

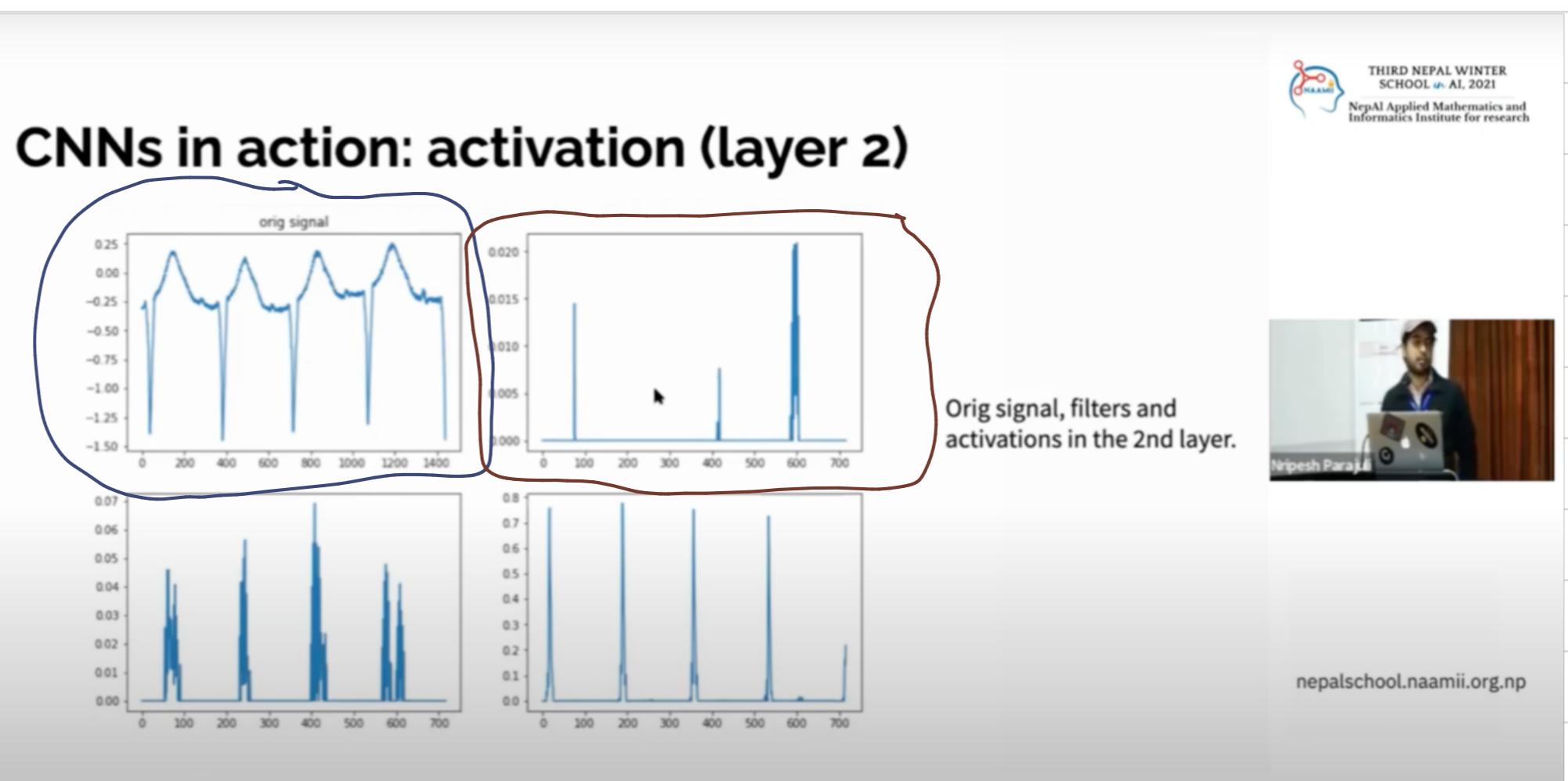
# Day 17, Oct -19, 2024.

