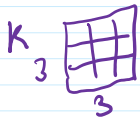Convolution Layer :-

Accepts a volume of size $W1 \times H1 \times D1$

K₃ [grid]
   3

$k \to$ # of filters
$S \to$ stride                    $F \to$ spacial extent
$P \to$ Amount of padding

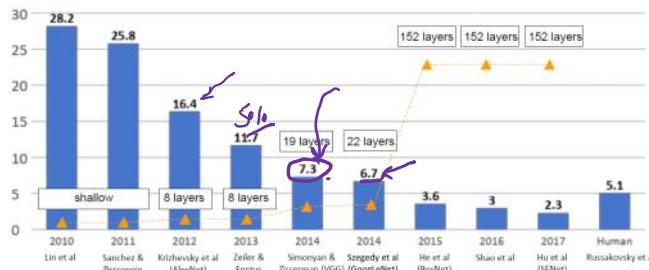Produces an output of size $W2 \times H2 \times D2$
                                                            $\underset{k}{}$

$$W2 = \frac{W1 - F + 2P}{S} + 1$$

$$H2 = \frac{H1 - F + 2P}{S} + 1$$

Weights per filter :-        $F \times F \times D1$

Image Net Dataset (ILSVRC)



ReLU  2f
  ↑    Net
 60   61.5m
  M   parameters    $227 \times 227 \times 3$
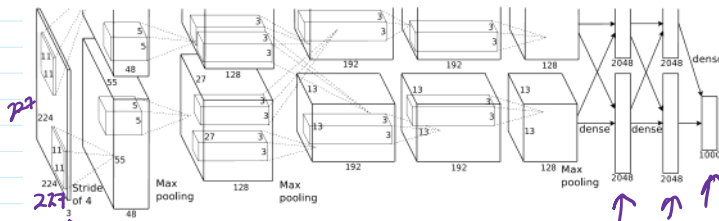
ALEXNET :-



Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

1000

$$\frac{227 - 11 + 2 \times 0}{4} + 1$$

$(227 \times 227 \times 3) \to$ Input
$(55 \times 55 \times 96) \to$ CONV 1 $\to$ 96 11×11 filters at stride 4, pad 0
$(27 \times 27 \times 96) \to$ MAXPOOL1  3×3 filters at stride 2
$(\quad '' \quad) \to$ NORM 1 $\leftarrow$ Not common anymore
$(27 \times 27 \times 256) \to$ CONV 2 $\to$ 256  5×5 filters stride 1, pad 2

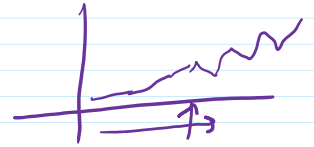$(13 \times 13 \times 256) \rightarrow$ MAXPOOL 2 $\rightarrow$ 3x3 stride 2

$(\quad " \quad) \rightarrow$ NORM2

$(13 \times 13 \times 384) \rightarrow$ CONV3 $\rightarrow$ 384 3x3 filters, stride 1, pad1

$(13 \times 13 \times 384) \rightarrow$ CONV4 $\rightarrow$ 384 3x3 filters, stride1, pad 1

$(13 \times 13 \times 256) \rightarrow$ CONV5 $\rightarrow$ 256 3x3 filters stride 1, pad1

$(6 \times 6 \times 256) \rightarrow$ MAXPOOL5 $\rightarrow$ 3x3 filters stride 2

$(6 \times 6 \times 256 \times 4096)$ $(4096) \rightarrow$ FC 6 $\rightarrow$ 4096 neurons

$(4096) \rightarrow$ FC 7 $\rightarrow$ 4096 neurons

$(1000) \rightarrow$ FC 8 $\rightarrow$ 1000 neurons $\leftarrow$ Softmax

| L1 | L2 | L3 | L4 | L5 | L6 | L7 | L8 | |
|----|----|----|----|----|----|----|----|----|
| 35k | 307k | 884k | 663k | 442k | 37M | 16M | 4M | $\rightarrow$ 60M |

## Details:—

① First use of ReLU

② Heavy data augmentation

③ Dropout 0.5

④ Batch size 128

⑤ SGD Momentum 0.9.

⑥ LR, 1e-2, Reduced manually where val. accuracy had plateus.

VGG Net :— (Visual Geometry Group) $\rightarrow$ Oxford University



AlexNet        VGG16        VGG19

```
INPUT: [224x224x3]      memory: 224*224*3=150K  params: 0      (not counting biases)
CONV3-64: [224x224x64]  memory: 224*224*64=3.2M  params: (3*3*3)*64 = 1,728
CONV3-64: [224x224x64]  memory: 224*224*64=3.2M  params: (3*3*64)*64 = 36,864
POOL2: [112x112x64]     memory: 112*112*64=800K  params: 0
CONV3-128: [112x112x128]  memory: 112*112*128=1.6M  params: (3*3*64)*128 = 73,728
CONV3-128: [112x112x128]  memory: 112*112*128=1.6M  params: (3*3*128)*128 = 147,456
POOL2: [56x56x128]      memory: 56*56*128=400K  params: 0
CONV3-256: [56x56x256]  memory: 56*56*256=800K  params: (3*3*128)*256 = 294,912
CONV3-256: [56x56x256]  memory: 56*56*256=800K  params: (3*3*256)*256 = 589,824
CONV3-256: [56x56x256]  memory: 56*56*256=800K  params: (3*3*256)*256 = 589,824
POOL2: [28x28x256]      memory: 28*28*256=200K  params: 0
CONV3-512: [28x28x512]  memory: 28*28*512=400K  params: (3*3*256)*512 = 1,179,648
CONV3-512: [28x28x512]  memory: 28*28*512=400K  params: (3*3*512)*512 = 2,359,296
CONV3-512: [28x28x512]  memory: 28*28*512=400K  params: (3*3*512)*512 = 2,359,296
POOL2: [14x14x512]      memory: 14*14*512=100K  params: 0
CONV3-512: [14x14x512]  memory: 14*14*512=100K  params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512]  memory: 14*14*512=100K  params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512]  memory: 14*14*512=100K  params: (3*3*512)*512 = 2,359,296
POOL2: [7x7x512]        memory: 7*7*512=25K  params: 0
FC: [1x1x4096]          memory: 4096  params: 7*7*512*4096 = 102,760,448
FC: [1x1x4096]          memory: 4096  params: 4096*4096 = 16,777,216
FC: [1x1x1000]          memory: 1000  params: 4096*1000 = 4,096,000
```

initial layers {

Later {

CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296
POOL2: [7x7x512] memory: 7*7*512=25K params: 0
FC: [1x1x4096] memory: 4096 params: 7*7*512*4096 = **102,760,448**
FC: [1x1x4096] memory: 4096 params: 4096*4096 = 16,777,216
FC: [1x1x1000] memory: 1000 params: 4096*1000 = 4,096,000

TOTAL memory: 24M * 4 bytes ~= 96MB / image (only forward! ~*2 for bwd)
TOTAL params: 138M parameters

Later