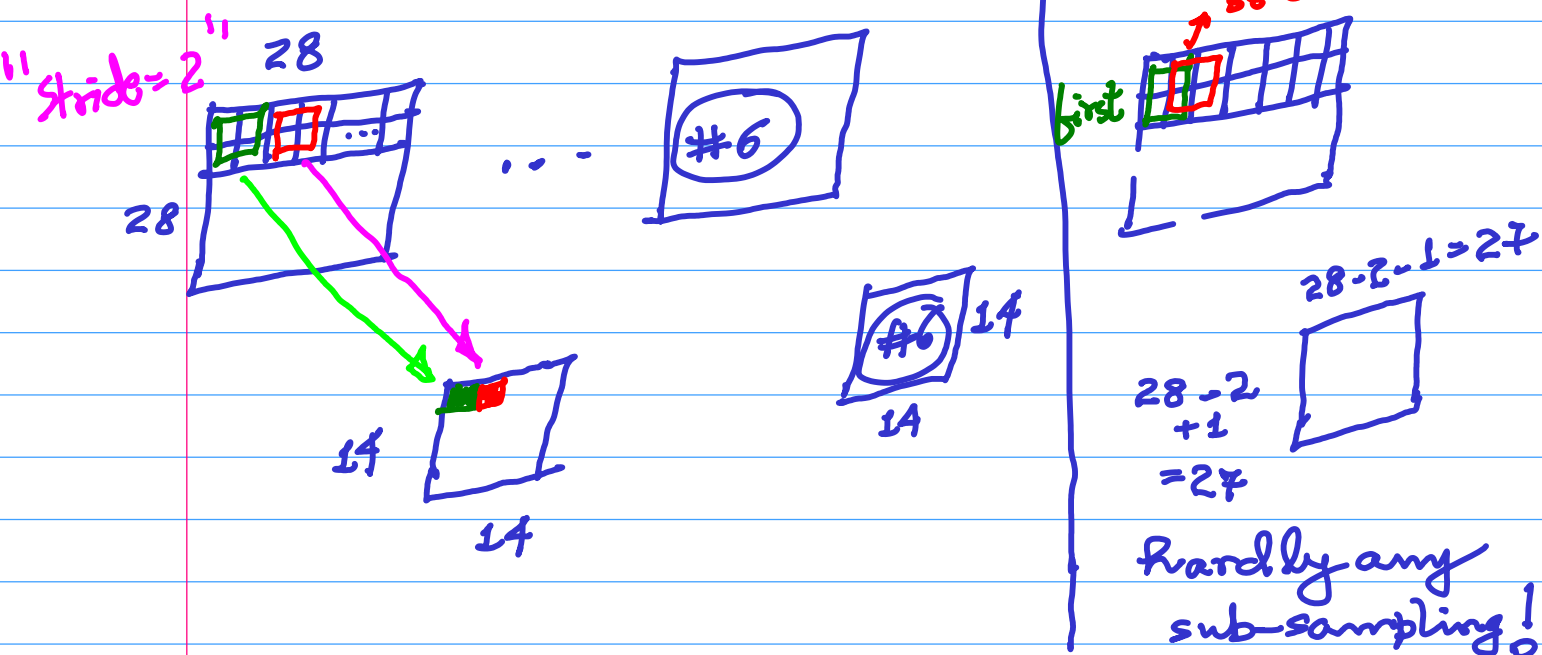
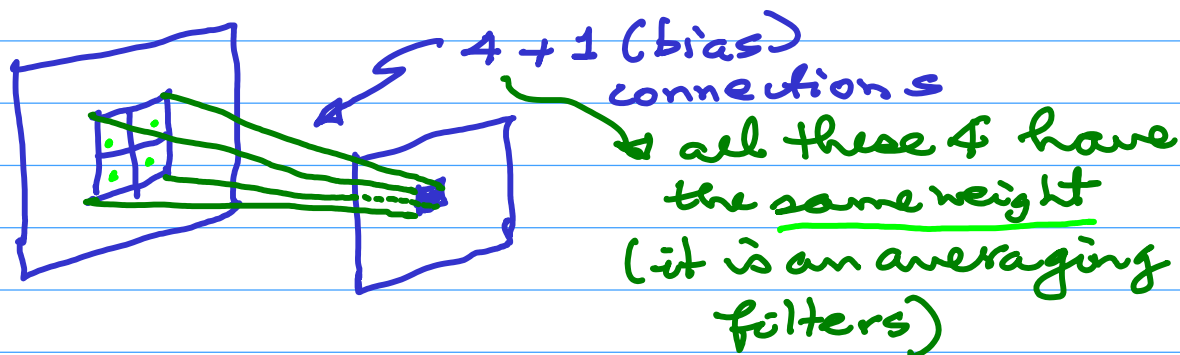


## Layer #2 (S2) Subsampling

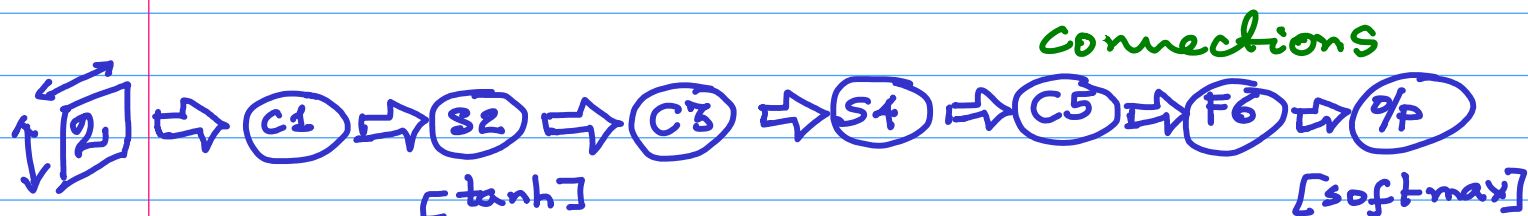


parameters & connections?



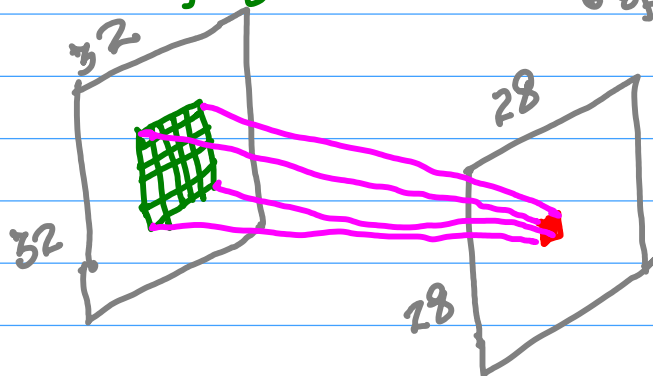
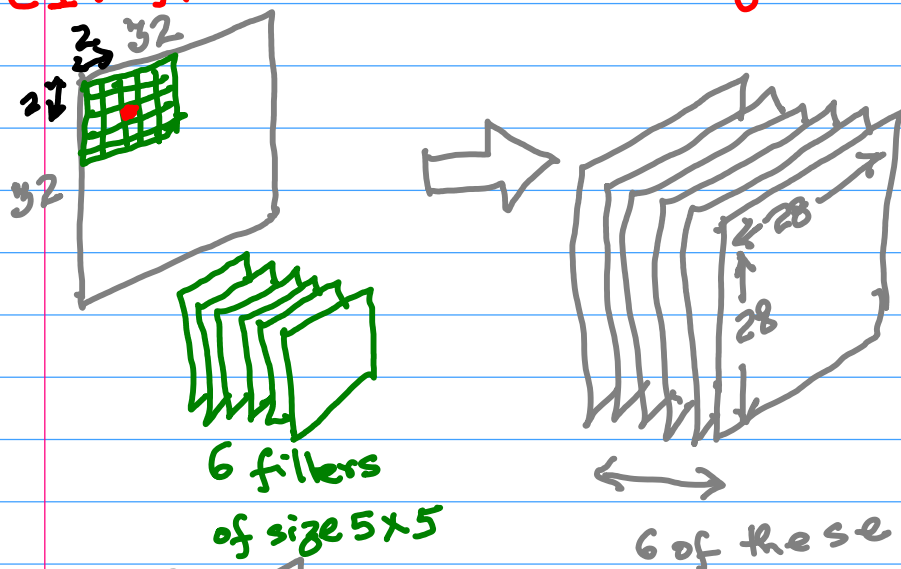
Parameters  $6 \times (1 + 1 (\text{bias})) = 12$  parameters

Connections  $6 \times (14 \times 14) \times (2 \times 2 + 1)$   
 $= 6 \times 14 \times 14 \times 5 = 5,880$