Continuing CNNs .

① Top left::

Original Input Signal,  
having multiple peaks &  
troughs. Can be

Sensor Signal (time-Series) .



② Top Right: and Bottom Graphs. → illustrates the activations at different points within the CNN

Specifically in the Second Layer

→ Spikes Suggest shot specific features or patterns in the  
Original Signals are detected

→ Activations are sparse, CNNs offer filtering - certain patterns

(like edges in images or specific frequencies in image)

audio signals).

# Input Signals → Filter → Activation at different points → Output

## # Convolutional Neural Networks Scope

- Try CNNs for time-series Signal Although we have RNN for the time-series and sequence input.
- CNNs looks for the reliable and repeatable patterns in input data. Images are periodic & have temporal pattern for processing and extraction.
- CNNs are not meant to do Time-Series Input. They are for Image Processing and Computer Vision, Video processing, Speech and Audio processing, AR/VR, Generative Models.

## # The Convolutional Operation

$$s(t) \Rightarrow (x * w)(t) = \sum_{a=-\infty}^{\infty} x(a) w(t-a)$$

where  $x(t)$  Signal

$w(t) \rightarrow$  Kernel

$t \Rightarrow$  time | position of op

$a \Rightarrow$  position or time-shift  
applied to Kernel

## # Discrete Convolutional Formula

$$s(i, j) = (K * I)(i, j) \Rightarrow \sum_m \sum_n I(i+m, j+n) K(m, n)$$

useful for Specific to Images, which have finite Known Sizes

where  $I(i, j) \rightarrow$  input Image

$K(m, n) \rightarrow$  Convolutional Kernel

$m, n \rightarrow$  indices for Kernel  
 $I(i+m, j+n) \rightarrow$  Input Image pixel at  
position  $(i+m, j+n)$   $i, j$  are local pixels.

for i in range():

→ Outer loop iterates

Over the rows of  
the images Storing

Index=1 and ending h-2.

for j in image

→ Inner loop iterates over the columns of the images , Storing from index 1  
and ending iw - 2 , ensures enough Space for Kernel 3x3 left and right.

③ im-ij ⇒ - - - → Extracts 3x3 Subregion of the image centered (ij)  
Slice i:i+3 selects 3 rows and 3 columns.

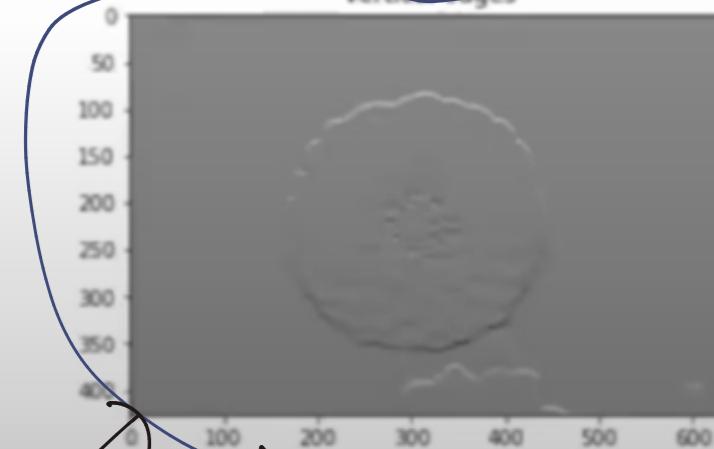
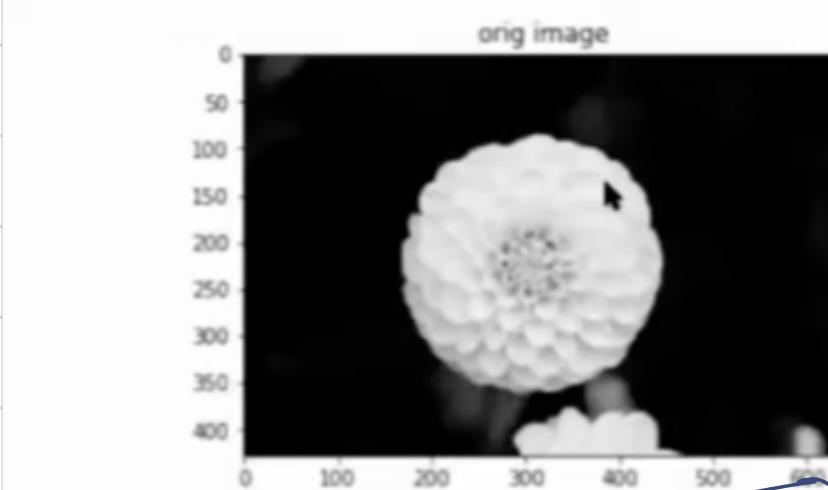
④ multiplies 3x3 Image element-wise with Kernel & Sum them (Dot product)



