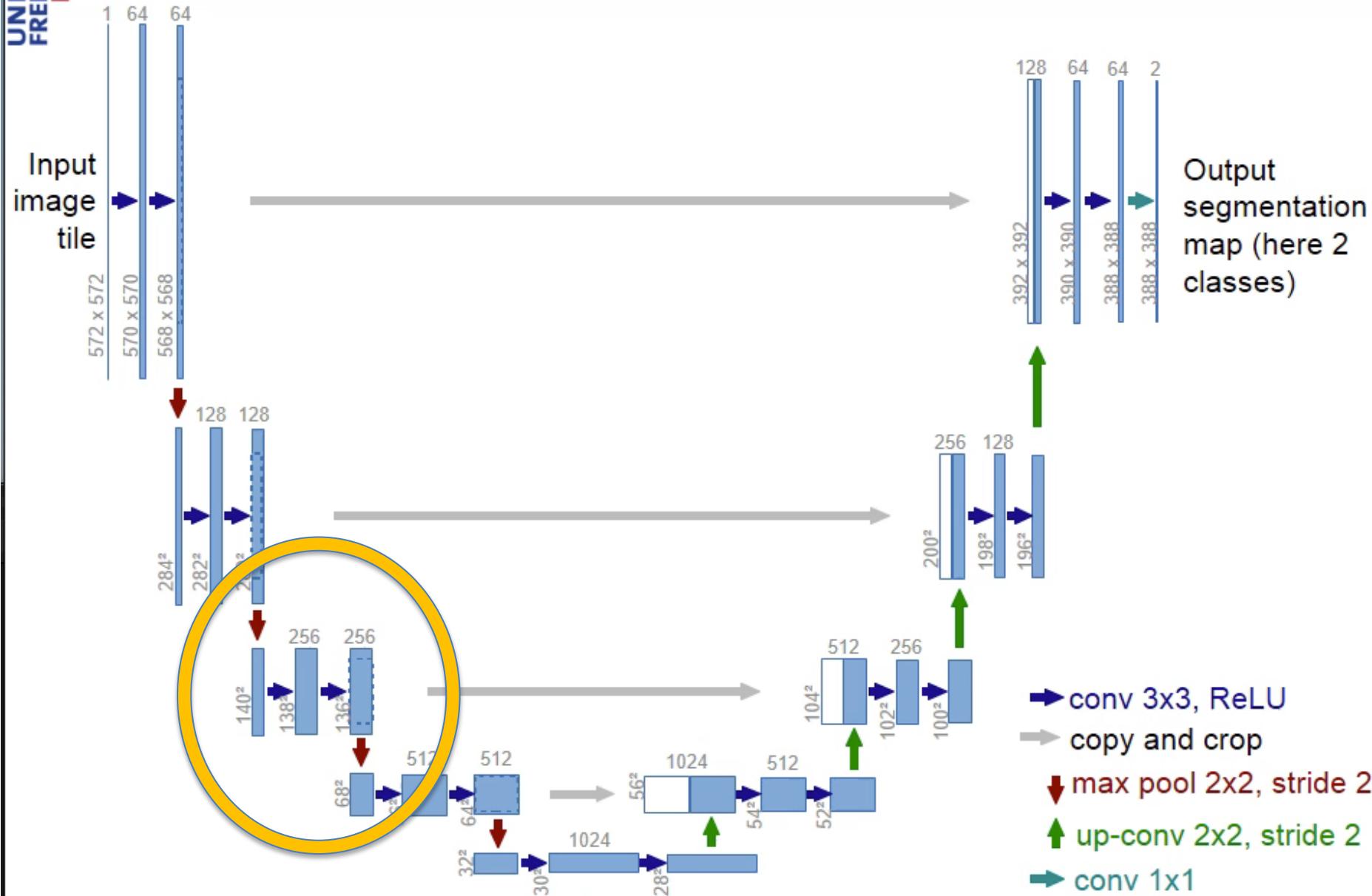
- conv 3x3, ReLU
- copy and crop
- ↓ max pool 2x2, stride 2
- ↑ up-conv 2x2, stride 2

→ conv 1x1

# Filter = 128, repeating the block!!

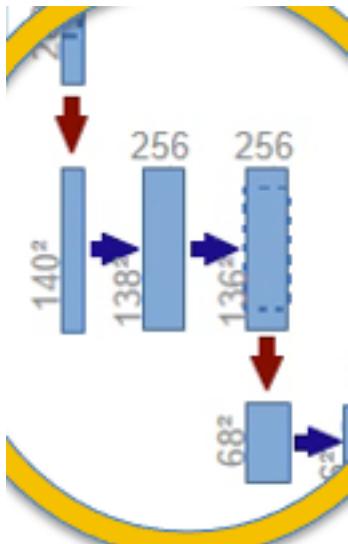
activation_3 (Activation)	(None, 40, 40, 64)	0	add[0] [0]
separable_conv2d_2 (SeparableCo	(None, 40, 40, 128)	8896	activation_3[0] [0]
batch_normalization_3 (BatchNor	(None, 40, 40, 128)	512	separable_conv2d_2[0]
activation_4 (Activation)	(None, 40, 40, 128)	0	batch_normalization_3
separable_conv2d_3 (SeparableCo	(None, 40, 40, 128)	17664	activation_4[0] [0]
batch_normalization_4 (BatchNor	(None, 40, 40, 128)	512	separable_conv2d_3[0]
max_pooling2d_1 (MaxPooling2D)	(None, 20, 20, 128)	0	batch_normalization_4
conv2d_2 (Conv2D)	(None, 20, 20, 128)	8320	add[0] [0]
add_1 (Add)	(None, 20, 20, 128)	0	max_pooling2d_1[0] [0] conv2d_2[0] [0]

# U-net Architecture



No of Filters = 256

20x20x128



20x20x256

Input matrix =  $20 \times 20 \times 128$

N = 256

Filter size = 3 x 3, stride = 1, same padding

X = layers.Activation("relu")(x)

X = layers.SeparableConv2D(256, 3, padding='same')(x)

X = layers.BatchNormalization()(x)

output matrix =  $20 \times 20 \times 256$

1<sup>st</sup> step : M = 128,  $(3 \times 3 \times 1) \times 128 = 1152$

2<sup>nd</sup> step : weights =  $(1 \times 1 \times 128) \times 256 = 32768$

Bias = 256

Total parameters = 34176

→ conv 3x3, ReLU

copy and crop

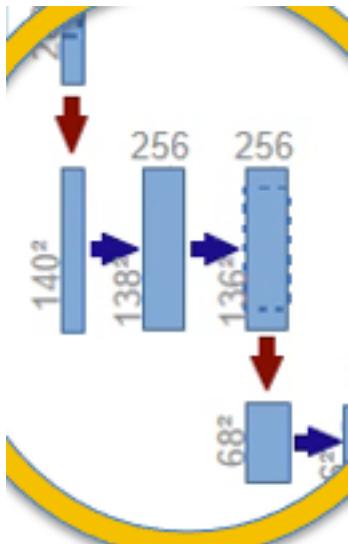
↓ max pool 2x2, stride 2

↑ up-conv 2x2, stride 2

→ conv 1x1

No of Filters = 256

20x20x256



20x20x256

Input matrix =  $20 \times 20 \times 256$

N = 256

Filter size =  $3 \times 3$ , stride = 1, same padding

X = layers.Activation("relu")(x)

X = layers.SeparableConv2D(256, 3, padding='same')(x)

X = layers.BatchNormalization()(x)

output matrix =  $20 \times 20 \times 256$

1<sup>st</sup> step : M = 256,  $(3 \times 3 \times 1) \times 256 = 2304$

2<sup>nd</sup> step : weights =  $(1 \times 1 \times 256) \times 256 = 65536$

Bias = 256

Total parameters = 68096

→ conv 3x3, ReLU

copy and crop

↓ max pool 2x2, stride 2

↑ up-conv 2x2, stride 2

→ conv 1x1

**X = layers.BatchNormalization()(x)**

Batch Normalization has 4 parameters per set of filters (feature)  
Number of parameters = 4 x number of filters (bias) = 1024

As you can read there, in order to make the batch normalization work during training, they need to keep track of the distributions of each normalized dimensions. To do so, since you are in mode=0 by default, they compute 4 parameters per feature on the previous layer. Those parameters are making sure that you properly propagate and backpropagate the information. These parameters are in fact [gamma weights, beta weights, moving\_mean(non-trainable), moving\_variance(non-trainable)], of the size of the input layer.

[https://en.wikipedia.org/wiki/Batch\\_normalization](https://en.wikipedia.org/wiki/Batch_normalization)

input matrix = 20 x 20 x 256



**X = layers.MaxPooling2D(3, stride= 2, padding='same')(x)**  
**Reduce size by half because of stride, paddling**

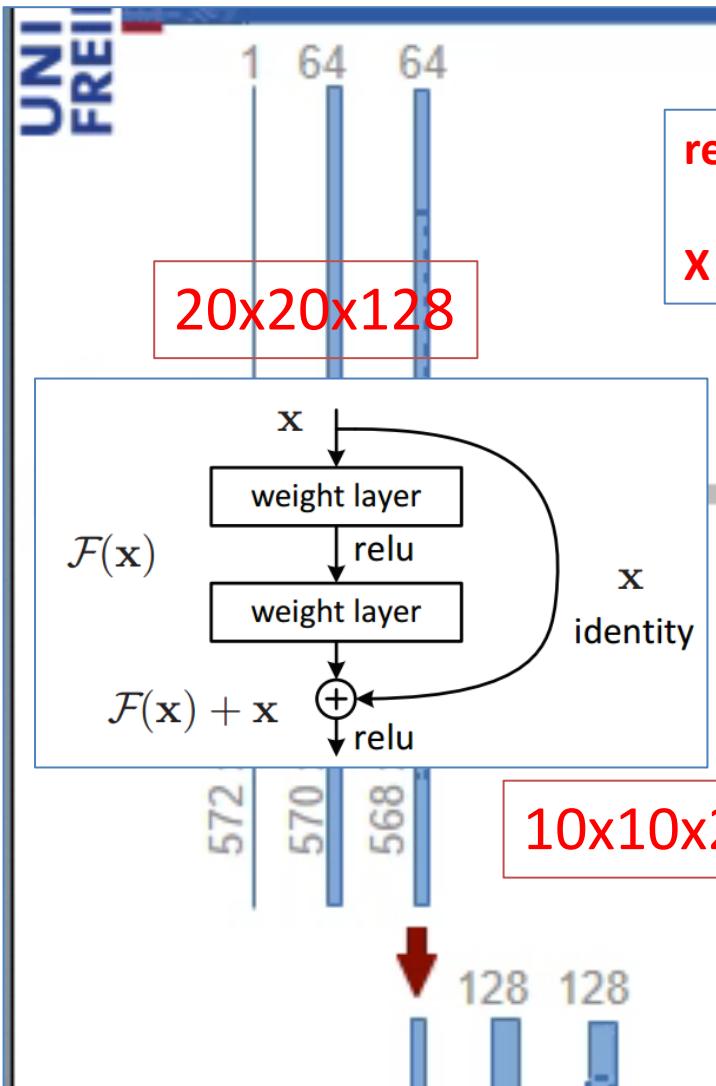


output matrix = 10 x 10 x 256

No of Filters = 256

Input matrix = 20 x 20 x 128

N = 256



Channel = 128, Filter size = 1 x1, stride = 2, same padding

```
residual = layers.Conv2D(256, 1, stride=2,  
padding='same')(previous_block_activation)  
X = layers.add(x, residual)
```

output matrix = 10 x 10 x 256

parameter,  $(1 \times 1 \times 128) \times 256 = 32768$   
Bias = 256

Total parameters = 33024

Add to the current X  
Size = 10x10x256

- conv 3x3, ReLU
- copy and crop
- ↓ max pool 2x2, stride 2
- ↑ up-conv 2x2, stride 2