# LeNet-5 DETAILS

## C1: first convolutional layer



$$a_j = \sum_{i=1}^{25} w_{ji}^{(1)} x_i + w_{j0} \text{ (bias)}$$

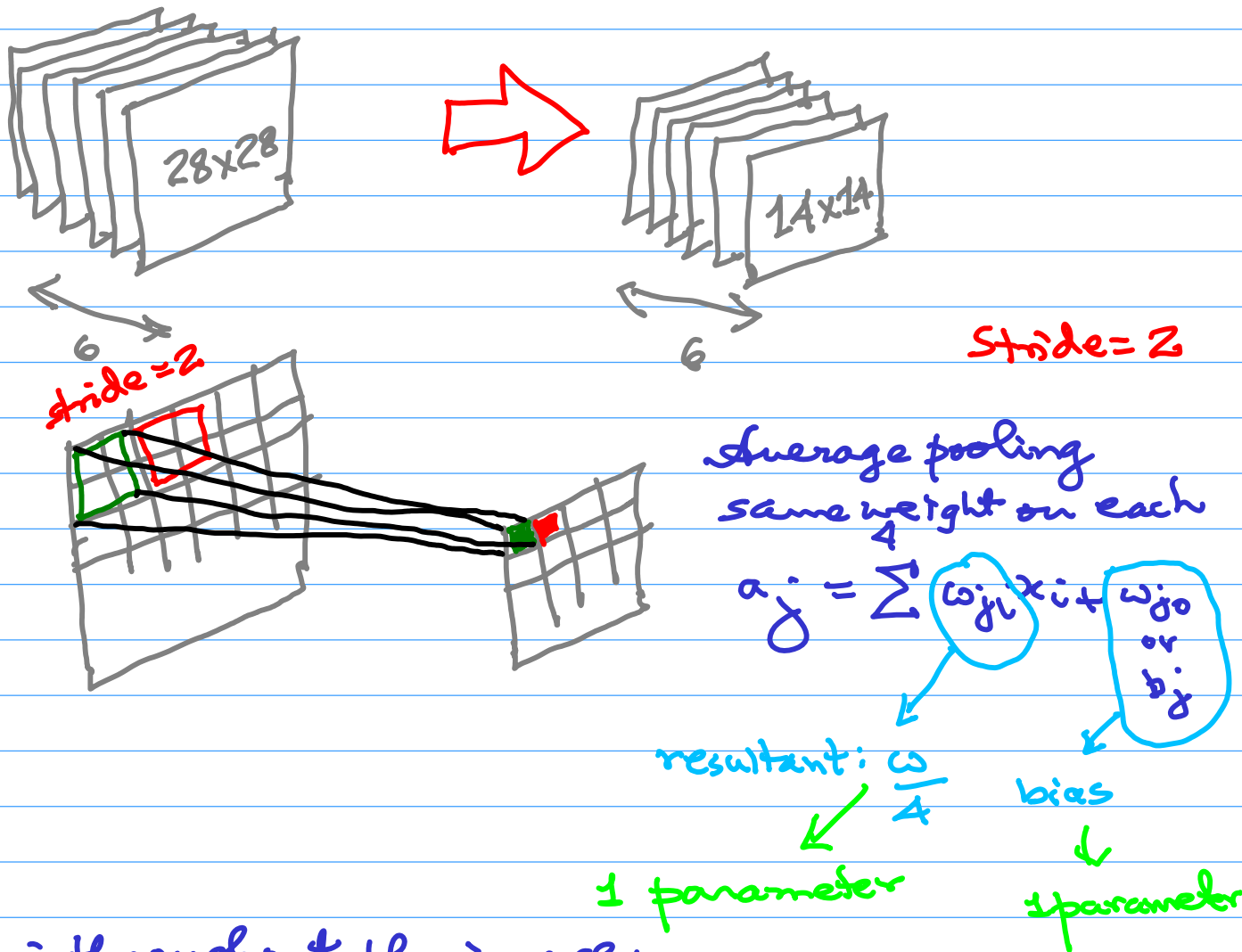
parameters

$$6 \times [5 \times 5 + 1] = 6 \times 26 = 156$$

connections: for each of the 6 filters, for each pixel in the resultant 28x28 image

$$28 \times 28 \times 6 \times \underbrace{[5 \times 5 + 1]}_{26} = 1,22,304$$

## Second layer (32) "Sub-sampling" or Average Pooling layer



this is throughout the image,  
is independent of the image size ( $14 \times 14$ )  
 $2 \times 6 = 12$  parameters

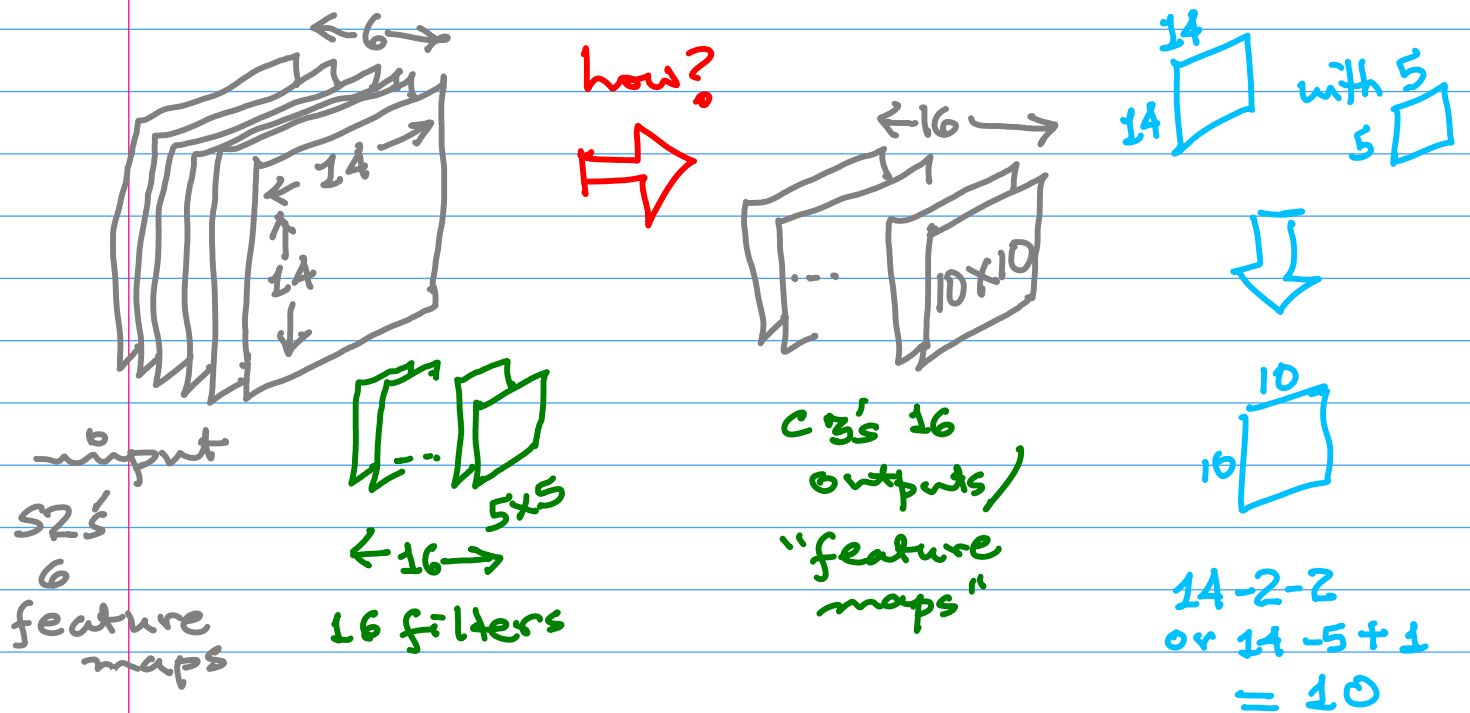
How many connections:

$$6 \times (4 + 1) \times 14 \times 14 = 5,880 \text{ connections}$$

connects image

for each filter/  
channel

## C3: 3rd Convolutional layer



How do 6 images of size 14x14  
map onto 16 images of size 10x10

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	x				x	x	x			x	x	x	x		x	x
1	x	x				x	x	x			x	x	x	x		x
2	x	x	x				x	x	x			x		x	x	x
3		x	x	x			x	x	x	x			x		x	x
4			x	x	x			x	x	x	x		x	x		x
5				x	x	x			x	x	x	x		x	x	x

"3 at a time"      "4 at a time"      "all 6"

Each column indicates which feature maps in S2 are combined by the units in a particular feature map of C3

- To break the symmetry in the network
- to keep the number of connections within reasonable bounds.