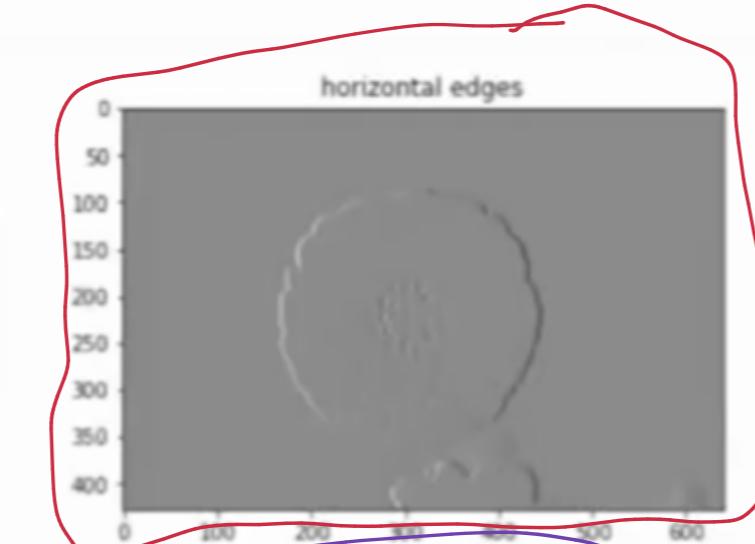
## The convolution operation: code

```
def conv_im_kernel(im, kernel):
    # hard-coded for 3 X 3 kernels
    k_h, k_w = kernel.shape
    im_h, im_w = im.shape
    new_im = np.zeros(im.shape)
    for i in range(1, im_h - 2):
        for j in range(1, im_w - 2):
            im_ij = im[i : i + 3, j : j + 3]
            new_im[i, j] = np.sum(np.multiply(im_ij, kernel))
    return new_im
```

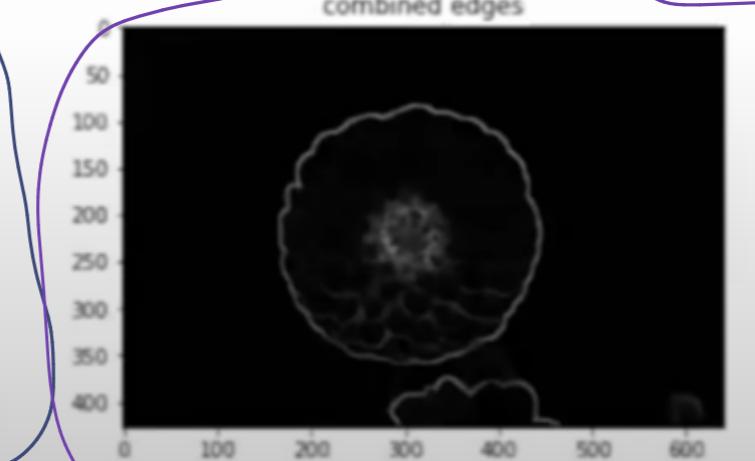
# Convolution example: edge detection



Horizontal Edges



→ vertical edges



vertical +  
horizontal edges

X



Combined edges

nepalschool.naamii.org.np

$$\begin{bmatrix} -1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & -1 \end{bmatrix}$$

Vertical

$$\begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix}$$

Horizontal

# Pooling : Is it differentiable | have Parameter? : + No.

# Stride : Number of pixels (or units) by which the filter or Kernel moves across the input image during the convolutional operation.