→ conv 1x1

# Filter = 256, repeating the block!!

activation_5 (Activation)	(None, 20, 20, 128)	0	add_1[0][0]
separable_conv2d_4 (SeparableCo	(None, 20, 20, 256)	34176	activation_5[0][0]
batch_normalization_5 (BatchNor	(None, 20, 20, 256)	1024	separable_conv2d_4[0]
activation_6 (Activation)	(None, 20, 20, 256)	0	batch_normalization_5
separable_conv2d_5 (SeparableCo	(None, 20, 20, 256)	68096	activation_6[0][0]
batch_normalization_6 (BatchNor	(None, 20, 20, 256)	1024	separable_conv2d_5[0]
max_pooling2d_2 (MaxPooling2D)	(None, 10, 10, 256)	0	batch_normalization_6
conv2d_3 (Conv2D)	(None, 10, 10, 256)	33024	add_1[0][0]
add_2 (Add)	(None, 10, 10, 256)	0	max_pooling2d_2[0][0] conv2d_3[0][0]

```
def get_model(img_size, num_classes):
    inputs = keras.Input(shape=img_size + (3,))

    ### [First half of the network: downsampling inputs] ###

    # Entry block
    x = layers.Conv2D(32, 3, strides=2, padding="same")(inputs)
    x = layers.BatchNormalization()(x)
    x = layers.Activation("relu")(x)

    previous_block_activation = x # Set aside residual

    # Blocks 1, 2, 3 are identical apart from the feature depth.
    for filters in [64, 128, 256]:
        x = layers.Activation("relu")(x)
        x = layers.SeparableConv2D(filters, 3, padding="same")(x)
        x = layers.BatchNormalization()(x)

        x = layers.Activation("relu")(x)
        x = layers.SeparableConv2D(filters, 3, padding="same")(x)
        x = layers.BatchNormalization()(x)

        x = layers.MaxPooling2D(3, strides=2, padding="same")(x)

    # Project residual
    residual = layers.Conv2D(filters, 1, strides=2, padding="same")(
        previous_block_activation
    )
    x = layers.add([x, residual]) # Add back residual
    previous_block_activation = x # Set aside next residual
```

```
input_dir = "images/"
target_dir = "annotat"
img_size = (160, 160)
num_classes = 3
batch_size = 32
```

```
### [Second half of the network: upsampling inputs] ###
```

```
for filters in [256, 128, 64, 32]:  
    x = layers.Activation("relu")(x)  
    x = layers.Conv2DTranspose(filters, 3, padding="same")(x)  
    x = layers.BatchNormalization()(x)
```

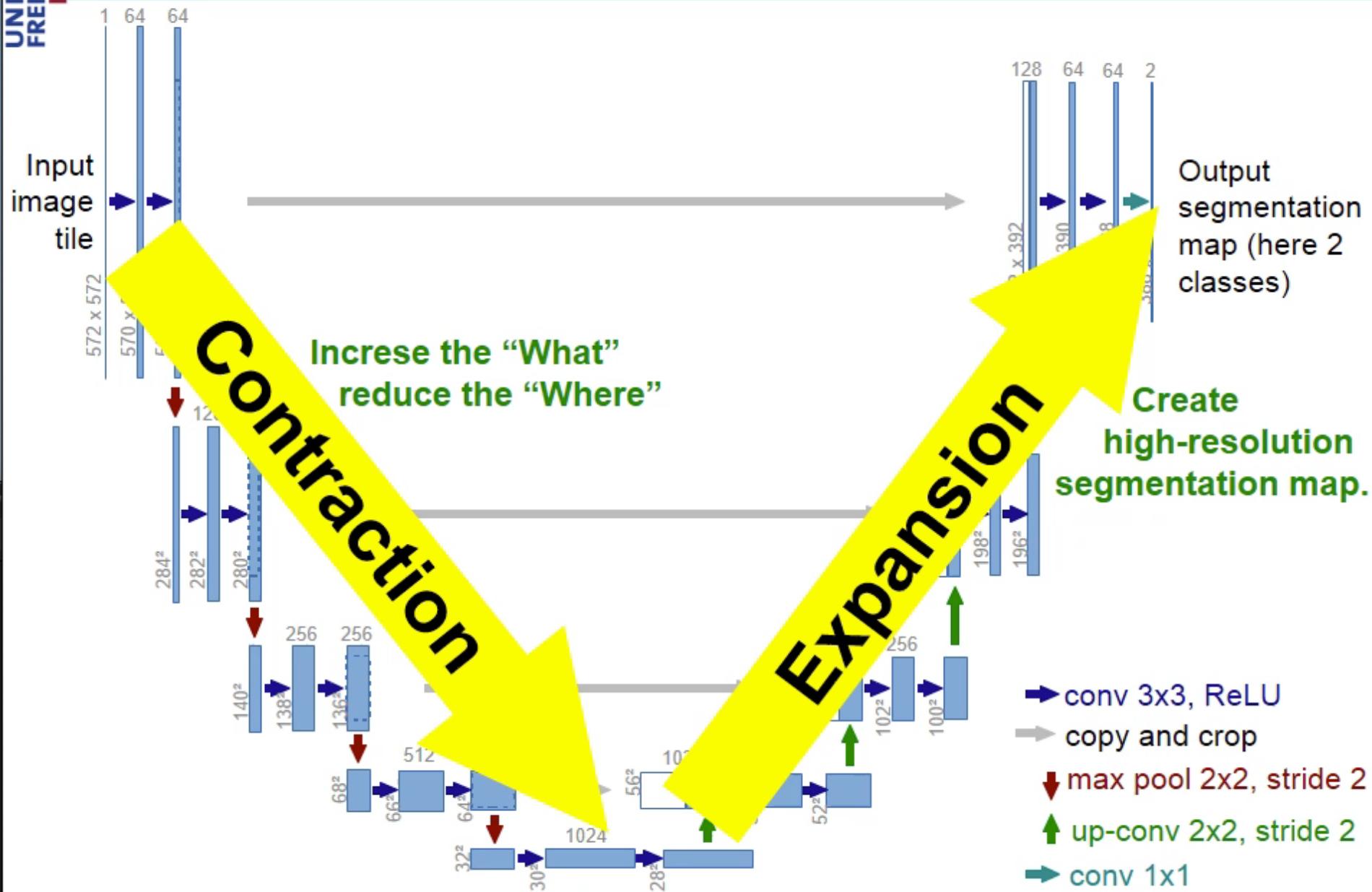
```
x = layers.Activation("relu")(x)  
x = layers.Conv2DTranspose(filters, 3, padding="same")(x)  
x = layers.BatchNormalization()(x)
```

```
x = layers.UpSampling2D(2)(x)
```

```
# Project residual  
residual = layers.UpSampling2D(2)(previous_block_activation)  
residual = layers.Conv2D(filters, 1, padding="same")(residual)  
x = layers.add([x, residual]) # Add back residual  
previous_block_activation = x # Set aside next residual
```

```
# Add a per-pixel classification layer  
outputs = layers.Conv2D(num_classes, 3, activation="softmax", padding="same")(x)  
  
# Define the model  
model = keras.Model(inputs, outputs)  
return model
```

# U-Net Architecture



# Segmentation

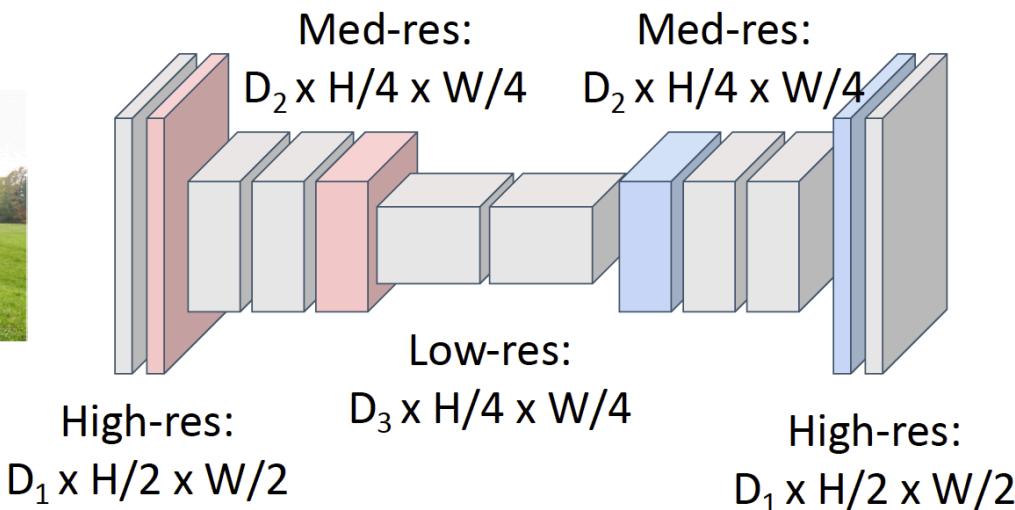
## Semantic Segmentation: Fully Convolutional Network

**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:**  
???



Predictions:  
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015  
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