16 filters of size  $5 \times 5$  give us 16 output feature maps of size  $10 \times 10$

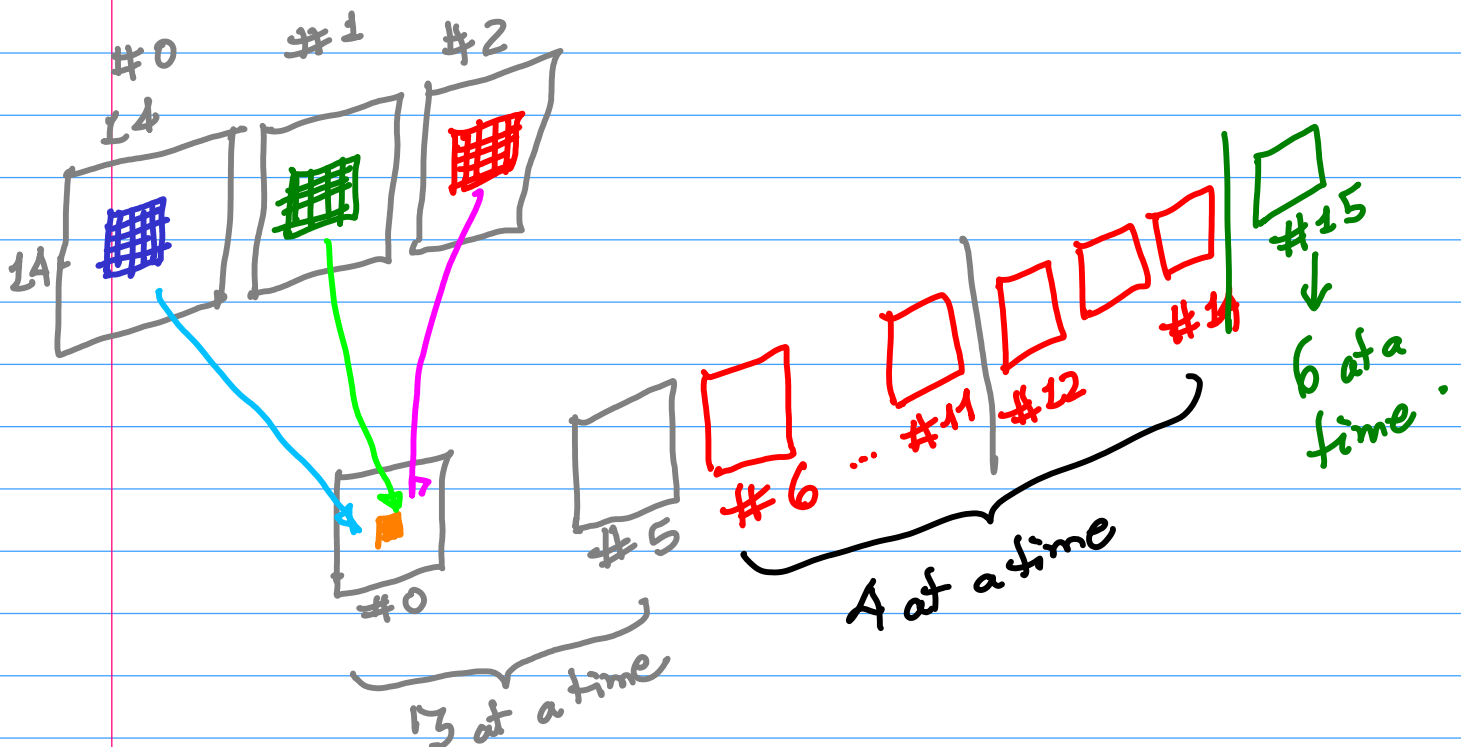
However, which ones of the 6 input feature maps of the previous layer do we consider, to take inputs from? Any one / some / all?

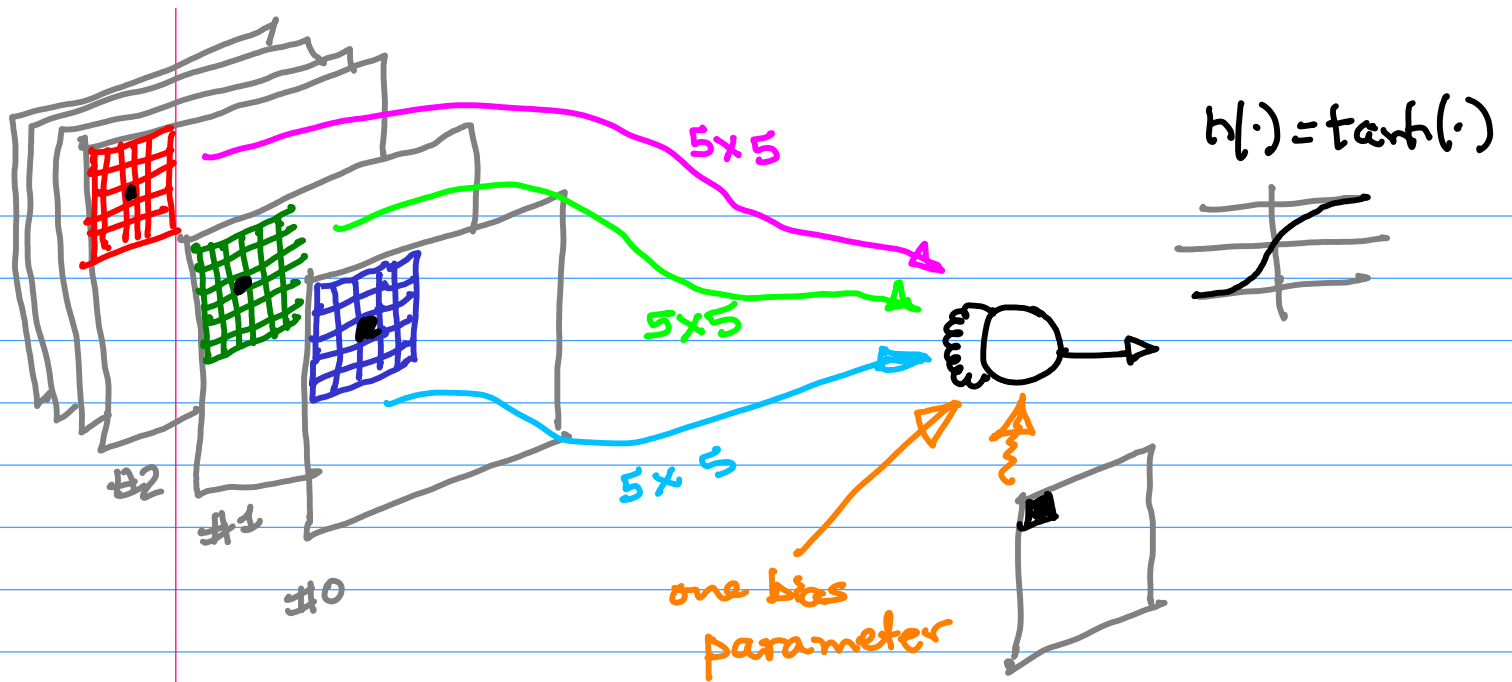
Answer: 3 at a time + 4 at a time + 6 at a time.

How many parameters?

"3 at a time"	"4 at a time"	"all 6"
(#0 to #5: 6 of these)	(#6 to #14: 9 of these)	(#15: 1 of this)
take input from 3	take input from 4	takes input from 6
$6 \times (3 \times 5 \times 5 + 1)$	$+ 9 \times (4 \times 5 \times 5 + 1)$	$+ 1 \times (6 \times 5 \times 5 + 1)$
$= 6 \times (75 + 1)$	$+ 9 \times (100 + 1)$	$+ 1 \times (50 + 1)$
$= 1516$		

How many connections? for each pixel in the  $10 \times 10$  image =  $1516 \times 100 = 1,51,600$  connections.