① Stride Value 1:  $\rightarrow$  filters moves one pixel at a time.

② Stride Value 2: filters moves 2 pixels at a time.

$$\text{O/p Size} \Rightarrow \left[ \frac{\text{Input Size} - \text{Kernel Size} + 2 \times \text{padding}}{\text{Stride}} \right] + 1$$

Large Strides  $\rightarrow$  Output Smaller feature Map.

Example:

Input Size  $\neq$  32

Kernel Size : 5

padding : 0

Strides : 2

$$\text{Op Side} = \left\lceil \frac{32 - 5 + 2 \times 0}{2} \right\rceil + 1$$

[Op Side 14]

Output Feature Map Size  $14 \times 14$ .

Translation Invariance:

$\rightarrow$  Ability of ANN to recognize features regardless of their exact position in the input.



• Purpose:  
Max Pooling

→ Reduce Spatial  
dimensions of the  
representation

→ Extract dominant  
features

→ Provide a form of translation invariance-

Example: Among

12	8
20	12

## Pooling (2)

12	8	0	30
20	12	0	8
23	40	0	21
108	32	22	12

2 X 2 max  
pooling

20	30
108	22

Final O/P using  
Max pooling

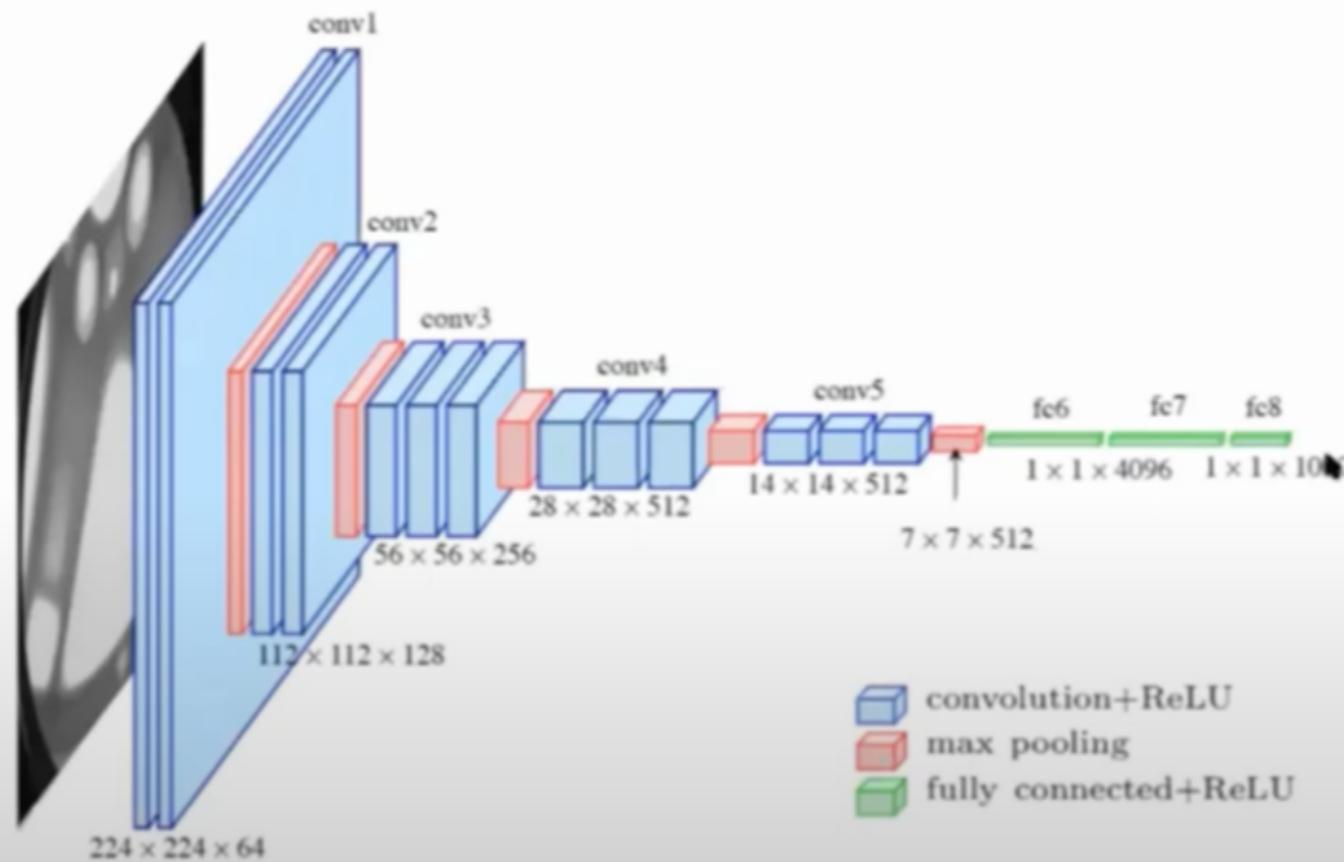
If image:

