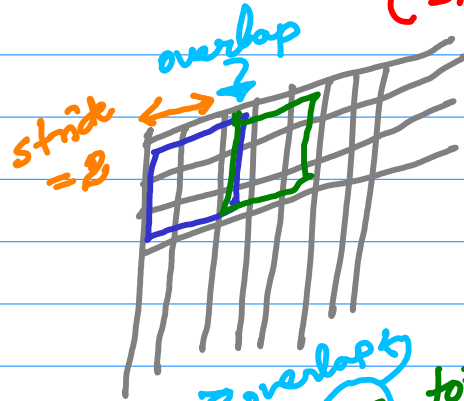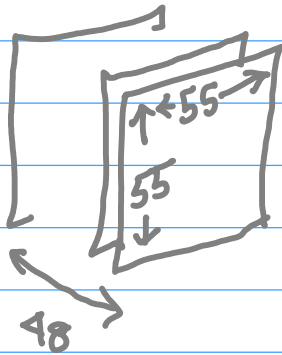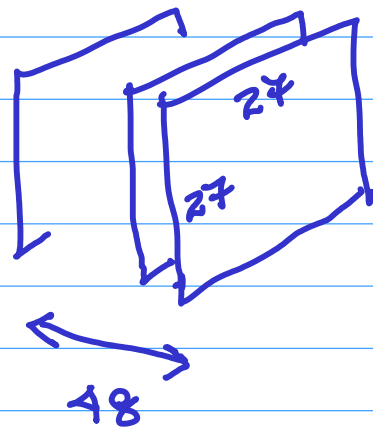# AlexNet (contd.)

② Second layer: 3x3 overlapping Max Pooling

(stride = 2)



(#0 to #2)
#0

(#2 to #4)
#1

(#4 to #6)
#2

... (#2k to #2k+2) ...
#k

(#52 to #54)

$$2k = 52 \Rightarrow k = 26$$

→ There are 27x27 outputs
48 of these, 2 groups



③ 3rd layer: Local Response Normalisation

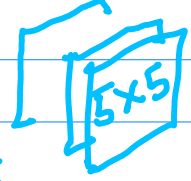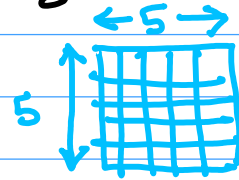→ only the values change, the size does not change.

# ④ Second Convolutional Layer
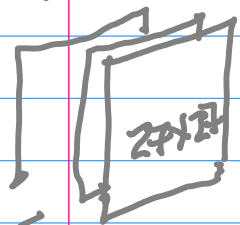
## 2 groups of 128 kernels of size 5×5×48

### stride = 1, pad = 2

(given)

Input:



27×27
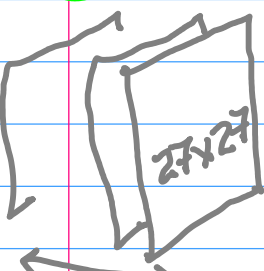
(48) ←————— one-to-one link ————→ (48)
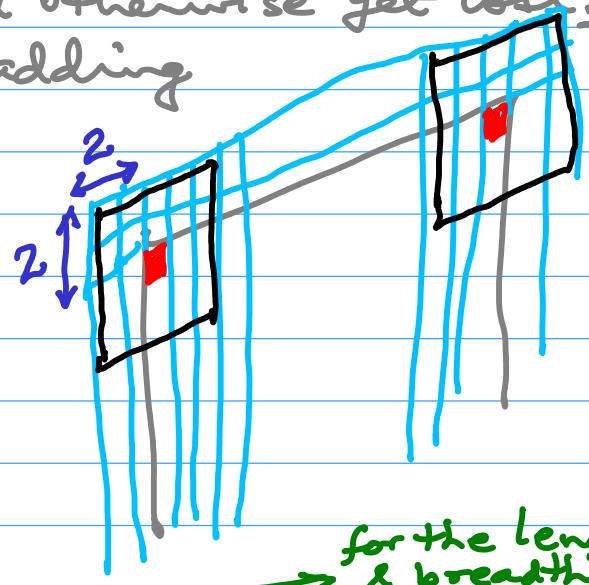
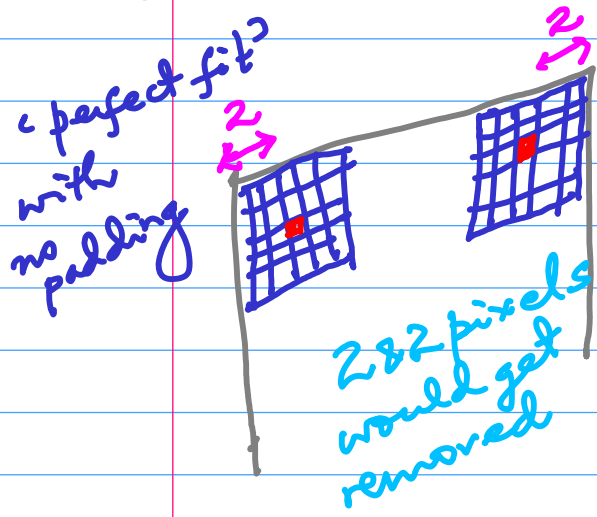5×5