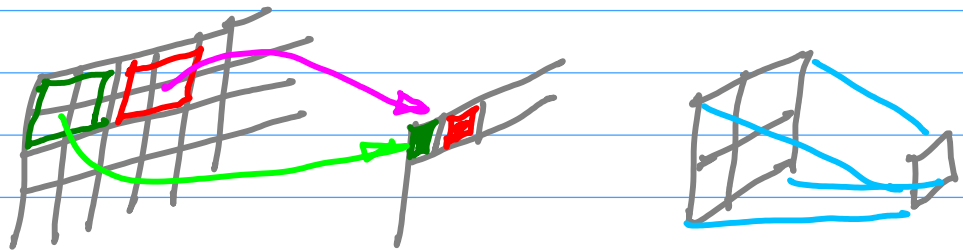
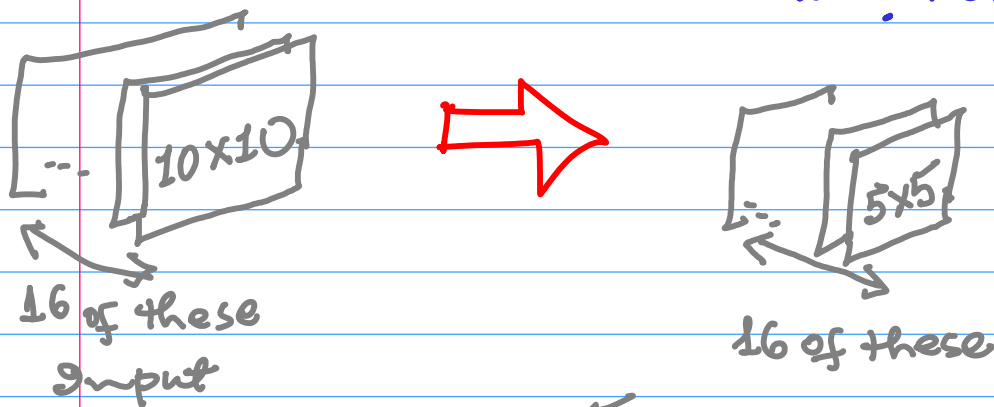




Example: for the first out of the "3 at a time" cases.

#### (4) "S4" the fourth (subsampling) layer

how? half the dimensions  
 $\Rightarrow$  stride = 2



Same weight for all 4 connections  
 + 1 bias = 2 parameters,  
 across the entire image, for each of the  
 16 output 5x5 images  
 $= (1+1) \times 16 = 32$ .

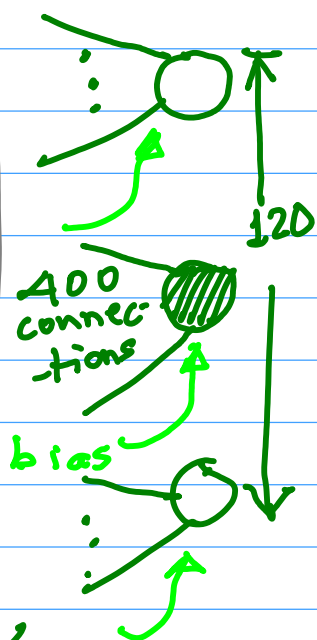
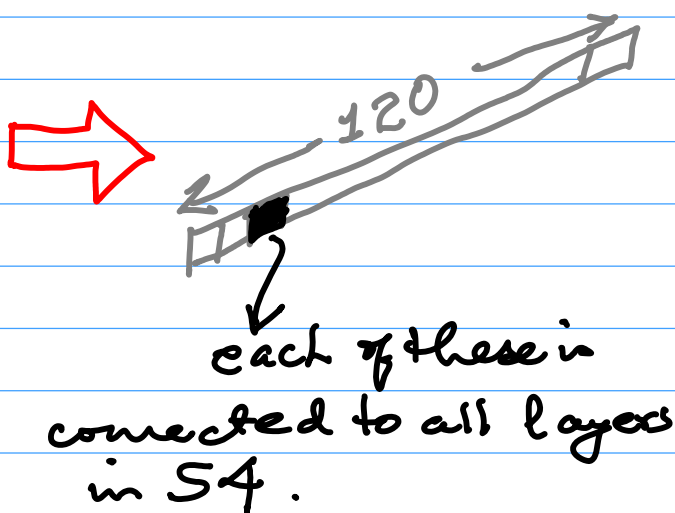
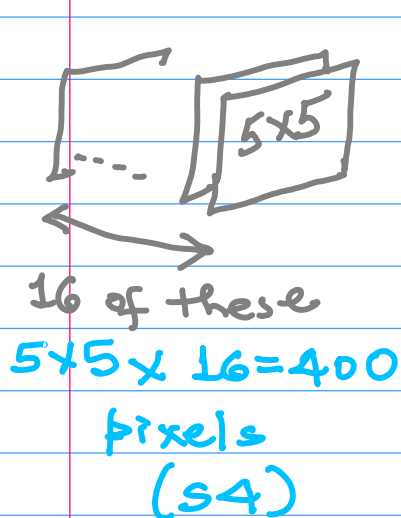
number of connections: The above structure repeats for each pixel in the  $5 \times 5$  output image: 16 of them:

$$\Rightarrow \underbrace{(5 \times 5)}_{\text{size}} \times \underbrace{16}_{\text{16 images}} \times \underbrace{(4 + 1)}_{\substack{\text{square:} \\ 2 \times 2}} = 5 \times 5 \times 80$$

bias = 2000

⑤ "C5" 5th layer, C  $\Rightarrow$  convolution

"Flattening"



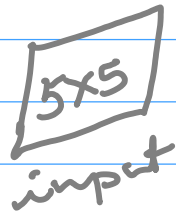
A fully connected layer has  
 no local receptive fields

parameters: For each of these 120 neurons e.g., the  $j$ th one: there are 400 weights + 1 bias term  
 $\Rightarrow 120 \times (400 + 1) = 120 \times 401 = 48,120$   
 parameters

connections: The number of connections is the same as the number of parameters. This is

a fully connected layer.

Why is this called a Convolutional layer



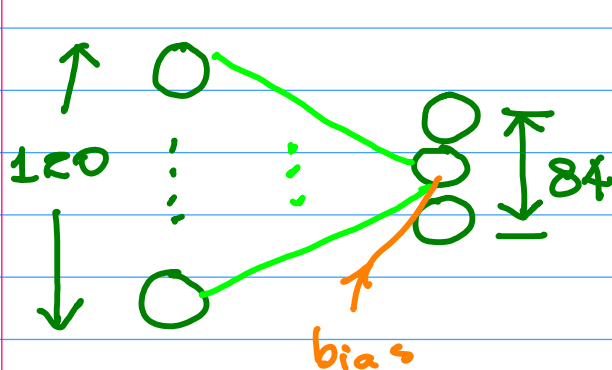
with a 5x5 mask

results in a 1x1 image;

there are 120 of these.

Why 120 & not any other number?  
(arbitrary!)

## ⑥ "F6" Fully connected layer



Parameters:

$$84 \times [120 + 1]$$

# of outputs

weights

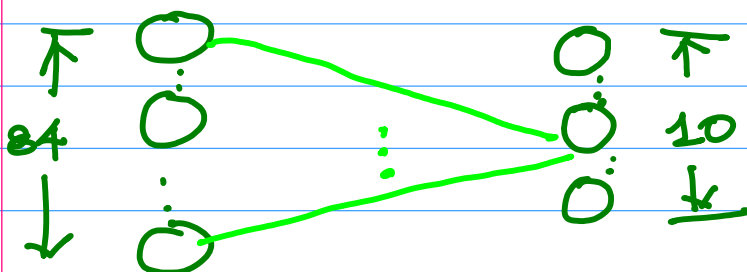
bias

$$= 84 \times 121$$

$$= 10,164$$

Connections: (same)

## ⑦ Output layer



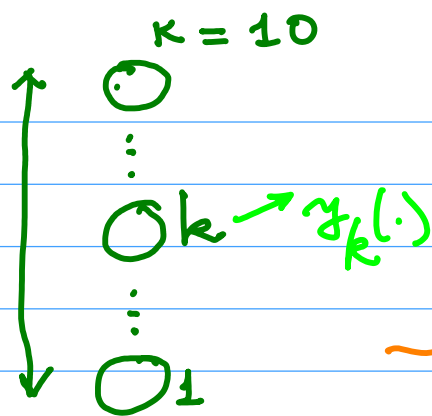
The activation functions at all other layers were the  $\tanh(\cdot)$  function

(smooth version of

the sigmoid function of the Perceptron)

For this output layer, it is the SOFTMAX





$$\sigma_k(y) \triangleq \frac{e^{y_k}}{\sum_l e^{y_l}}$$

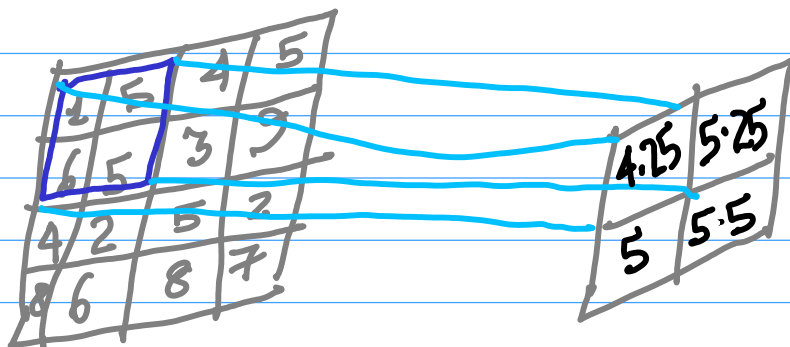
why?

normalisation: getting a probability.  $\rightarrow$  like a predictor.  
 why did we raise a number to an exponential?  $\rightarrow$  large #, & not a small ratio.