detect features regardless of their exact position -

→ 20 is Max, and Similarly for other 2x2  
this shows that max pooling allows to

# # Recap of CNN's Architecture.

## Understanding the network diagrams



This is a typical CNN with some FC layers towards the end:

- Input image -  $(224 \times 224 \times 3)$  ✓
- First layer is convolution + ReLU - padded to have same shape - 64 channels. ✓
- $(2 \times 2)$  maxpooling applied after this  $\rightarrow$  new size  $112 \times 112$  ✓
- Repeated 4 more times. ✓
- The 2D features are vectorized eventually. ✓
- 3 fully connected layers are added. ✓



nepalschool.naamii.org.np

Src: <https://www.researchgate.net/figure/Fig-A1-The-standard-VGG-16-network-architecture>

f. Decreasing Spatial Space and Increasing Features Map (Extract Rich Features)

ReLU [0, ∞)

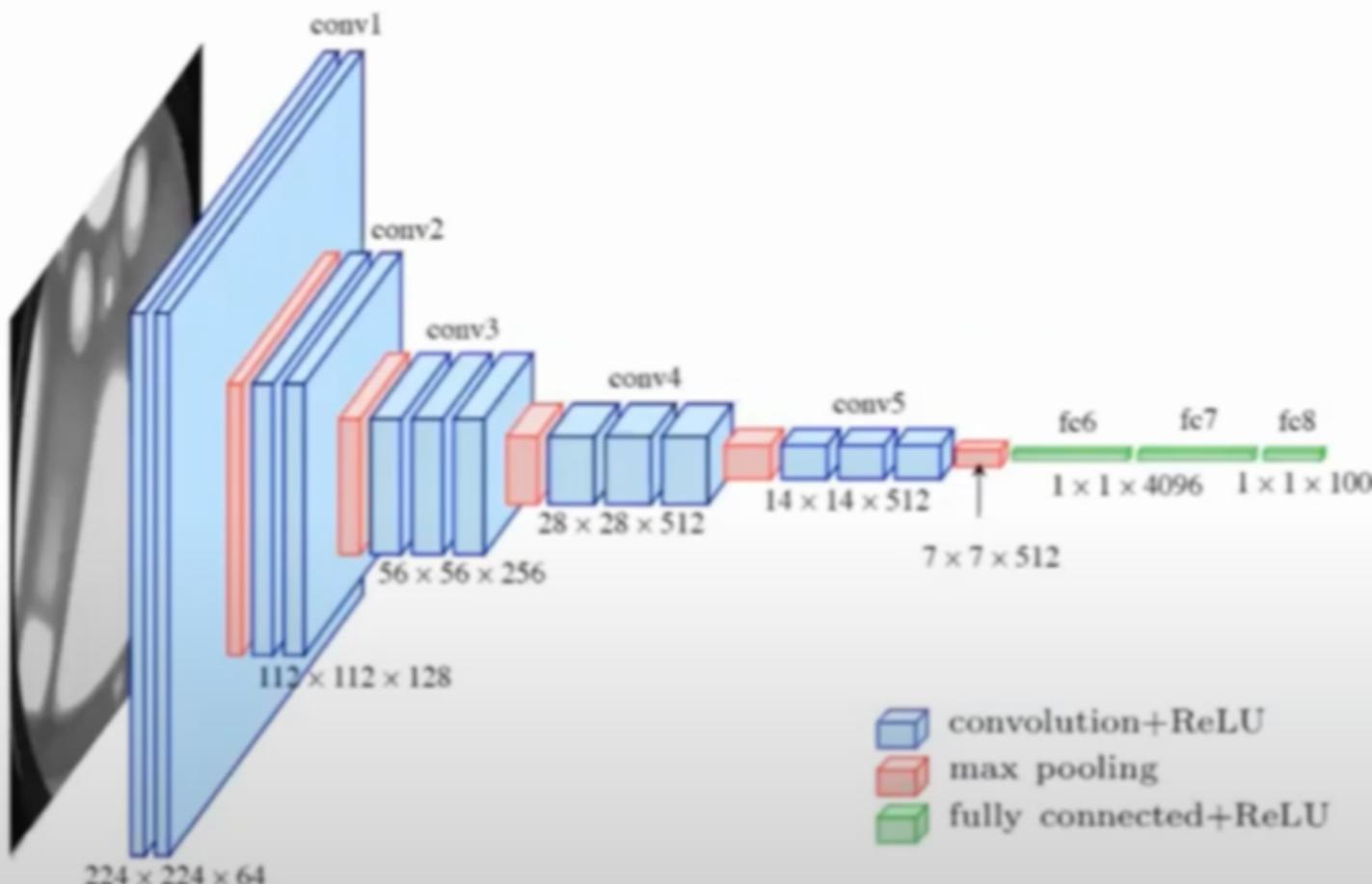
# Counting number of parameters

*Layer 1:*  
 64-filters  
 $(3 \times 3) \rightarrow \text{Kernel}$   
 $\times 3 \rightarrow \text{RGB}$   
 $\times 64 \rightarrow \text{bias}$