27×27

(48) ←——→ (48)

5×5

### stride = 1

⇒ Result can get smaller with a 'perfect fit' convolution, but here

### pad = 2 → doesn't in this case, since the number

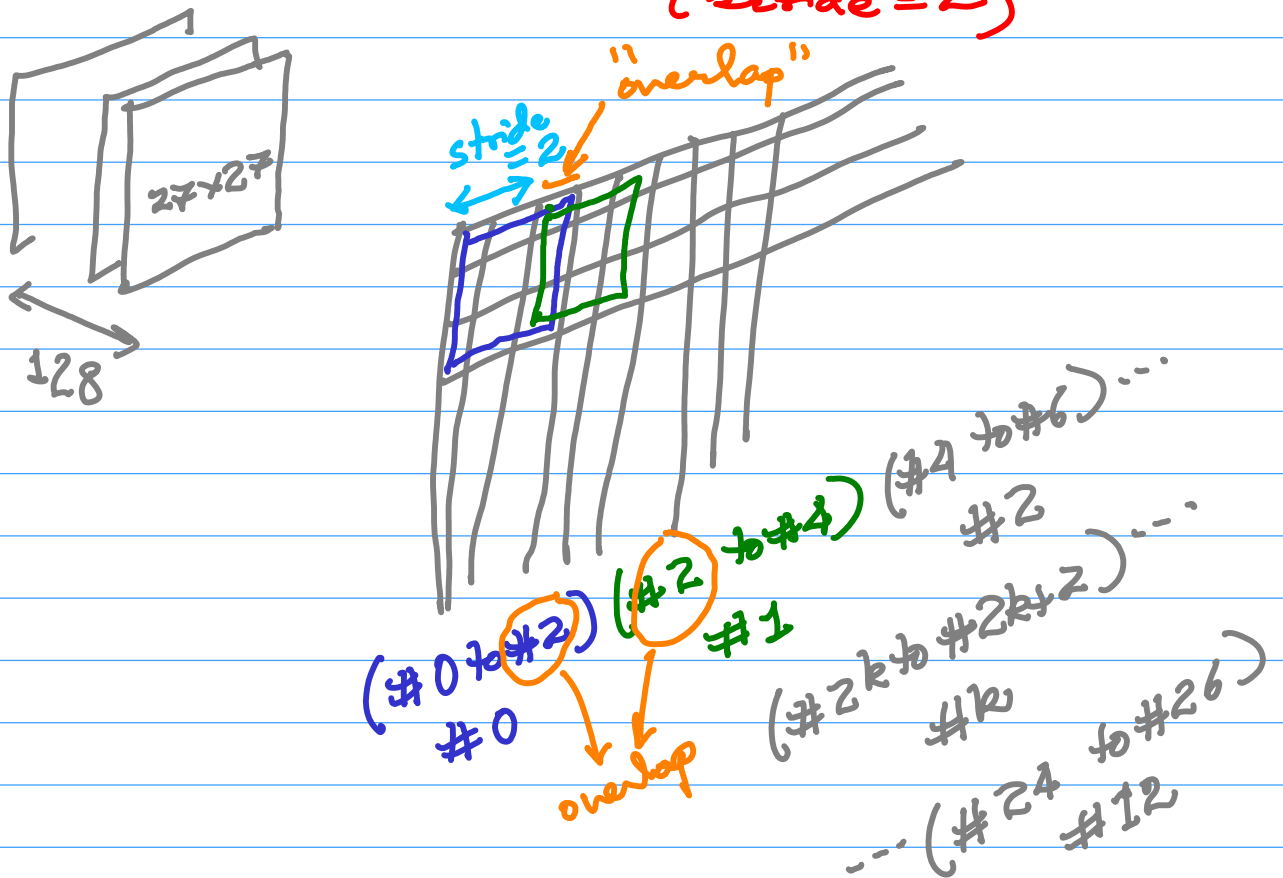of pixels which would otherwise get lost, get made up with the padding

'perfect fit' with no padding

2    2

2·82 pixels would get removed

⟹ with padding

2    2

padding is an extra 2 pixel layer

for the length & breadth

→ The resultant image size remains 27 X 27
128 filters/kernels in 2 groups
→ Output 27 X 27 X 128 X 2 groups.