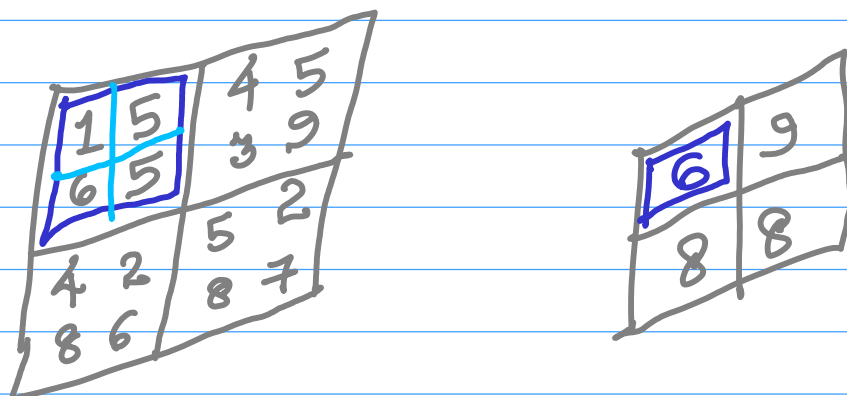
## SOME IMPORTANT CONCEPTS

### POOLING:

#### 1) Average Pooling (LeNet)



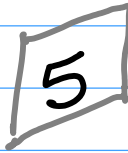
#### 2) Max Pooling (AlexNet)



### 5) Global Average Pooling

1	5	4	5
6	5	3	9
4	2	5	2
8	6	8	7

The entire frame into one  
 $\Sigma = 80, 80/4 \times 4 = 5$

⇒  (16 pixels:  
4x4  
grid)

Why Pooling? — contraction

— Invariance. Invariance to small transformations, distortions, translation. A small distortion in the input will not change the outcome of pooling drastically since we take the max / average value in a local neighbourhood

### LOCAL RESPONSE NORMALISATION (LRN)

from neuro-biology

(AlexNet)

→ 'lateral inhibition' → an excited neuron inhibits

its neighbours → creates contrast in the area and increases sensory perception

→ No solid mathematical background

→ motivated by biology (evolution), shown statistically.

## Batch Normalisation (2015) ["BN"]

[Ioffe & Szegedy, 2015]

→ Addresses the problem of Internal Covariate Shift (initially believed) current opinion: works because it smooths the objective function.

At initialisation: actually induces severe gradient explosion, which is only alleviated by skip connections in residual networks

the basic issue (historically) Each layer of a neural network has inputs with a corresponding distribution. This is affected during the training process by the randomisation

[ - the parameter initialisation  
- the input

- During training, as parameters of the preceding layers change, the distribution of inputs to the current layer changes accordingly
- the current layer needs to constantly adjust to new distributions.
- small changes in the initial layers amplify, and result in a significant shift in deeper hidden layers

### BENEFITS:

- Reduce unwanted shifts to speed up training
- permits a higher learning rate without vanishing/exploding gradients
- Regularisation effect: unnecessary to use 'dropout' to mitigate overfitting
- Robustness to different initialisation schemes and learning rates

Stochastic Optimisation  $\rightarrow$  not all training data is taken together  $\rightarrow$  normalisation is done in batches. Batch  $B$  has size  $m$

$$\mu_B = \frac{1}{m} \sum_{n=1}^m x^{(n)}; \sigma_B^2 = \frac{1}{m} \sum_{n=1}^m [x^{(n)} - \mu_B]^2$$

$x$ :  $x$  is a scalar  $\rightarrow$  one of the input dimensions call it  $x_i$ , say.

$$\hat{x}_i \triangleq \frac{x_i - \mu_{B_i}}{\sqrt{\sigma_{B_i}^2 + \epsilon}}$$

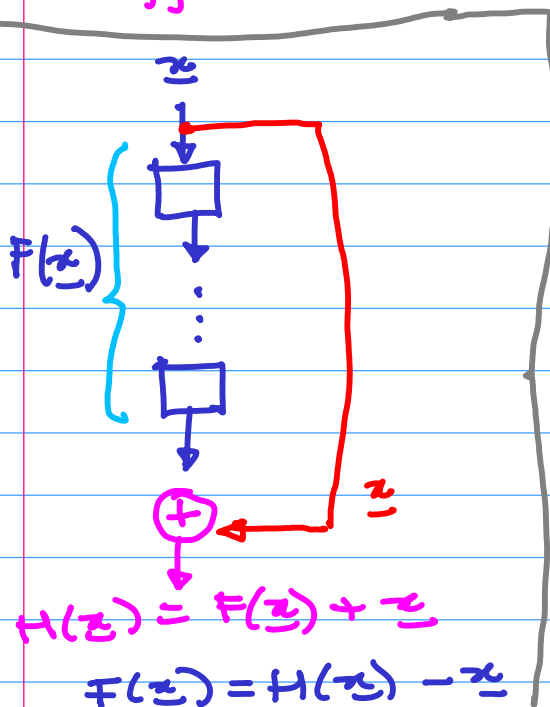
small positive constant used for numerical stability

$$y_i \triangleq \gamma_i \hat{x}_i + \beta_i$$

tunable parameters

## Residual Connections / Skip Connections / Highway networks

Trend towards deeper neural networks  $\rightarrow$  better accuracy or performance  
trade-off: harder for training to converge



Used widely  $\rightarrow$

ResNet (Computer Vision),  
Transformer (NLP, Computer vision)  
AlphaZero (RL),  
AlphaFold (Protein Structure prediction)

### Issues:

- Difficult to learn an identity mapping across layers
- Training a deep network is difficult because of exploding and vanishing gradients
- Observation: convolutional layers are often better at learning the residual rather than learning the feature map directly.

## (\*) Residual connections/skip connections/ highway connections (contd.)

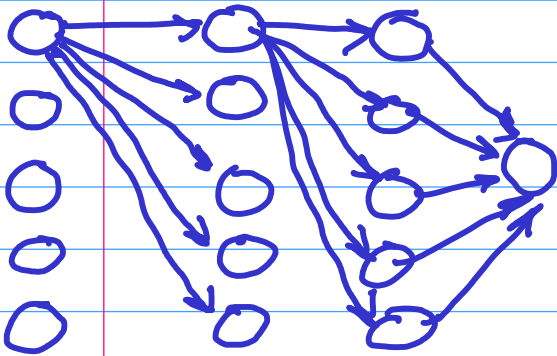
(\*) The magnitude of the problem: →

- AlexNet (2012) had 5 convolutional layers
- 2014: VGG, GoogleNet
  - 19 layers
  - 22 layers
- ResNet (34 → 50 layers deep architectures)  
deeper but had overall lower complexity.

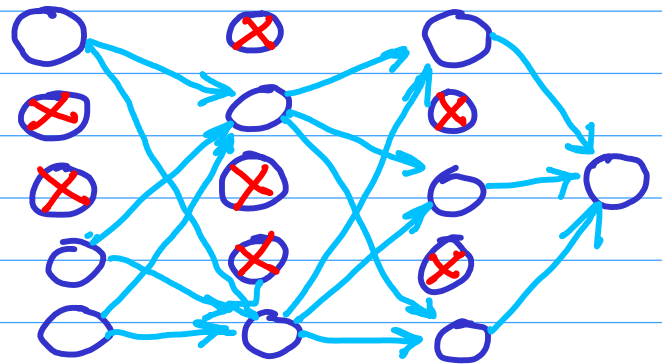
## (\*) 'DROPOUT' (AlexNet 2012)

Heuristic: applied at the training phase (FC layers)

Dropout is a kind of regularisation technique to reduce overfitting.



Usual FC scenario



Scenario with dropout

AlexNet :  $p = 0.5$  at the first two fully connected layers

Neuron : has a probability not to contribute to the feedforward phase & participate in the backpropagation. → Each neuron can have a larger chance to be trained, and not depend on some 'strong' neuron. No dropout at the test time

## (\*) Residual connections/skip connections/highway connections (contd.)

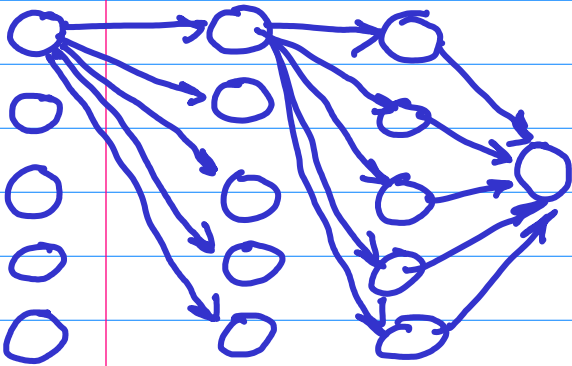
(\*) The magnitude of the problem: →

- AlexNet (2012) had 5 convolutional layers
- 2014: VGG, GoogleNet
  - 19 layers
  - 22 layers
- ResNet (34 → 50 layers deep architectures) deeper but had overall lower complexity.