<https://web.eecs.umich.edu/~justincj/teaching/eecs498/FA2020/schedule.html>

# In-Network Upsampling: “Unpooling”

## Bed of Nails

1	2
0	0
3	4



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

Input  
 $C \times 2 \times 2$

Output  
 $C \times 4 \times 4$

## Nearest Neighbor

1	2
3	4



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

Input  
 $C \times 2 \times 2$

Output  
 $C \times 4 \times 4$

## In-Network Upsampling: “Max Unpooling”

**Max Pooling:** Remember which position had the max

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



5	6
7	8

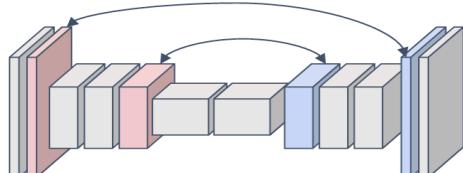
Rest  
of  
net



1	2
3	4



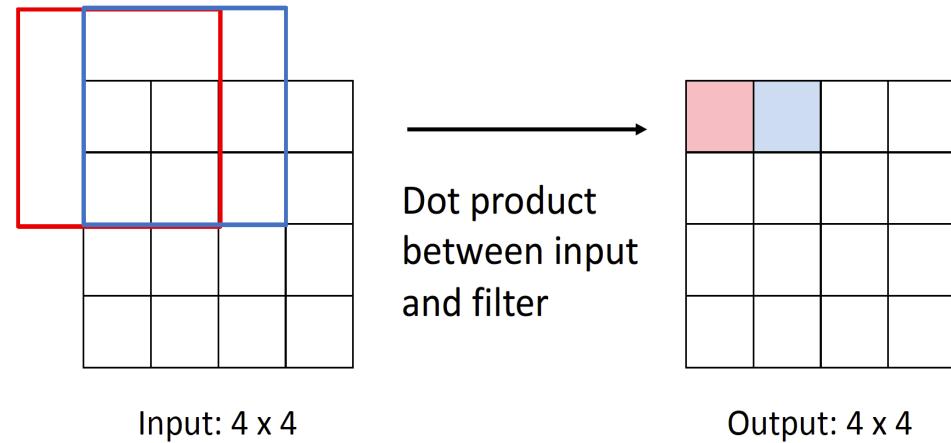
0	0	2	0
0	1	0	0
0	0	0	0
3	0	0	4



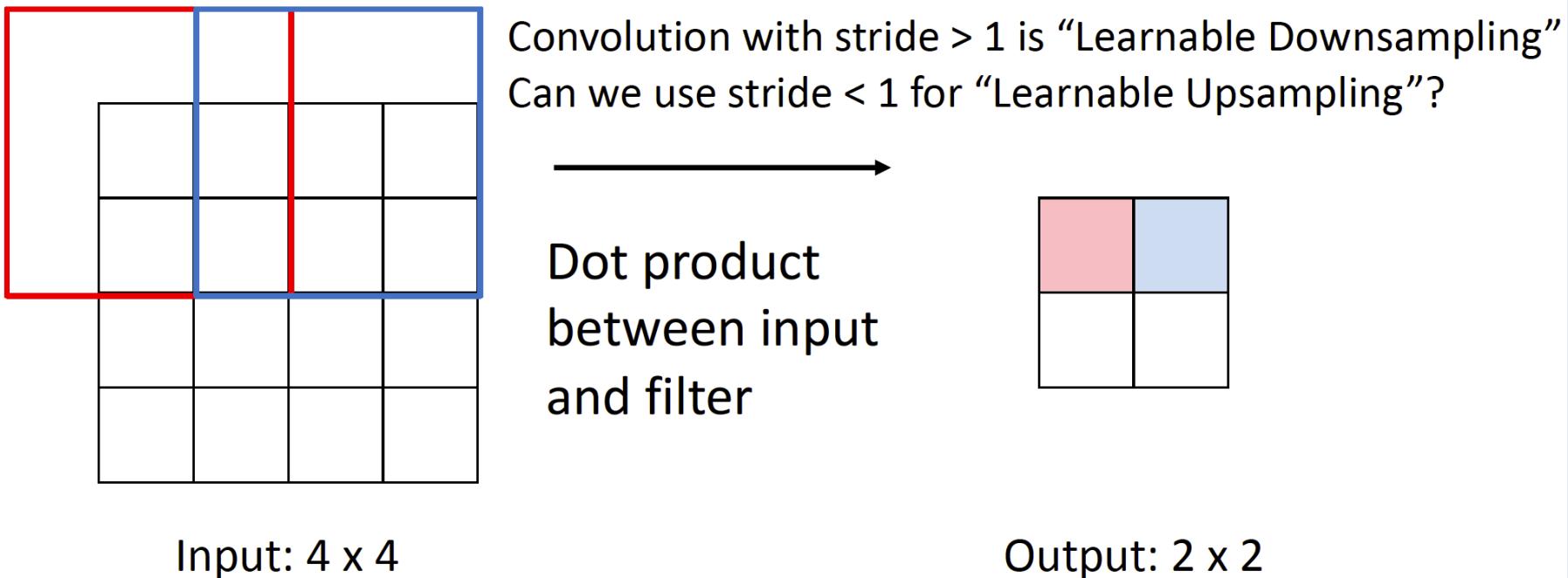
Pair each downsampling layer with an upsampling layer

## Learnable Upsampling: Transposed Convolution

**Recall:** Normal  $3 \times 3$  convolution, stride 1, pad 1



**Recall:** Normal  $3 \times 3$  convolution, stride 2, pad 1

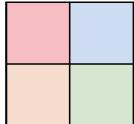


# MM in Segmentation Algorithm

## Learnable Upsampling: Transposed Convolution

3 x 3 convolution transpose, stride 2

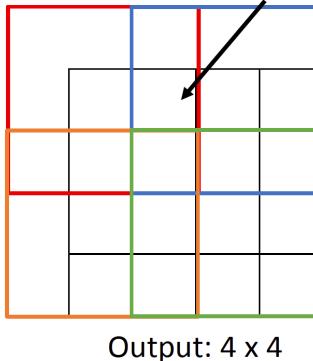
This gives 5x5 output – need to trim one pixel from top and left to give 4x4 output



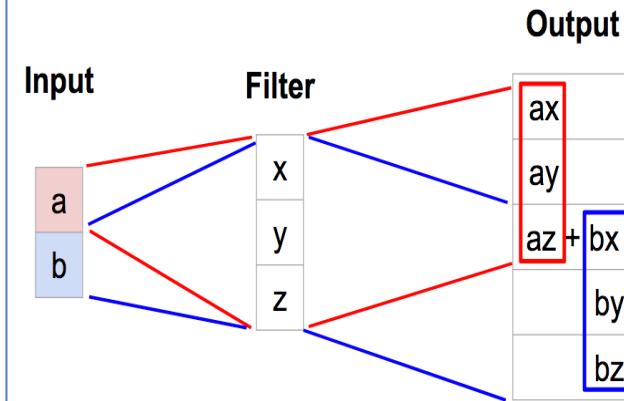
Weight filter by input value and copy to output

Input: 2 x 2

Sum where output overlaps



## Transpose Convolution: 1D Example



<https://www.tensorflow.org/tutorials/images/segmentation>

## Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & x & 0 & 0 & 0 \\ 0 & x & y & x & 0 & 0 \\ 0 & 0 & x & y & x & 0 \\ 0 & 0 & 0 & x & y & x \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ax + by + cz \\ ay + bz \\ bx + cy + dz \\ cx + dy \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=1, padding=1

Convolution transpose multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{bmatrix} x & 0 & 0 & 0 \\ y & x & 0 & 0 \\ z & y & x & 0 \\ 0 & z & y & x \\ 0 & 0 & z & y \\ 0 & 0 & 0 & z \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ax \\ ay + bx \\ az + by + cx \\ bz + cy + dx \\ cz + dy \\ dz \end{bmatrix}$$

When stride=1, convolution transpose is just a regular convolution (with different padding rules)

## Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ax + bz \\ ay + bz \\ bx + cy + dz \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=2, padding=1

Convolution transpose multiplies by the transpose of the same matrix:

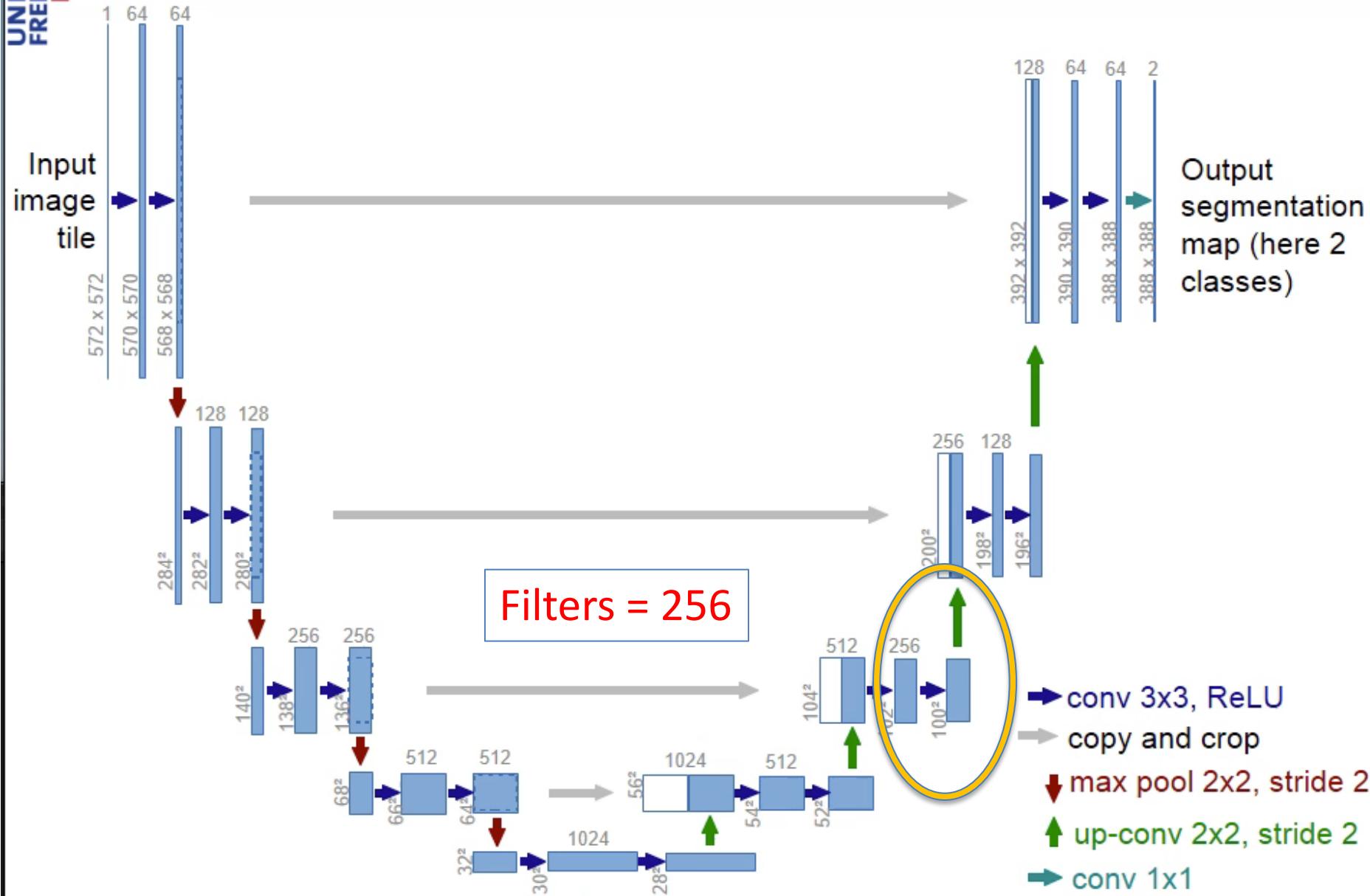
$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{bmatrix} x & 0 \\ y & 0 \\ z & x \\ 0 & y \\ 0 & z \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ax \\ ay \\ az + bx \\ by \\ bz \\ 0 \end{bmatrix}$$

When stride>1, convolution transpose is no longer a normal convolution!