THIRD NEPAL WINTER  
SCHOOL on AI, 2021

NepAI Applied Mathematics and  
Informatics Institute for research



- Layer 1:
  - $64 \times (3 \times 3) \times 3 + 64 = 1792$
  - $64 \times (3 \times 3) \times 64 + 64 = 36,928$
- Layer 2:
  - $128 \times (3 \times 3) \times 64 + 128 = 73,856$
  - $128 \times (3 \times 3) \times 128 + 128 = 147,584$
- .....
- FC6  $\rightarrow 7 \times 7 \times 512 \times 1024 = 102,760,448$
- FC7  $\rightarrow 4096 \times 1000 = 25,088,000$
- Total  $\approx 138$  millions



nepalschool.naamii.org.np

Src: <https://www.researchgate.net/figure/Fig-A1-The-standard-VGG-16-network-architecture>

*Layer 2:*

①  $[128 \rightarrow \text{Filters} | \times (3 \times 3) \rightarrow \text{Kernel} | \times 64 \rightarrow \text{channel of } | + 128 \rightarrow \text{bias} |$   
 for features

*Layer 1*

Layer 1:

64 x (3x3) x 3 + 64  $\Rightarrow$  1792

Annotations:

- 64: No. of filters or O/p channels
- (3x3): Kernel size
- 3: RGB
- 64: bias's

64 x (3x3) x 64 + 64  $\Rightarrow$  36,928

Annotations:

