

Published as a conference paper at ICLR 2015

---

# VERY DEEP CONVOLUTIONAL NETWORKS FOR LARGE-SCALE IMAGE RECOGNITION

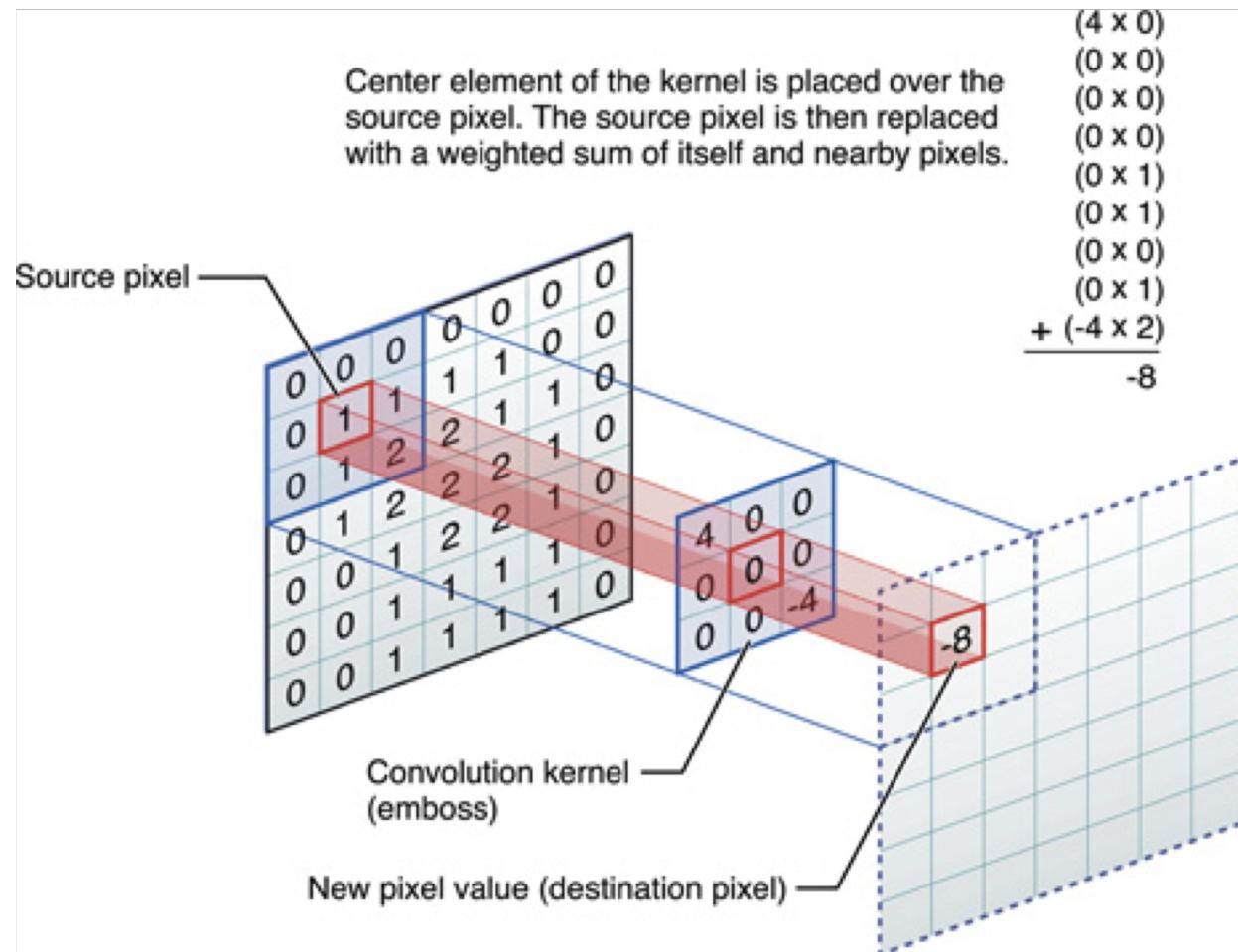
**Karen Simonyan\*** & **Andrew Zisserman<sup>+</sup>**