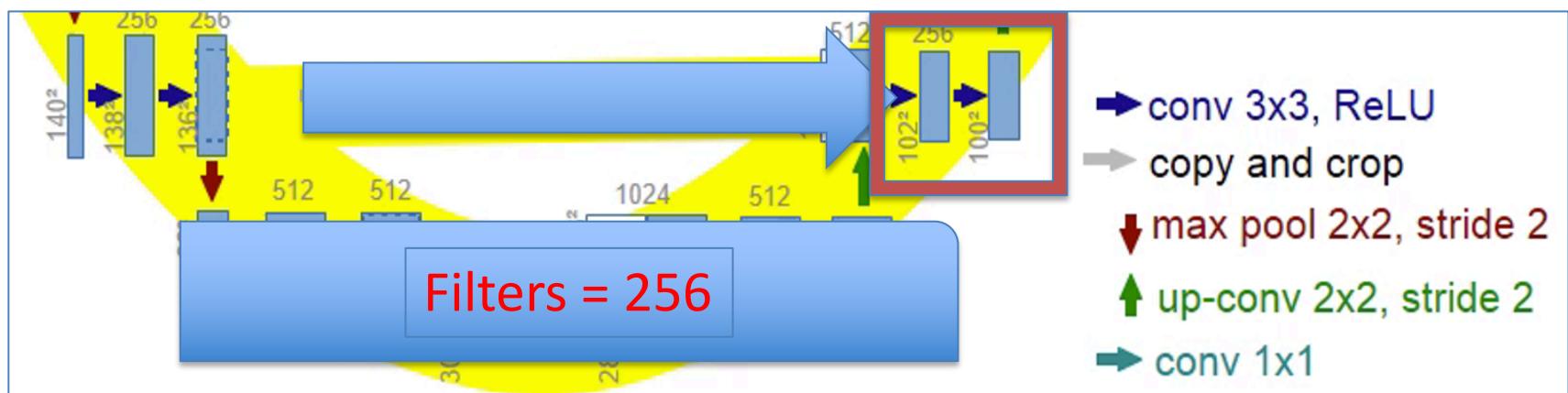
<https://web.eecs.umich.edu/~justincj/teaching/eecs498/FA2020/schedule.html>

# U-net Architecture

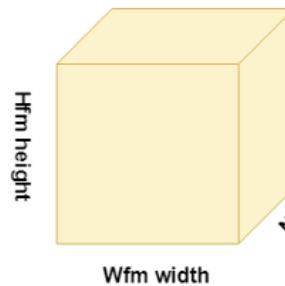


```
### [Second half of the network: upsampling inputs] ###
```

```
for filters in [256, 128, 64, 32]:  
    x = layers.Activation("relu")(x)  
    x = layers.Conv2DTranspose(filters, 3, padding="same")(x)  
    x = layers.BatchNormalization()(x)  
  
    x = layers.Activation("relu")(x)  
    x = layers.Conv2DTranspose(filters, 3, padding="same")(x)  
    x = layers.BatchNormalization()(x)  
  
    x = layers.UpSampling2D(2)(x)  
  
    # Project residual  
    residual = layers.UpSampling2D(2)(previous_block_activation)  
    residual = layers.Conv2D(filters, 1, padding="same")(residual)  
    x = layers.add([x, residual]) # Add back residual  
    previous_block_activation = x # Set aside next residual  
  
# Add a per-pixel classification layer  
outputs = layers.Conv2D(num_classes, 3, activation="softmax", padding="same")(x)  
  
# Define the model  
model = keras.Model(inputs, outputs)  
return model
```



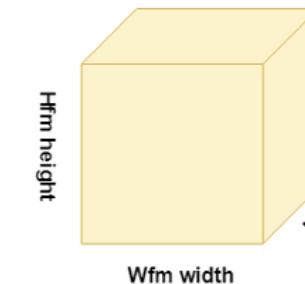
Input Channel = 256



N=256  
Size = (3x3 x 256)

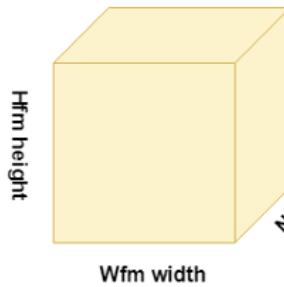


output Channel = 256



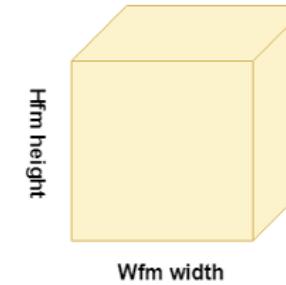
$$\text{Height} = 10 ; \text{width} = 10; N = 256. \Rightarrow \text{parameters} = (3 \times 3 \times 256) \times 256 + 256 = 590080$$

Input Channel = 256



N=256  
Size = (3x3 x 256)

output Channel = 256



$$\text{Height} = 10 ; \text{width} = 10; N = 256. \Rightarrow \text{parameters} = (3 \times 3 \times 256) \times 256 + 256 = 590080$$

```
### [Second half of the network: upsampling inputs] ###
```

```
for filters in [256, 128, 64, 32]:  
    x = layers.Activation("relu")(x)  
    x = layers.Conv2DTranspose(filters, 3, padding="same")(x)  
    x = layers.BatchNormalization()(x)  
  
    x = layers.Activation("relu")(x)  
    x = layers.Conv2DTranspose(filters, 3, padding="same")(x)  
    x = layers.BatchNormalization()(x)  
  
    x = layers.UpSampling2D(2)(x)
```

```
# Project residual  
residual = layers.UpSampling2D(2)(previous_block_activation)  
residual = layers.Conv2D(filters, 1, padding="same")(residual)  
x = layers.add([x, residual]) # Add back residual  
previous_block_activation = x # Set aside next residual
```

```
# Add a per-pixel classification layer  
outputs = layers.Conv2D(num_classes, 3, activation="softmax", padding="same")(x)  
  
# Define the model  
model = keras.Model(inputs, outputs)  
return model
```

Previous : 10x10x256, Upsampling2D(2)  
Reverse of max pooling of 2x2

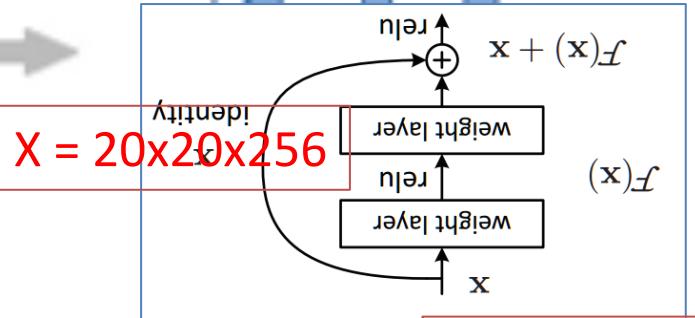
Input matrix = 20 x 20 x 256

N = 256

Channel = 256, Filter size = 1 x1, stride = 1, same padding

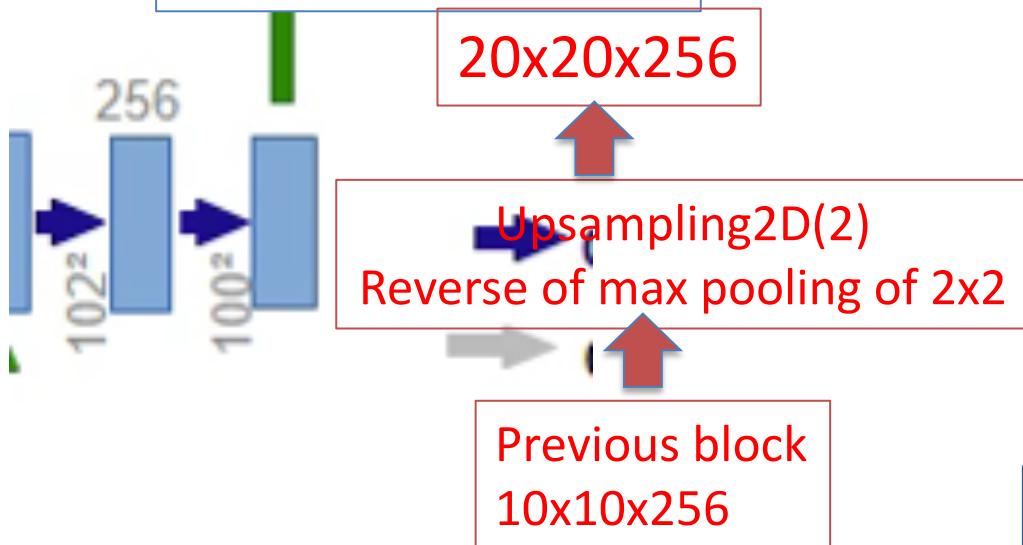
No of Filters= 256

```
residual = layers.Upsampling2D(2)(previous_block_activation)
residual = layers.Conv2D(256, 1, padding='same')(residual)
X = layers.add(x, residual)
previous_block_activation = x
```