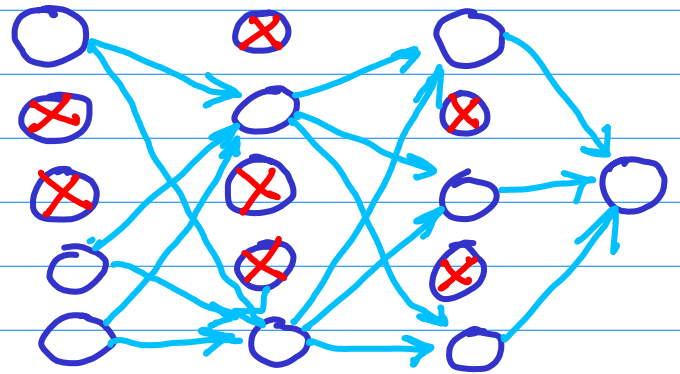
## (\*) "DROPOUT" (AlexNet 2012)

Heuristic: applied at the training phase (FC layers)

Dropout is a kind of regularisation technique to reduce over-fitting.



Usual FC scenario



Scenario with dropout

AlexNet:  $p = 0.5$  at the first two fully connected layers

Neuron: has a probability not to contribute to the feedforward phase & participate in the backpropagation. → Each neuron can have a larger chance to be trained, and not depend on some 'strong' neuron. No dropout at the test time

## Second class of Successful deep Architectures: AlexNet (2012) and CaffeNet

↳ used 2 GPUs

(two parallel paths)

↳ used 1 GPU

What & Why? ImageNet: 15 million labelled  
high-resolution images, with 22K categories  
(22,000)

2012: ImageNet Large-Scale  
Visual Recognition Competition

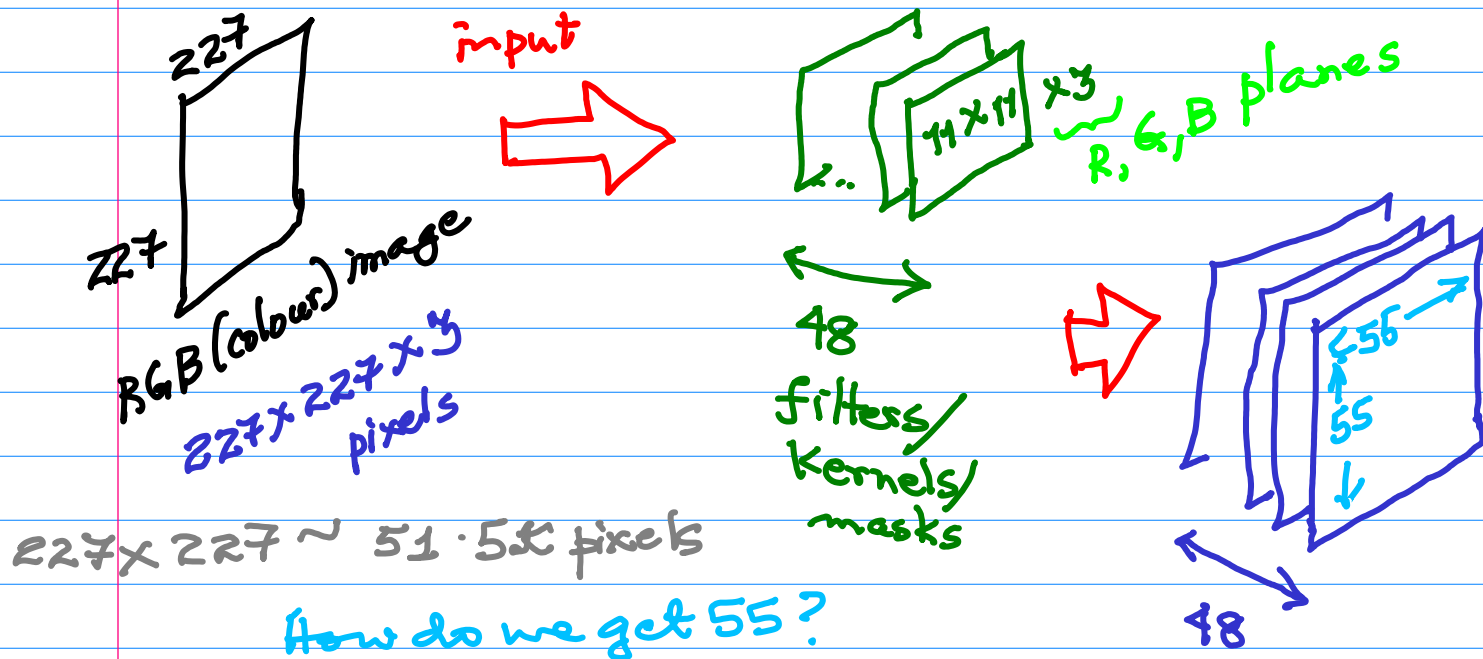
dataset: subset of ImageNet

\* 1K images in each of the 1K categories/  
classes.

[ 1.2 million training images  
50K validation images  
1.5K testing images

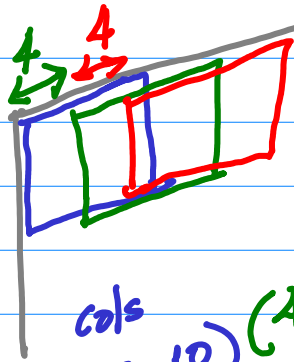
AlexNet: 8 layers (5 conv layers)

### (1) 1st Convolutional layer





Stride = 4



cols

(0 to 10)  
#0

(4 to 14)  
#1

(8 to 18)  
#2

(4k to 4k+10)  
#k

(216 to 226)  
#?

$$4k = 216$$

$$k = \frac{216}{4}$$

$$k = 54$$

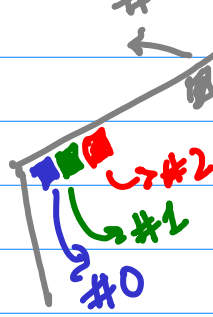
#54

#k

#54

55x55

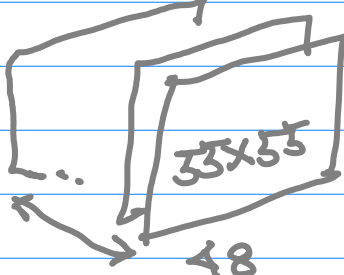
Result



AlexNet has 2 groups of 48 filters/kernels/masks each, each group given to one of the two GPUs used.

[CaffeNet is a 1 GPU version,  
96 filters / kernels / masks]

→ Result: 2 groups of 48 outputs,  
(AlexNet) each of size 55x55



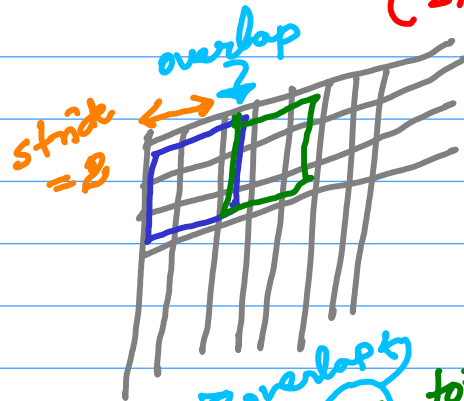
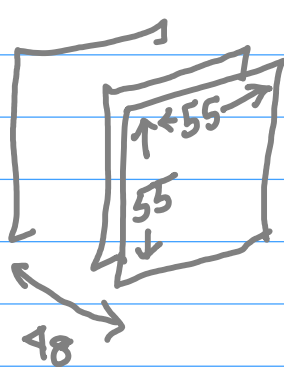
group #0  
(GPU #0)



group #1  
(GPU #1)

## AlexNet (contd.)

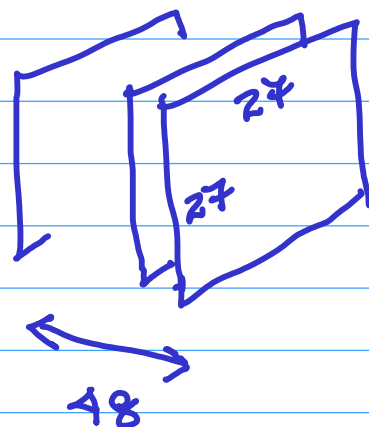
② Second layer:  $3 \times 3$  overlapping Max Pooling  
(stride = 2)