Visual Geometry Group, Department of Engineering Science, University of Oxford

{karen,az}@robots.ox.ac.uk

# Convolution Neural Net (CNN)

## Convolution in computer vision



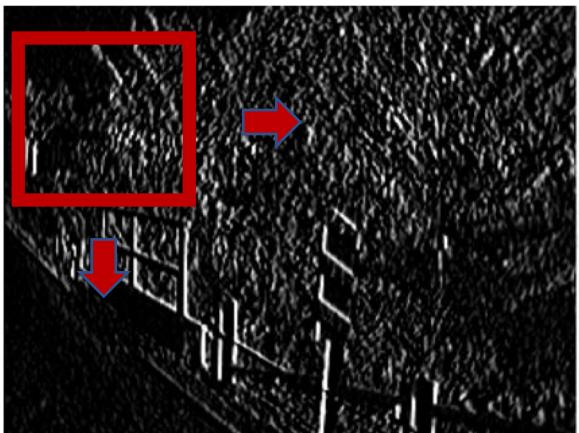
# Convolution Neural Net (CNN)

## Convolution in computer vision



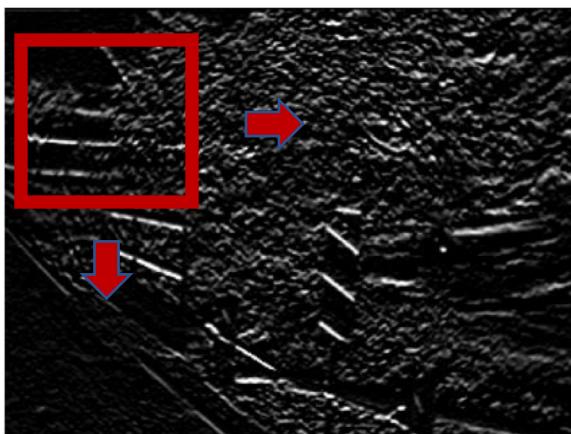
### Sobel Edge Detection: Gradient Approximation

*Note anisotropy of edge finding*



1	0	-1
2	0	-2
1	0	-1

Horizontal diff.