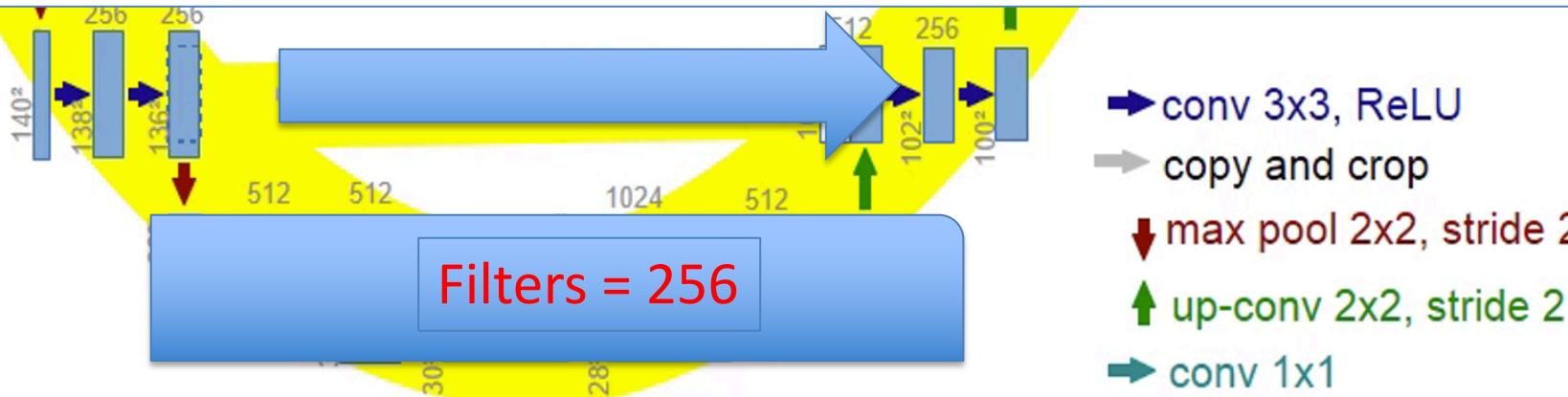
output matrix = 20 x 20 x 256

parameter,  $(1 \times 1 \times 256) \times 256 + 256$   
Total parameters = 65792



Add to the current X  
Size = 20x20x256

- conv 3x3, ReLU
- copy and crop
- ↓ max pool 2x2, stride 2
- ↑ up-conv 2x2, stride 2
- conv 1x1

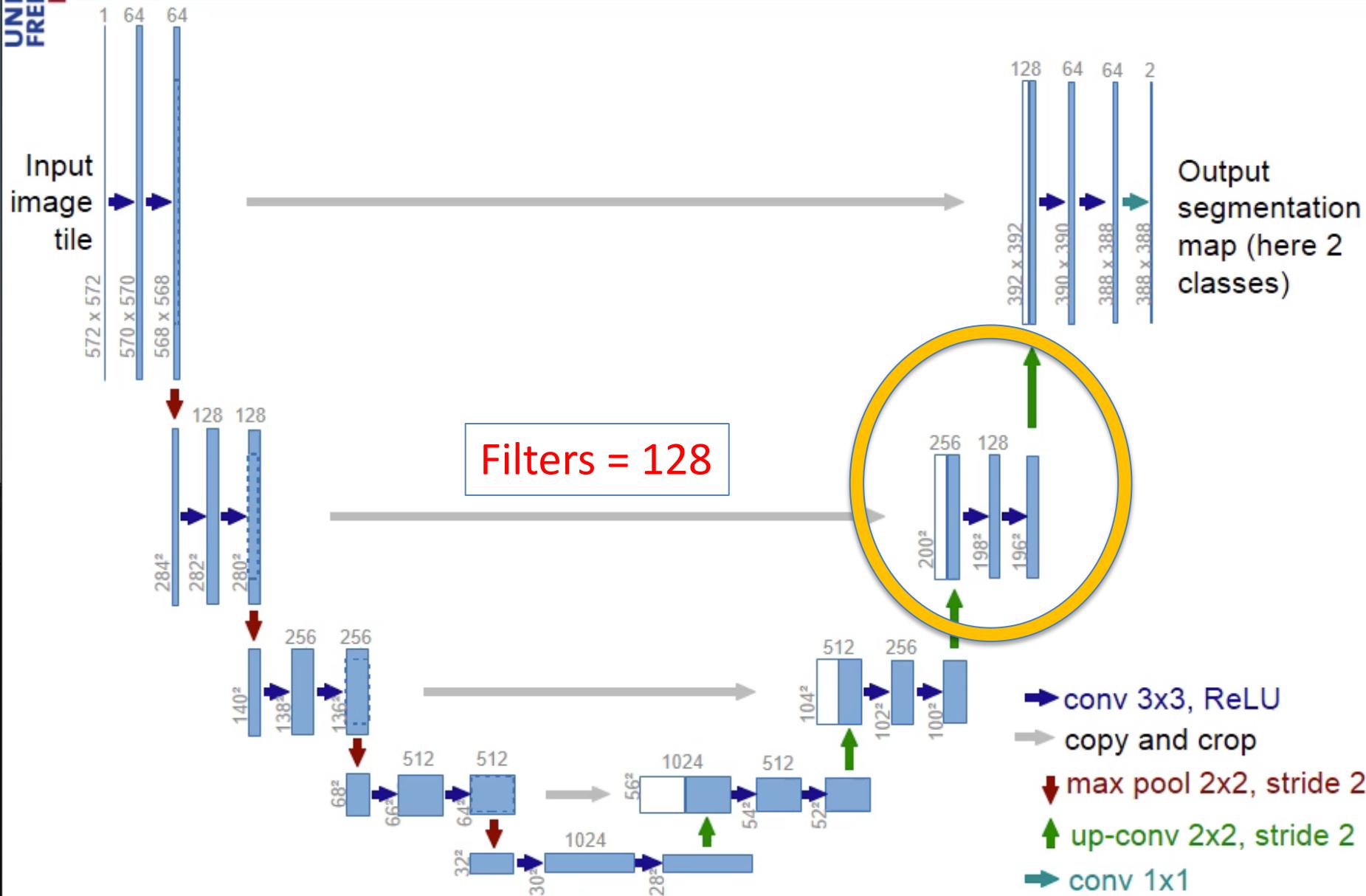


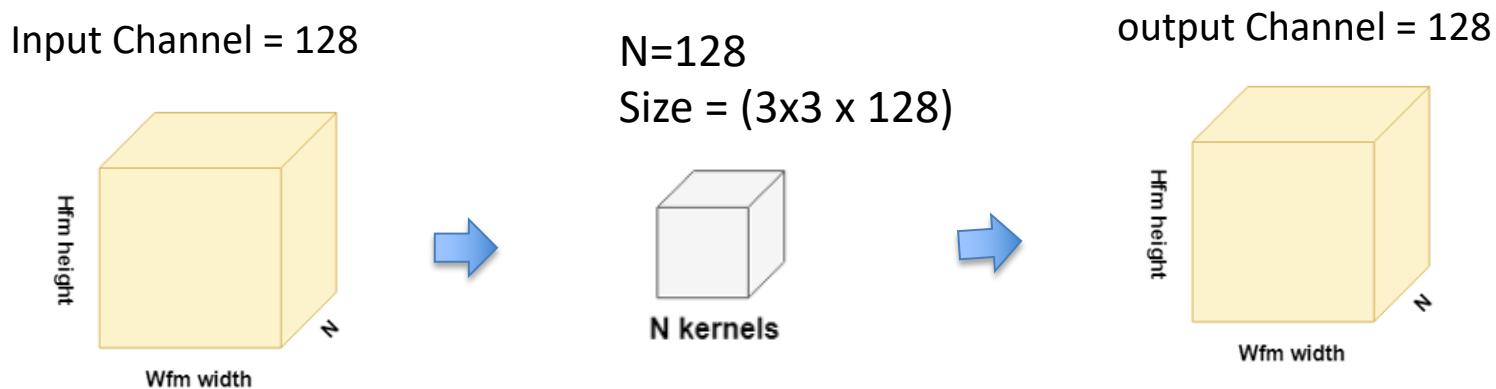
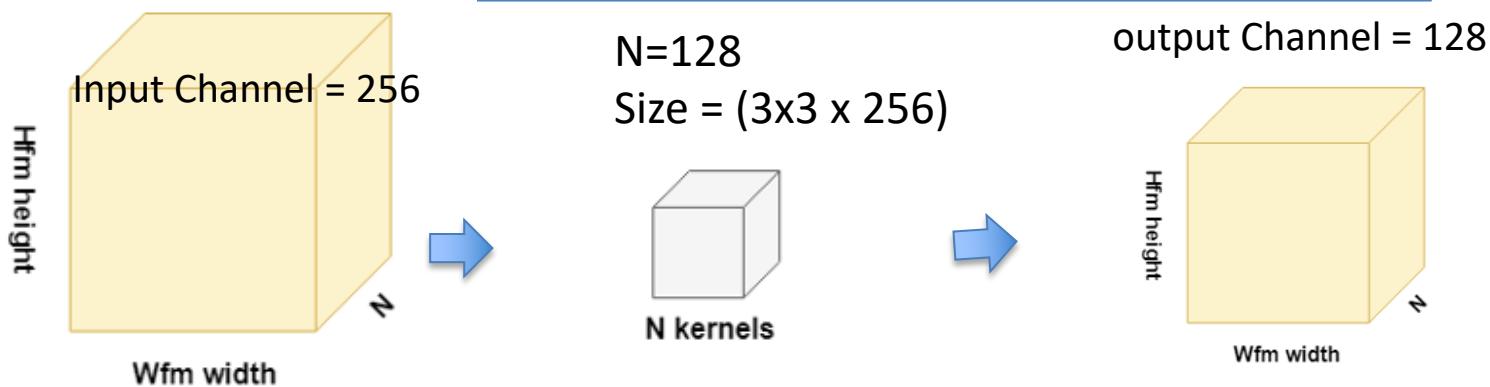
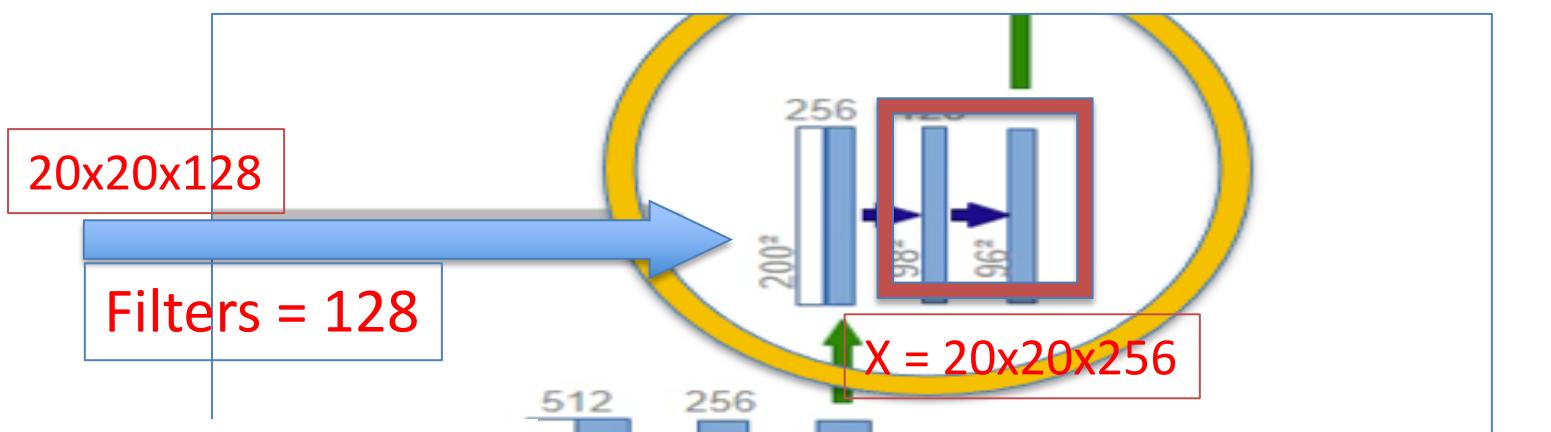
activation_7 (Activation)	(None, 10, 10, 256)	0	add_2[0][0]
conv2d_transpose (Conv2DTranspo)	(None, 10, 10, 256)	590080	activation_7[0][0]
batch_normalization_7 (BatchNor)	(None, 10, 10, 256)	1024	conv2d_transpose[0][0]
activation_8 (Activation)	(None, 10, 10, 256)	0	batch_normalization_7
conv2d_transpose_1 (Conv2DTrans)	(None, 10, 10, 256)	590080	activation_8[0][0]
batch_normalization_8 (BatchNor)	(None, 10, 10, 256)	1024	conv2d_transpose_1[0]
up_sampling2d_1 (UpSampling2D)	(None, 20, 20, 256)	0	add_2[0][0]
<b>Reverse Max Pooling =&gt; no parameters</b>			
up_sampling2d (UpSampling2D)	(None, 20, 20, 256)	0	batch_normalization_8
conv2d_4 (Conv2D)	(None, 20, 20, 256)	65792	up_sampling2d_1[0][0]
add_3 (Add)	(None, 20, 20, 256)	0	up_sampling2d[0][0] conv2d_4[0][0]

```
### [Second half of the network: upsampling inputs] ###
```

```
for filters in [256, 128, 64, 32]:  
    x = layers.Activation("relu")(x)  
    x = layers.Conv2DTranspose(filters, 3, padding="same")(x)  
    x = layers.BatchNormalization()(x)  
  
    x = layers.Activation("relu")(x)  
    x = layers.Conv2DTranspose(filters, 3, padding="same")(x)  
    x = layers.BatchNormalization()(x)  
  
    x = layers.UpSampling2D(2)(x)  
  
    # Project residual  
    residual = layers.UpSampling2D(2)(previous_block_activation)  
    residual = layers.Conv2D(filters, 1, padding="same")(residual)  
    x = layers.add([x, residual]) # Add back residual  
    previous_block_activation = x # Set aside next residual  
  
# Add a per-pixel classification layer  
outputs = layers.Conv2D(num_classes, 3, activation="softmax", padding="same")(x)  
  
# Define the model  
model = keras.Model(inputs, outputs)  
return model
```