(#0 to #2) (#2 to #4) (#4 to #6)  
#0 #2  
... (#2k to #2k+2) ... (#52 to #54)  
#k

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

→ There are  $27 \times 27$  outputs  
48 of these, 2 groups



③ 3rd layer: Local Response Normalisation

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

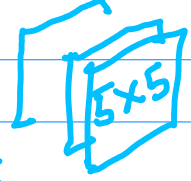
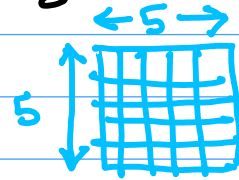
## ④ Second Convolutional Layer

2 groups of 128 kernels of size  $5 \times 5 \times 48$

Input:

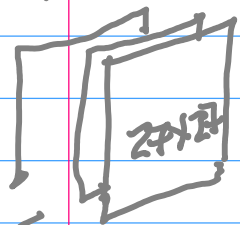
stride = 1, pad = 2

← given ←

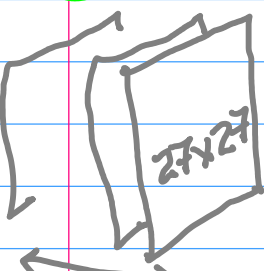


48

one-to-one link



48



48



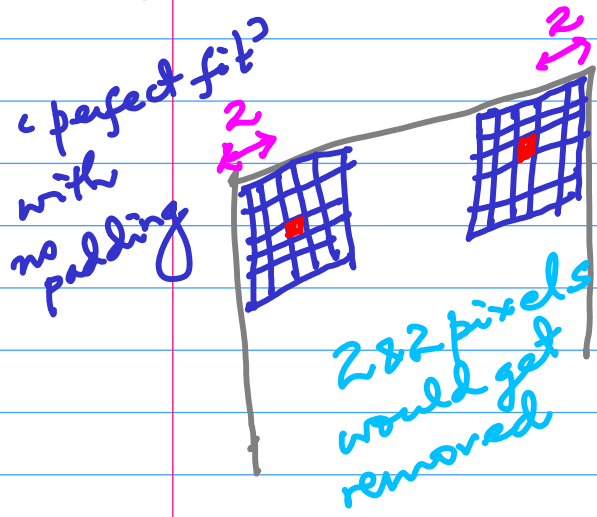
48

stride = 1

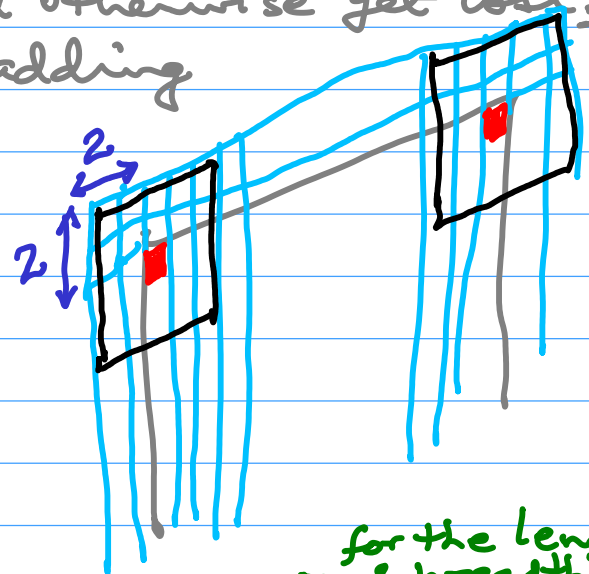
⇒ Result can get smaller with a 'perfect fit' convolution, but here

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

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



with padding



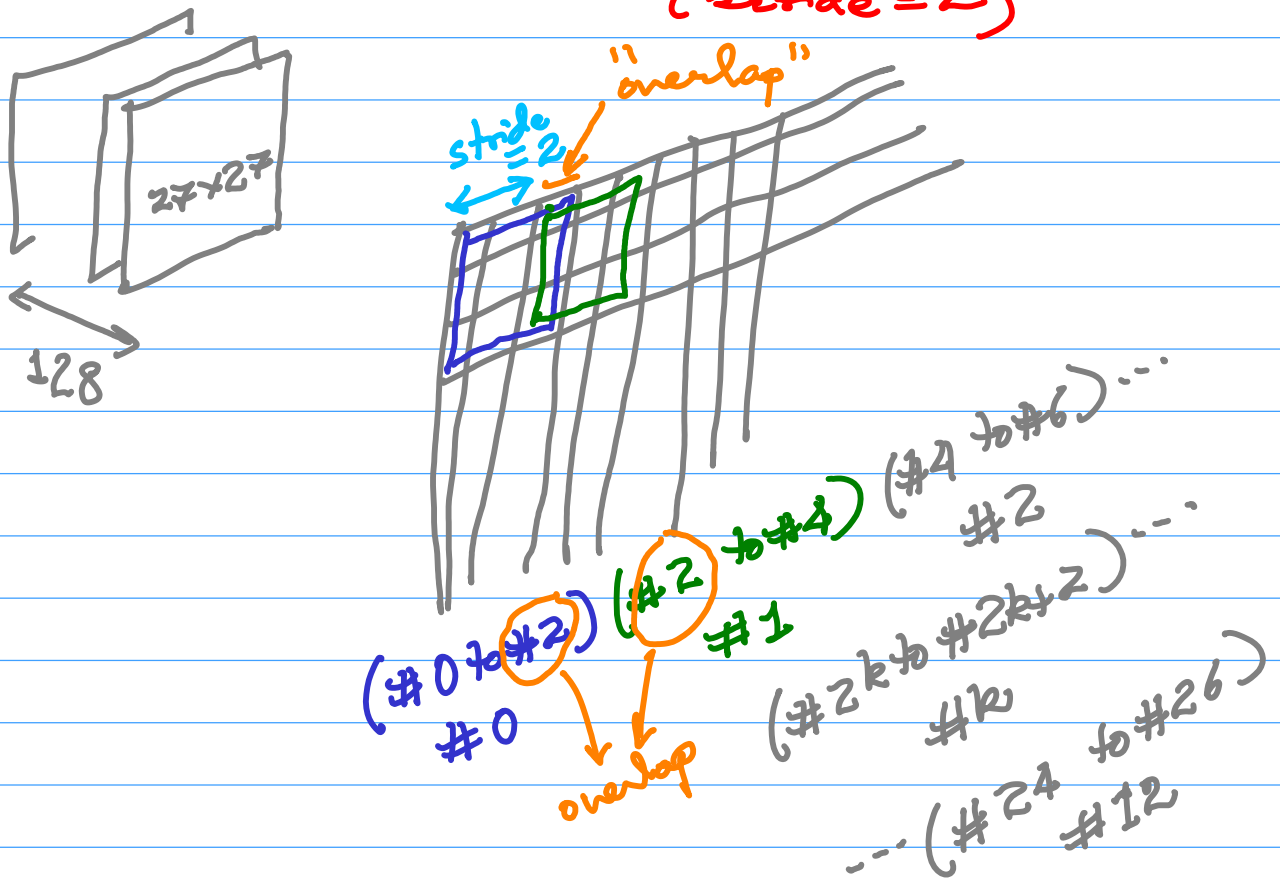
padding is an extra 2 pixel layer for the length & breadth

→ The resultant image size remains  $27 \times 27$   
128 filters/kernels in 2 groups

→ output  $27 \times 27 \times 128 \times 2$  groups.

⑤ next layer:  $3 \times 3$  overlapping Max Pooling

(stride = 2)



→ there are  $13 \times 13 \times 128$  outputs  
in 2 groups.

⑥ Local Response Normalisation

(This does not change the output size: size-preserving, but changes the values)

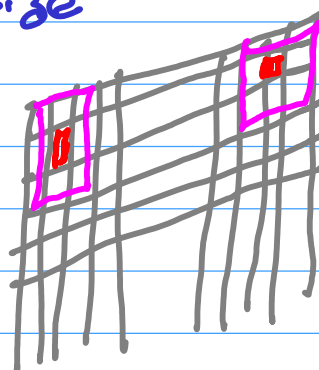
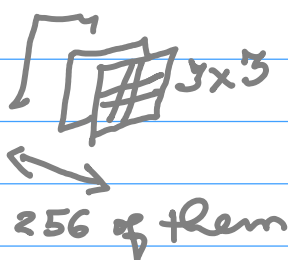
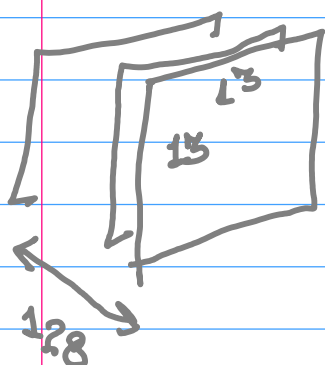
## AlexNet (contd).