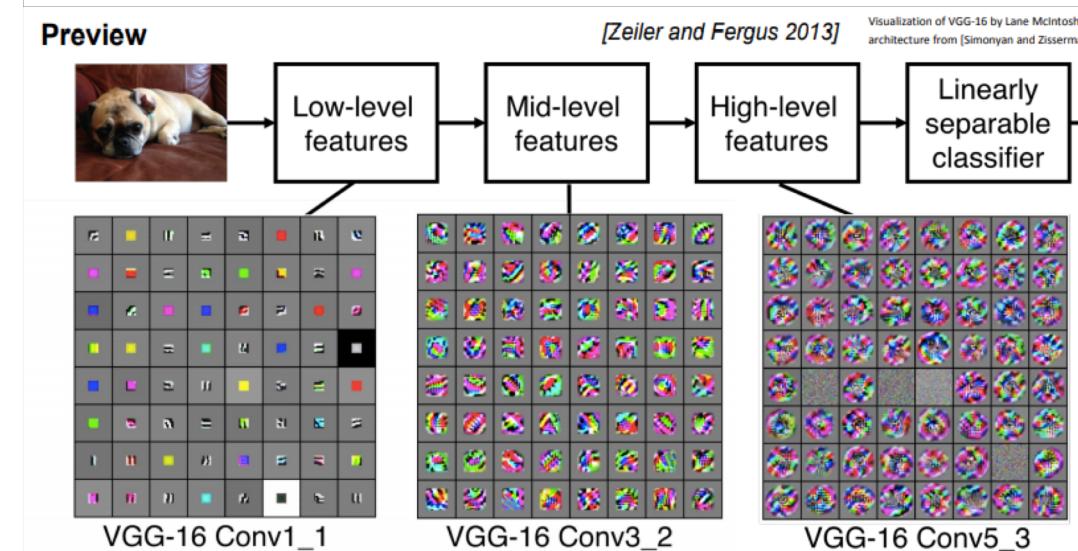
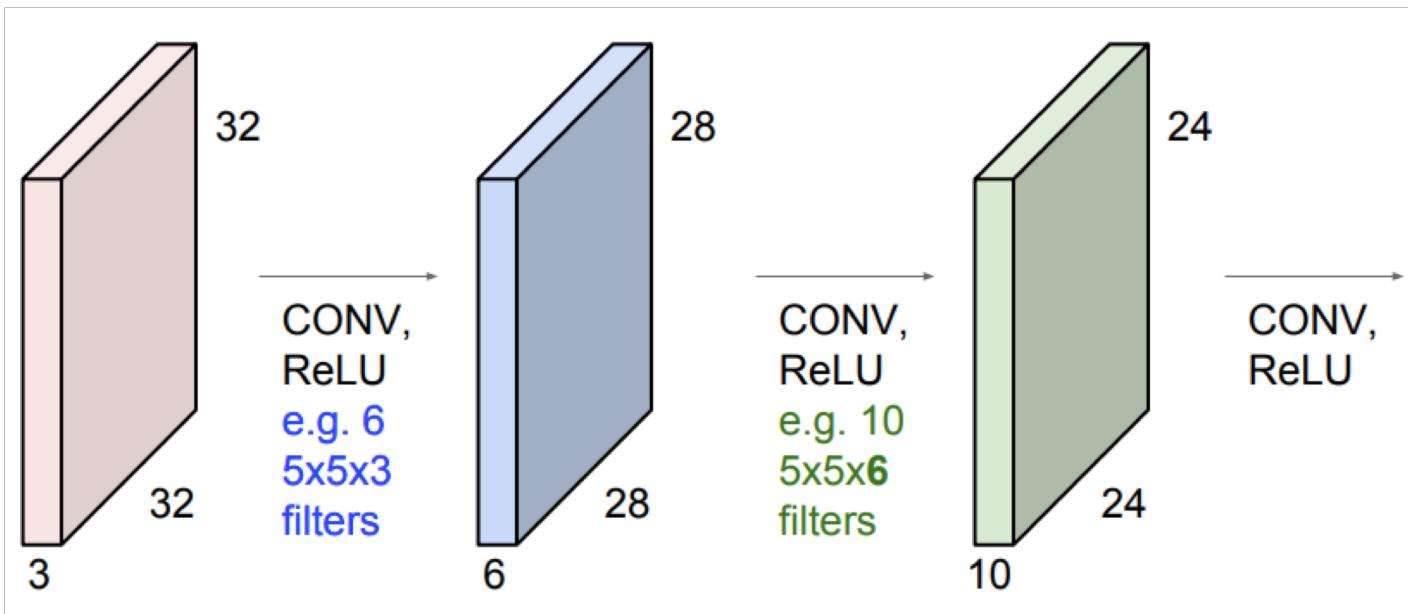
1	2	1
0	0	0
-1	-2	-1

Vertical diff.

Computer Vision :

- Extract feature  
ex) edge, ridge, etc...

# Convolution Neural Net (CNN)



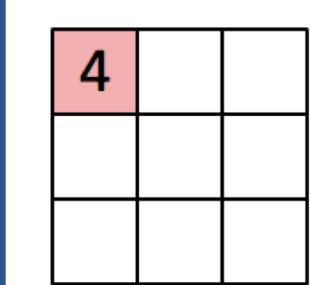
One filter => one activation map  
6,  $5 \times 5$  filter => get 6 activation map