# U-net Architecture





**Height = 20 ; width = 20; N = 128 => parameters =  $(3 \times 3 \times 128) \times 128 + 128 = 147587$**

```
### [Second half of the network: upsampling inputs] ###
```

```
for filters in [256, 128, 64, 32]:  
    x = layers.Activation("relu")(x)  
    x = layers.Conv2DTranspose(filters, 3, padding="same")(x)  
    x = layers.BatchNormalization()(x)  
  
    x = layers.Activation("relu")(x)  
    x = layers.Conv2DTranspose(filters, 3, padding="same")(x)  
    x = layers.BatchNormalization()(x)  
  
    x = layers.UpSampling2D(2)(x)
```

```
# Project residual  
residual = layers.UpSampling2D(2)(previous_block_activation)  
residual = layers.Conv2D(filters, 1, padding="same")(residual)  
x = layers.add([x, residual]) # Add back residual  
previous_block_activation = x # Set aside next residual
```

```
# Add a per-pixel classification layer  
outputs = layers.Conv2D(num_classes, 3, activation="softmax", padding="same")(x)  
  
# Define the model  
model = keras.Model(inputs, outputs)  
return model
```

Previous : 20x20x256, Upsampling2D(2)  
Reverse of max pooling of 2x2

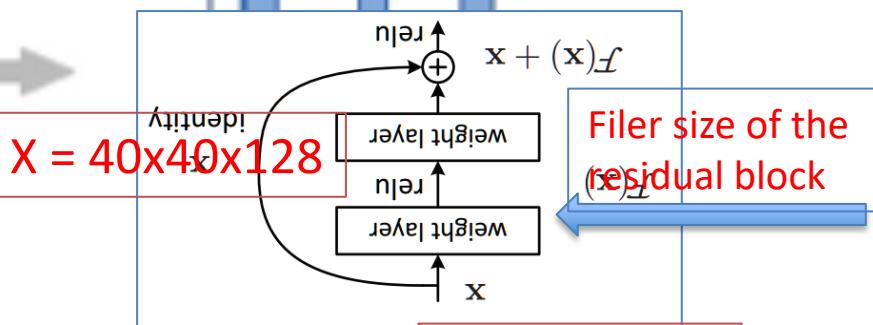
Input matrix = 40 x 40 x 256

N = 128

Channel = 256, Filter size = 1 x1, stride = 1, same padding

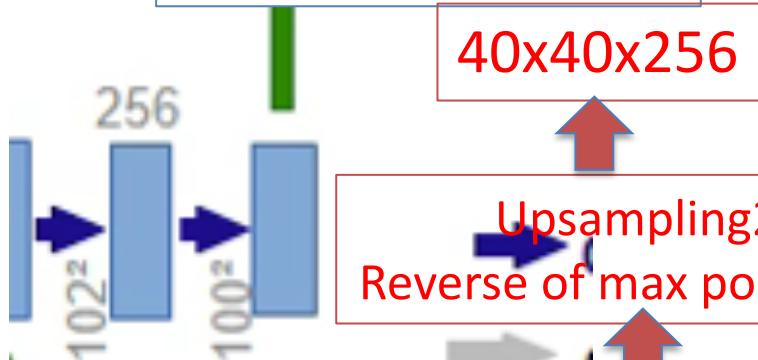
No of Filters = 128

residual = layers.Upsampling2D(2)(previous\_block\_activation)  
residual = layers.Conv2D(128, 1, padding='same')(residual)  
 $X = \text{layers.add}(x, \text{residual})$   
previous\_block\_activation = x



output matrix = 40 x 40 x 128

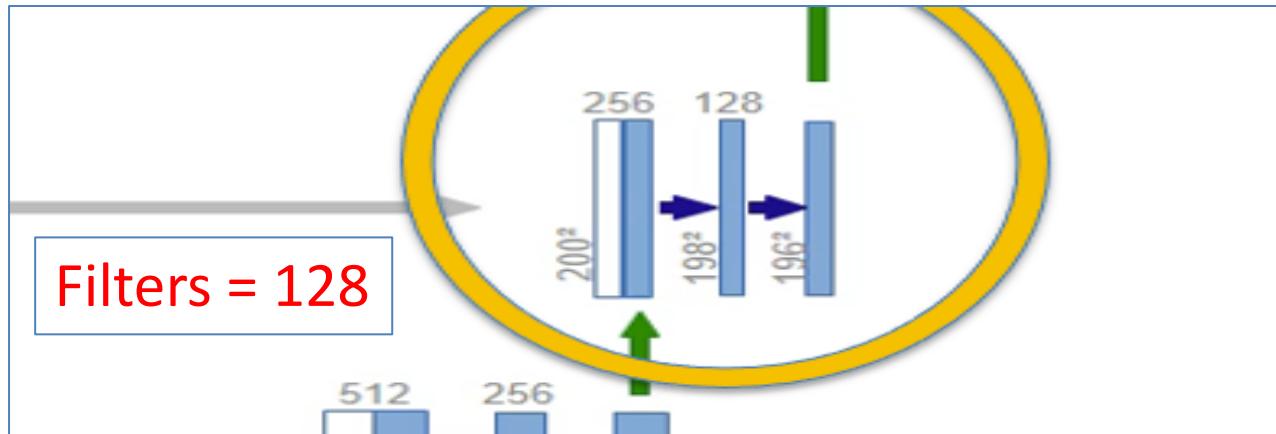
parameter,  $(1 \times 1 \times 256) \times 128 + 128$   
Total parameters = 32890



Previous block  
20x20x256

Add to the current X  
Size = 40x40x128

→ conv 3x3, ReLU  
→ copy and crop  
↓ max pool 2x2, stride 2  
↑ up-conv 2x2, stride 2  
→ conv 1x1