### ⑦ Third convolutional layer:

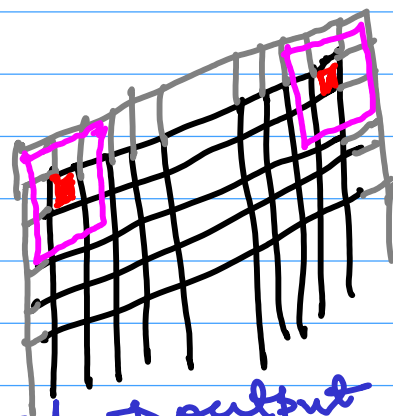
2 groups of 192 kernels of size  $3 \times 3 \times 256$   
stride = 1, pad = 1

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

first, let us consider  
the size



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



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

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

To resolve 128, 256 and 192

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

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.

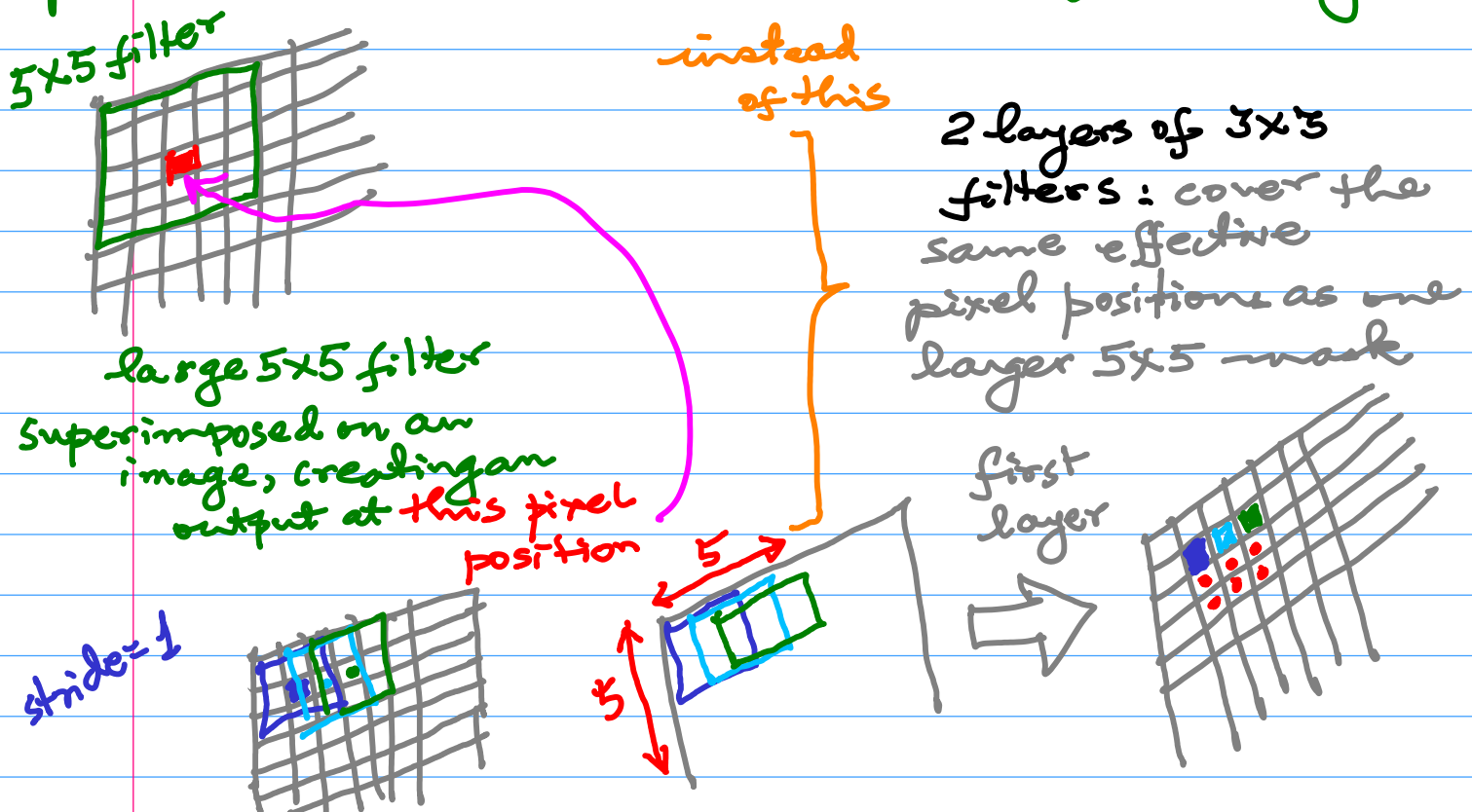
VGG - 16/-19 "Visual Geometry Group"  
at the University of Oxford

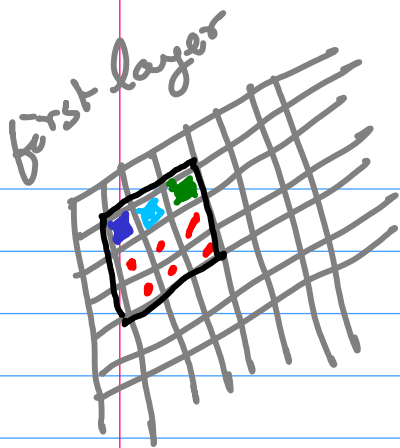
Basic Concept:

Karen Simonyan,  
Andrew Zisserman (2014)

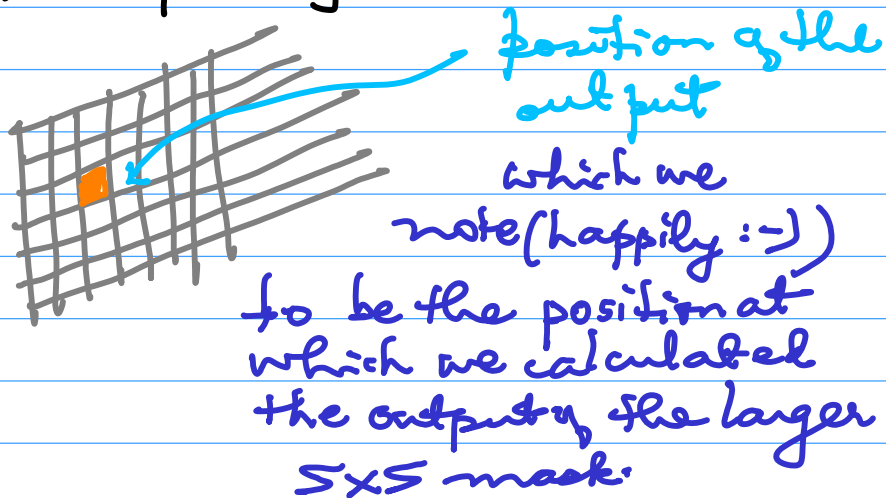
The use of  $3 \times 3$  filters/kernels  
in place of larger  $11 \times 11$  or  $7 \times 7$  filters.

Result: Simpler architecture with a smaller no. of  
parameters, but an increased depth: (16 - 19 layers)



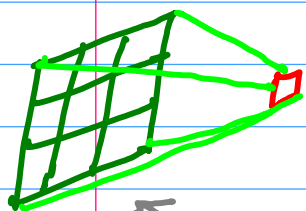


we note that the pixel positions in the first layer allow a  $3 \times 3$  mask/ filter/ kernel at another layer to be put right here



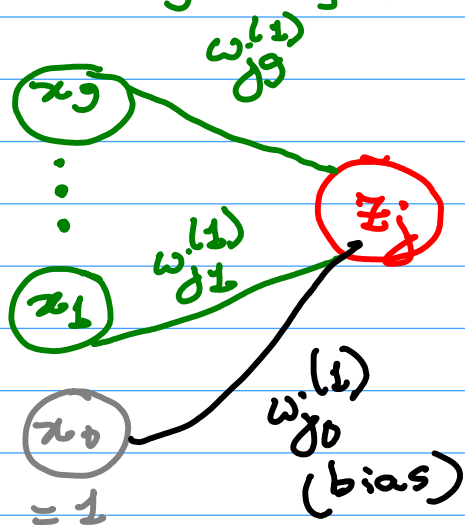
→ Advantage: there is a smaller # of parameters  
(\*)  $5 \times 5$  filter:  $5 \times 5 + 1$  (bias) = 26 parameters

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



$3 \times 3$  inputs

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



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

(activation)