⑤ next layer: 3 X 3 overlapping Max Pooling
(stride = 2)



→ there are 13 X 13 X 128 outputs
in 2 groups.

⑥ Local Response Normalisation
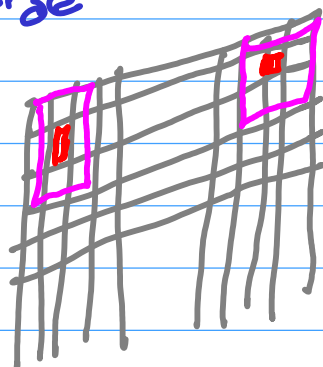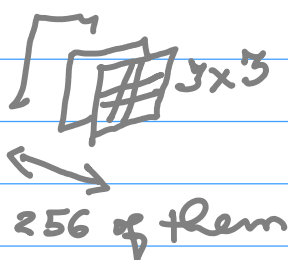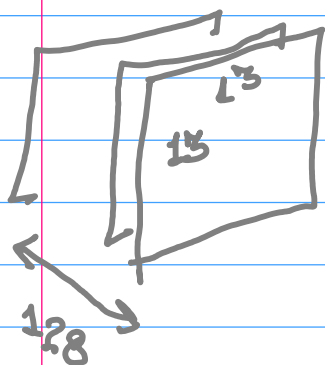(This does not change the output size: size −
preserving, but changes the values)
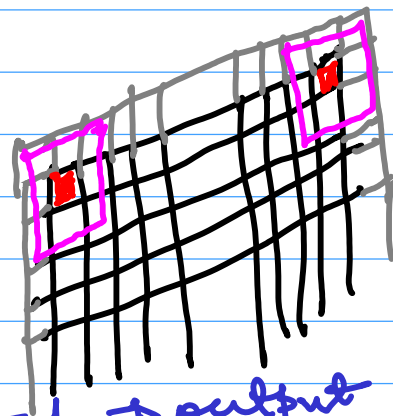
(7) Third Convolutional layer:

2 groups of 192 kernels of size $3 \times 3 \times 256$

stride = 1, pad = 1

Input : $13 \times 13$ images $\times$ 128 in 2 groups.

First, let us consider the size

256 of them

'just fit' $\rightarrow$ $12 \times 12$

pad = 1 $\Rightarrow$ output $13 \times 13$

Now, how do we account for the number?
$\longrightarrow$ heuristic!

To resolve 128, 256 and 192

128 $\xrightarrow{\quad + 128/2 \quad}$ 192 $\xrightarrow{\quad +128/2 \quad}$ 256

nothing mentioned about pooling, so possibly
2 kernels each for the 128 to give 256 and then

some selection and pooling to give 192.
now that we have an understanding of padding,
convolutions, stride & pooling,
we will recognise that there are many heuristics
to get actual numbers.
→ Try to look for conceptual ideas from
    different families of successful
                                architectures.

## VGG - 16/-19     "Visual Geometry Group"
                    at the University of Oxford

### Basic Concept:                 Karen Simonyan,
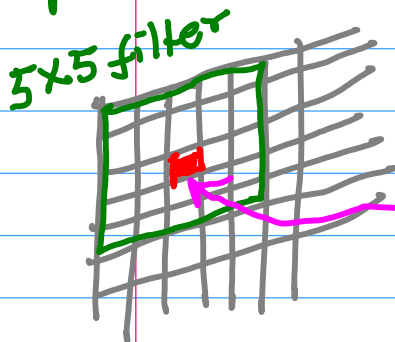                                   Andrew Zisserman (2014)

The use of 3×3 filters/kernels
in place of larger 11×11 or 7×7 filters.
Result: Simpler architecture with a smaller no. of
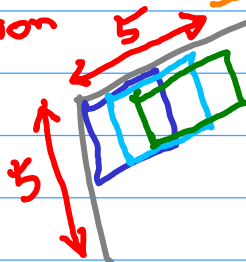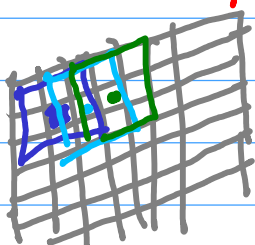parameters, but an increased depth: (16 - 19 layers)