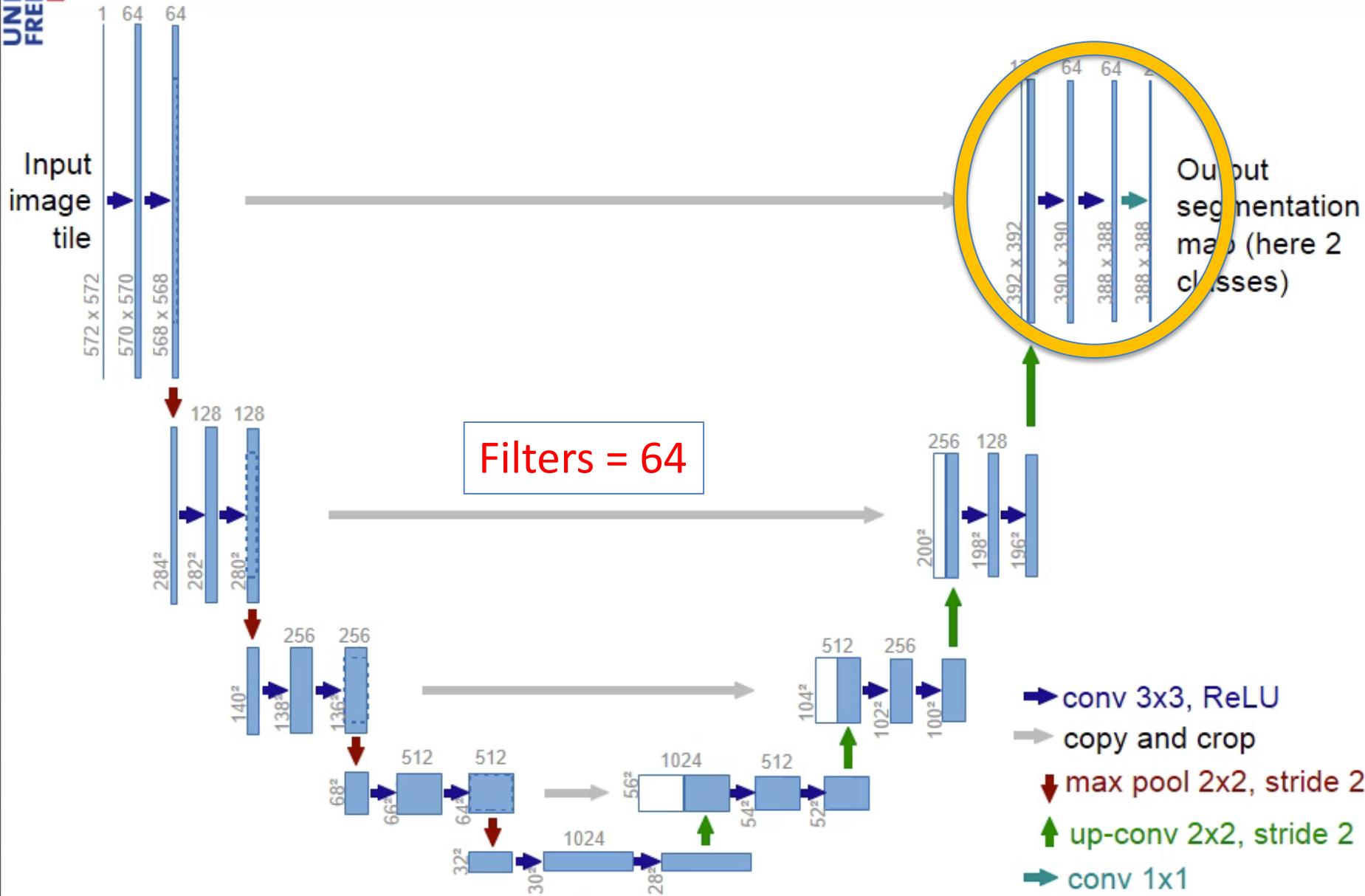


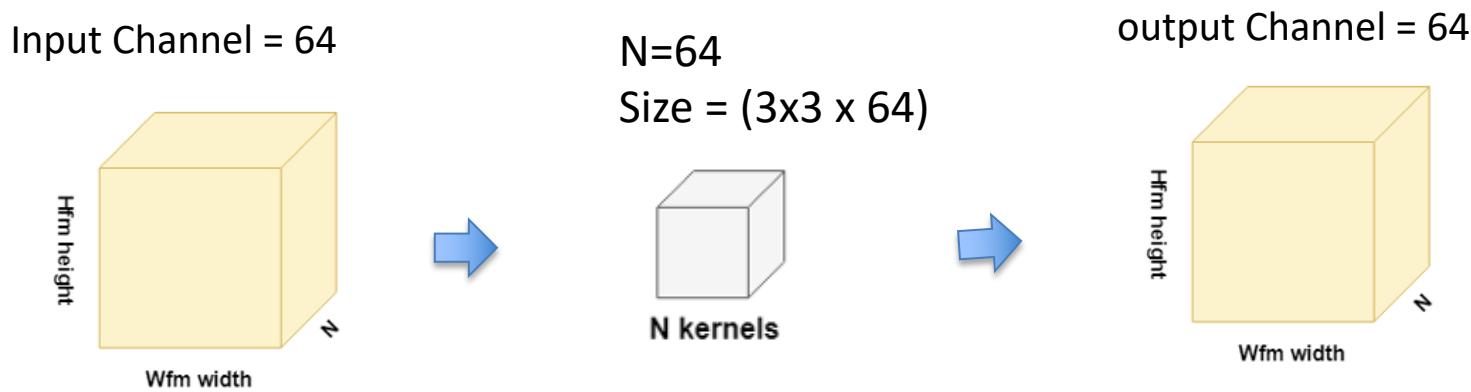
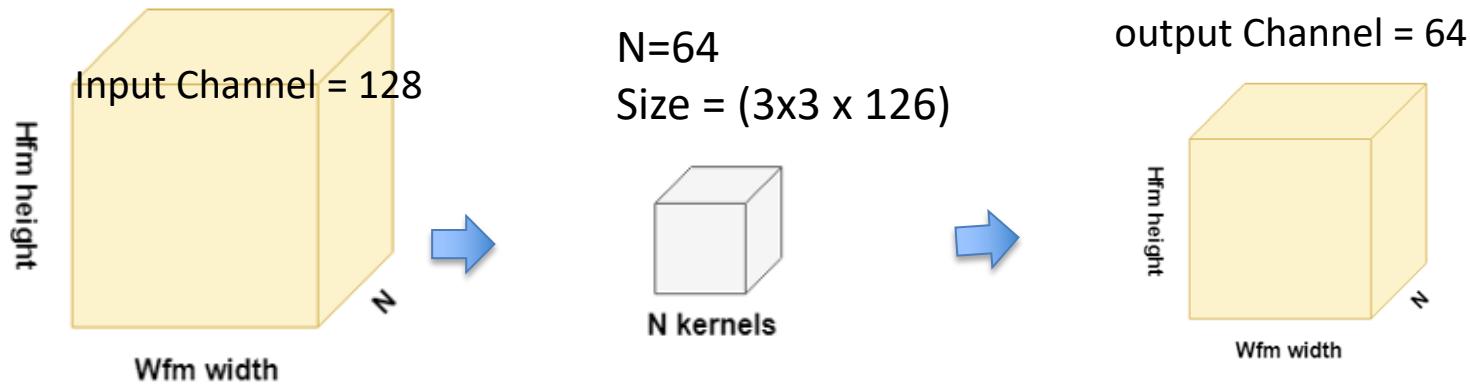
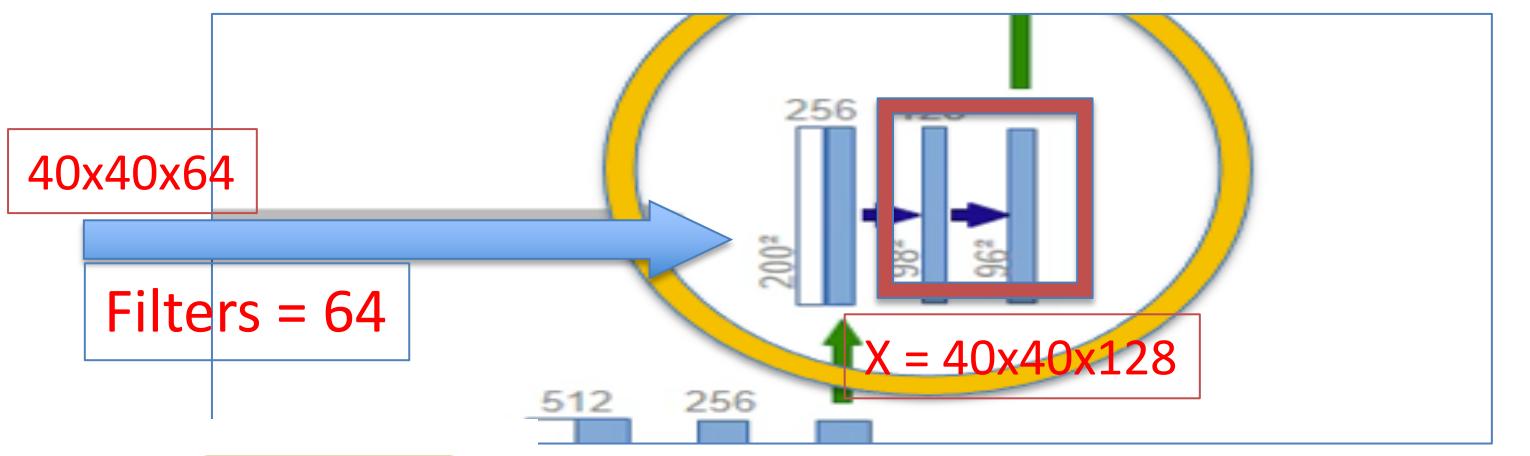
activation_9 (Activation)	(None, 20, 20, 256)	0	add_3[0] [0]
conv2d_transpose_2 (Conv2DTrans)	(None, 20, 20, 128)	295040	activation_9[0] [0]
batch_normalization_9 (BatchNor)	(None, 20, 20, 128)	512	conv2d_transpose_2[0]
activation_10 (Activation)	(None, 20, 20, 128)	0	batch_normalization_9
conv2d_transpose_3 (Conv2DTrans)	(None, 20, 20, 128)	147584	activation_10[0] [0]
batch_normalization_10 (BatchNo)	(None, 20, 20, 128)	512	conv2d_transpose_3[0]
up_sampling2d_3 (UpSampling2D)	(None, 40, 40, 256)	0	add_3[0] [0]
up_sampling2d_2 (UpSampling2D)	(None, 40, 40, 128)	0	batch_normalization_1
conv2d_5 (Conv2D)	(None, 40, 40, 128)	32896	up_sampling2d_3[0] [0]
add_4 (Add)	(None, 40, 40, 128)	0	up_sampling2d_2[0] [0]
			conv2d_5[0] [0]

```
### [Second half of the network: upsampling inputs] ###
```

```
for filters in [256, 128, 64, 32]:  
    x = layers.Activation("relu")(x)  
    x = layers.Conv2DTranspose(filters, 3, padding="same")(x)  
    x = layers.BatchNormalization()(x)  
  
    x = layers.Activation("relu")(x)  
    x = layers.Conv2DTranspose(filters, 3, padding="same")(x)  
    x = layers.BatchNormalization()(x)  
  
    x = layers.UpSampling2D(2)(x)  
  
    # Project residual  
    residual = layers.UpSampling2D(2)(previous_block_activation)  
    residual = layers.Conv2D(filters, 1, padding="same")(residual)  
    x = layers.add([x, residual]) # Add back residual  
    previous_block_activation = x # Set aside next residual  
  
# Add a per-pixel classification layer  
outputs = layers.Conv2D(num_classes, 3, activation="softmax", padding="same")(x)  
  
# Define the model  
model = keras.Model(inputs, outputs)  
return model
```

# U-net Architecture





```
### [Second half of the network: upsampling inputs] ###
```

```
for filters in [256, 128, 64, 32]:  
    x = layers.Activation("relu")(x)  
    x = layers.Conv2DTranspose(filters, 3, padding="same")(x)  
    x = layers.BatchNormalization()(x)  
  
    x = layers.Activation("relu")(x)  
    x = layers.Conv2DTranspose(filters, 3, padding="same")(x)  
    x = layers.BatchNormalization()(x)  
  
    x = layers.UpSampling2D(2)(x)
```

```
# Project residual  
residual = layers.UpSampling2D(2)(previous_block_activation)  
residual = layers.Conv2D(filters, 1, padding="same")(residual)  
x = layers.add([x, residual]) # Add back residual  
previous_block_activation = x # Set aside next residual
```

```
# Add a per-pixel classification layer  
outputs = layers.Conv2D(num_classes, 3, activation="softmax", padding="same")(x)  
  
# Define the model  
model = keras.Model(inputs, outputs)  
return model
```