- 64: No. of filters or O/p channels
- (3x3): Kernel size
- 64: Channel of layer 1 bias
- 64: bias's
- 64: Total parameters

Layer 2:

128 x (3x3) x 64 + 128  $\Rightarrow$  73,856 & Total parameters

Annotations:

- 128: No. of filters
- (3x3): Kernel Size
- 64: Input from previous O/p channel layers
- 128: bias's

# FC6 ( Fully Connected 6 ) :

$$7 \times 7 \times 512 \times 4096 = 102,768,448$$

$7 \times 7 \times 512 \rightarrow$  Size of Input feature Map

$4096 \Rightarrow$  no. of neurons in this Layer No Separate bias

term shown , but typically there would be 4096 bias term .

# FC7 ( Fully Connected 7 ) :



$1 \times 1 \times 4096$



$1 \times 1 \times 1000$

$$4096 \times 1000 \Rightarrow 4,096,000$$

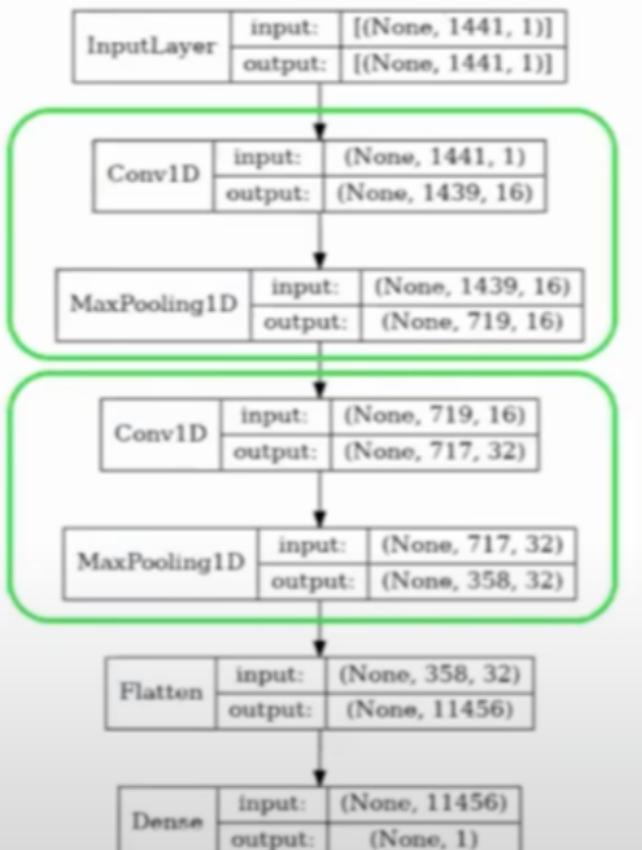
\*  $4096 \rightarrow$  no. of input neurons from FC6

\*  $1000 \rightarrow$  no. of output neurons ( would 1000 bias terms ) .

Total  $\rightarrow \sim 138$  millions ( Approx ) .

# Coding Implementation Part:

## Implementation: what does it look like?



A simple  
CNN in code

```
from tensorflow import keras
from tensorflow.keras import layers

#create model
cnn_model = keras.Sequential(
    [
        keras.Input(shape=X_train.shape[1], 1),
        layers.Conv1D(16, kernel_size=3, activation="relu"),
        layers.MaxPool1D(2),
        layers.Conv1D(32, kernel_size=3, activation="relu"),
        layers.MaxPool1D(2),
        layers.Flatten(),
        layers.Dense(
            1, activation="sigmoid",
        ),
    ]
)

cnn_model.compile(
    optimizer='adam',
    loss='binary_crossentropy',
    metrics=[
        metrics.BinaryAccuracy(),
    ],
)

cnn_model.fit(
    X_train_exp, y_train,
    validation_data=(X_test_exp, y_test),
    epochs=8,
    callbacks=[model_checkpoint_callback],
)
```



nepalschool.naamii.org.np

Hyperparameters are not in the training dataset (e.g: number of layers, batch size, K in KNN)