5×5 filter

large 5×5 filter
superimposed on an
image, creating an
output at this pixel
position
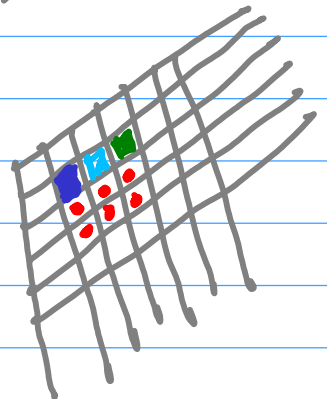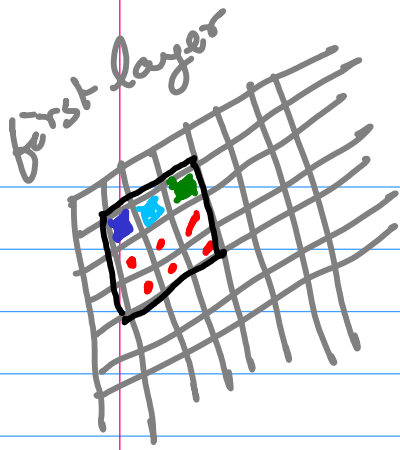
instead
of this

2 layers of 3×3
filters: cover the
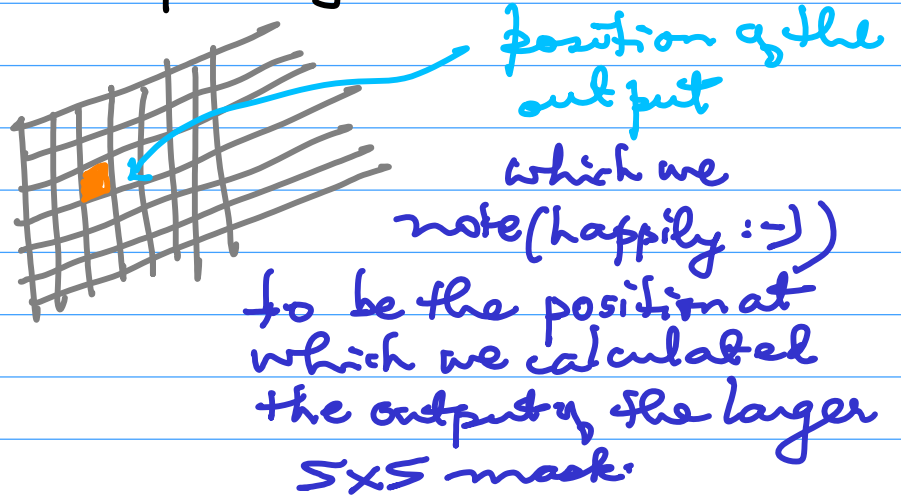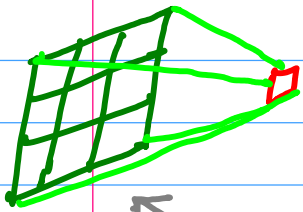same effective
pixel positions as one
larger 5×5 mask

first
layer

stride=1

5

5

we note that the pixel positions in the first layer allow a 3×3 mask/ filter/ kernel at another layer to be put right here
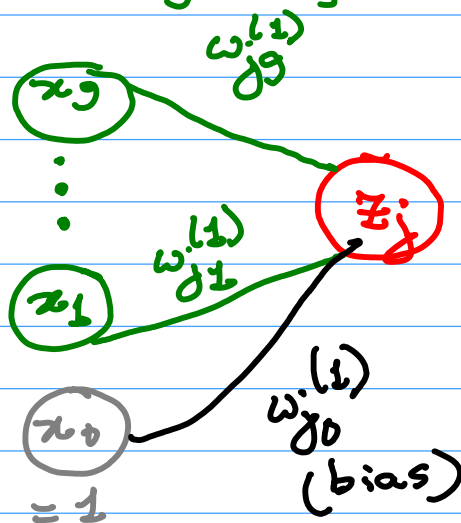
position of the output

which we note(happily :-))

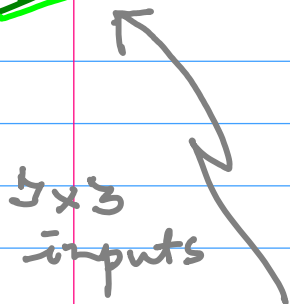to be the position at which we calculate the output of the larger 5×5 mask.

→ Advantage: there is a smaller # of parameters

(**) 5×5 filter : $5 \times 5 + 1$ (bias) $= 26$ parameters

2 layers of 3×3 : $2 \times \left( 3 \times 3 + 1 \right) = 2 \times 10 = 20$ parameters

bias

These weights (parameters): they are the relative-position-invariant weights of the local receptive field

3×3 inputs

$x_9$

$\omega_{j9}^{(1)}$

$\omega_{j1}^{(1)}$

$z_j$

$x_1$

$x_0$
$= 1$

$\omega_{j0}^{(1)}$
(bias)

$$z_j = \sum_{i=1}^{9} \omega_{ji}^{(1)} x_i + \omega_{j0}$$

(activation)