Previous : 40x40x128, Upsampling2D(2)  
Reverse of max pooling of 2x2

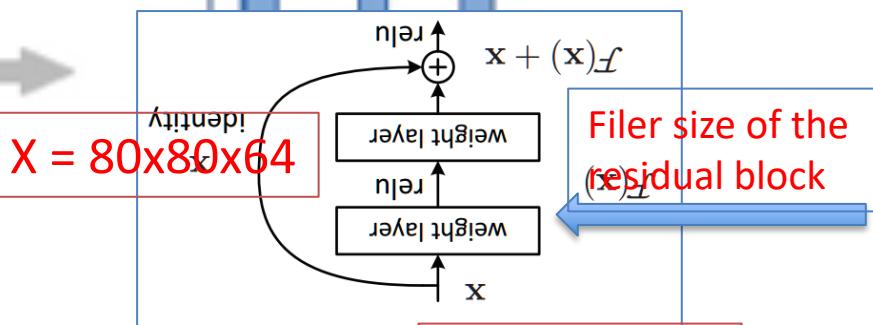
Input matrix = 80 x 80 x 128

N = 64

Channel = 256, Filter size = 1 x1, stride = 1, same padding

No of Filters = 64

residual = layers.Upsampling2D(2)(previous\_block\_activation)  
residual = layers.Conv2D(128, 1, padding='same')(residual)  
 $X = \text{layers.add}(x, \text{residual})$   
previous\_block\_activation = x



output matrix = 80 x 80 x 64

parameter,  $(1 \times 1 \times 128) \times 64 + 64$   
Total parameters = 8256

Add to the current X  
Size = 80x80x64

Upsampling2D(2)  
Reverse of max pooling of 2x2

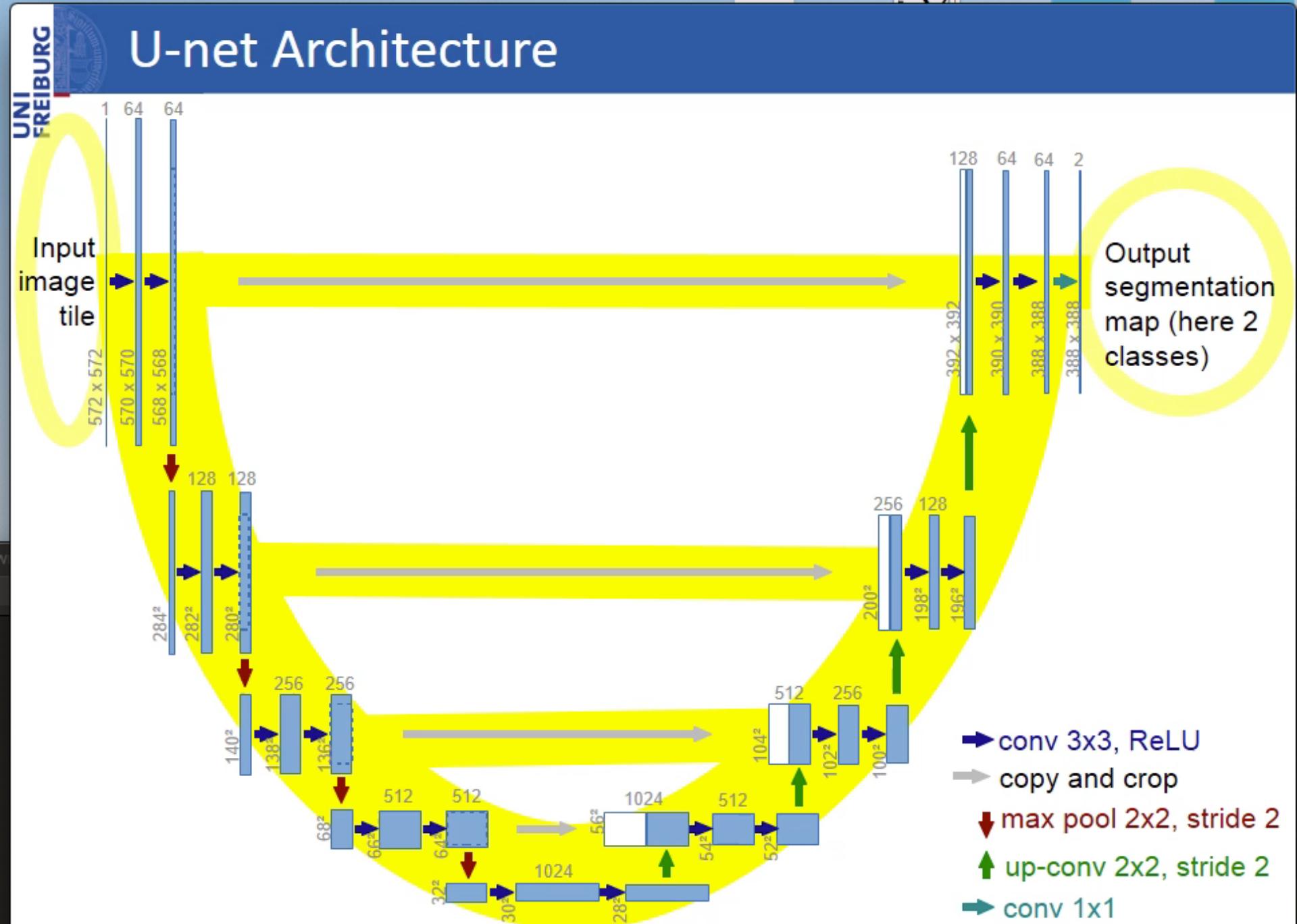
Previous block  
40x40x128

→ conv 3x3, ReLU  
→ copy and crop  
↓ max pool 2x2, stride 2  
↑ up-conv 2x2, stride 2

→ conv 1x1

activation_11 (Activation)	(None, 40, 40, 128) 0	add_4[0][0]	
conv2d_transpose_4 (Conv2DTrans)	(None, 40, 40, 64) 73792	activation_11[0][0]	
batch_normalization_11 (BatchNo)	(None, 40, 40, 64) 256	conv2d_transpose_4[0]	
activation_12 (Activation)	(None, 40, 40, 64) 0	batch_normalization_1	Total params: 2,058,979
conv2d_transpose_5 (Conv2DTrans)	(None, 40, 40, 64) 36928	activation_12[0][0]	Trainable params: 2,055,203
batch_normalization_12 (BatchNo)	(None, 40, 40, 64) 256	conv2d_transpose_5[0]	Non-trainable params: 3,776
up_sampling2d_5 (UpSampling2D)	(None, 80, 80, 128) 0	add_4[0][0]	
up_sampling2d_4 (UpSampling2D)	(None, 80, 80, 64) 0	batch_normalization_1	
conv2d_6 (Conv2D)	(None, 80, 80, 64) 8256	activation_13 (Activation)	(None, 80, 80, 64) 0 add_5[0][0]
add_5 (Add)	(None, 80, 80, 64) 0	conv2d_transpose_6 (Conv2DTrans)	(None, 80, 80, 32) 18464 activation_13[0][0]
<p><b>output = layers.Conv2D(num_class=3, 1, activation = "softmax", padding='same')(x)</b></p>		batch_normalization_13 (BatchNo)	(None, 80, 80, 32) 128 conv2d_transpose_6[0]
<p>Filer size of the final conv2D block</p>		activation_14 (Activation)	(None, 80, 80, 32) 0 batch_normalization_1
<p>parameter, (3x3x32) x 3 + 3 Total parameters = 8256</p>		conv2d_transpose_7 (Conv2DTrans)	(None, 80, 80, 32) 9248 activation_14[0][0]
		batch_normalization_14 (BatchNo)	(None, 80, 80, 32) 128 conv2d_transpose_7[0]
		up_sampling2d_7 (UpSampling2D)	(None, 160, 160, 64) 0 add_5[0][0]
		up_sampling2d_6 (UpSampling2D)	(None, 160, 160, 32) 0 batch_normalization_1
		conv2d_7 (Conv2D)	(None, 160, 160, 32) 2080 up_sampling2d_7[0][0]
		add_6 (Add)	(None, 160, 160, 32) 0 up_sampling2d_6[0][0] conv2d_7[0][0]
		conv2d_8 (Conv2D)	(None, 160, 160, 3) 867 add_6[0][0]

# U-net Architecture

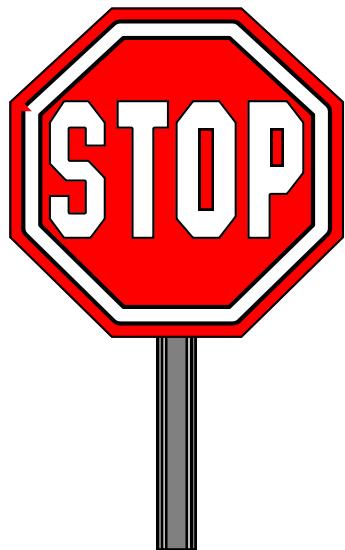


# Acknowledgements and References

This portion of the tutorial contains many extracted materials from many online websites and courses. Listed below are the major sites. This portion of the materials is not intended for public distribution. Please visit the sites websites for detail contents. If you are beginner users of DNN, I suggest to read the following list of websites in its order.

- 1) <http://neuralnetworksanddeeplearning.com>
- 2) <https://codelabs.developers.google.com/codelabs/cloud-tensorflow-mnist>
- 3) <https://www.deeplearning.ai/deep-learning-specialization/>
- 4) <http://cs231n.stanford.edu/>
- 5) MIT 6.S191, <https://www.youtube.com/watch?v=njKP3FqW3Sk>
- 6) <https://livebook.manning.com/book/deep-learning-with-python/about-this-book/>
- 7) <https://www.fast.ai/>
- 8) <https://petewarden.com/2015/04/20/why-gemm-is-at-the-heart-of-deep-learning/>
- 9) <https://www.nersc.gov/users/training/gpus-for-science/gpus-for-science-2020/>
- 10) <https://oneapi-src.github.io/oneDNN/>
- 11) <http://www.cs.cornell.edu/courses/cs4787/2020sp/>
- 12) More sites and free books are listed in [www.jics.utk.edu/actia](http://www.jics.utk.edu/actia) → ML
- 13) <https://machinelearningmastery.com>

# The End



- The End!