# Convolution Neural Net (CNN)

## Convolution Layer

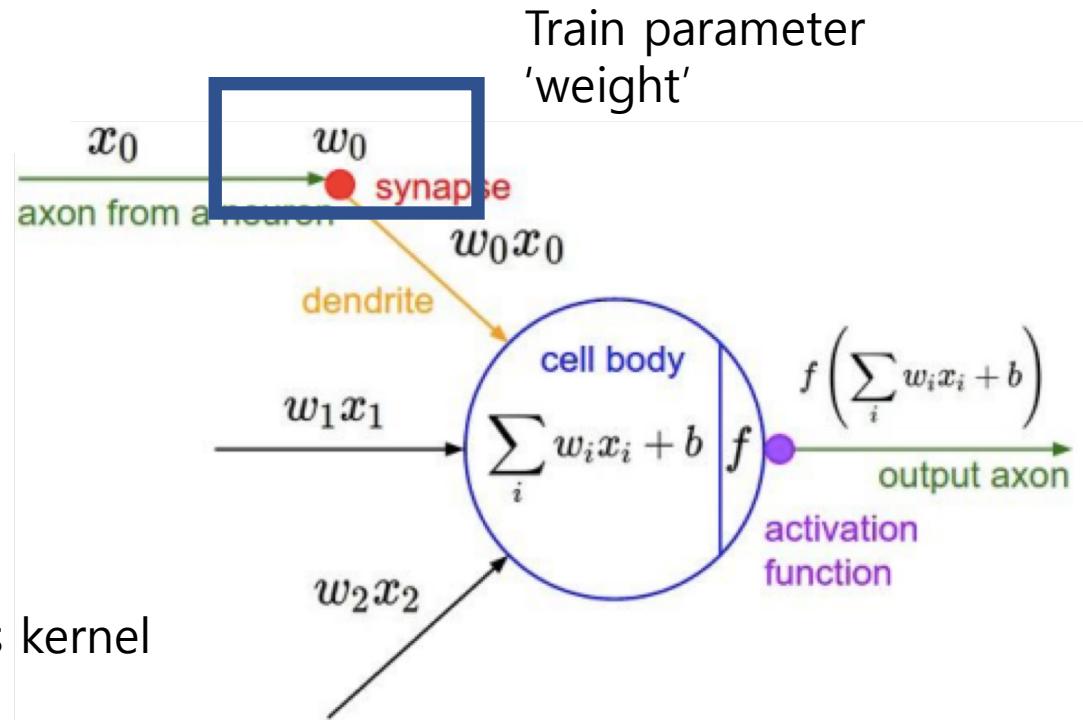
1 <small><math>\times 1</math></small>	1 <small><math>\times 0</math></small>	1 <small><math>\times 1</math></small>	0	0
0 <small><math>\times 0</math></small>	1 <small><math>\times 1</math></small>	1 <small><math>\times 0</math></small>	1	0
0 <small><math>\times 1</math></small>	0 <small><math>\times 0</math></small>	1 <small><math>\times 1</math></small>	1	1
0	0	1	1	0
0	1	1	0	0

Image



Convolved Feature

Train this kernel



# Convolution Neural Net (CNN)

## Convolution Layer

1 <small>x1</small>	1 <small>x0</small>	1 <small>x1</small>	0	0
0 <small>x0</small>	1 <small>x1</small>	1 <small>x0</small>	1	0
0 <small>x1</small>	0 <small>x0</small>	1 <small>x1</small>	1	1
0	0	1	1	0
0	1	1	0	0

Image

Convolved  
Feature

4		

= Feature Map

# Convolution Neural Net (CNN)

## Convolution Layer

1 <small>x1</small>	1 <small>x0</small>	1 <small>x1</small>	0	0
0 <small>x0</small>	1 <small>x1</small>	1 <small>x0</small>	1	0
0 <small>x1</small>	0 <small>x0</small>	1 <small>x1</small>	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved Feature

Feature Map x Activation function (Relu)  
= Activate Map (Convolution layer output)

# Convolution Neural Net (CNN)

---

## Pooling Layer

Activation Map

12	20	30	0
8	12	2	0
34	70	37	7
112	100	22	12

Max Pooling

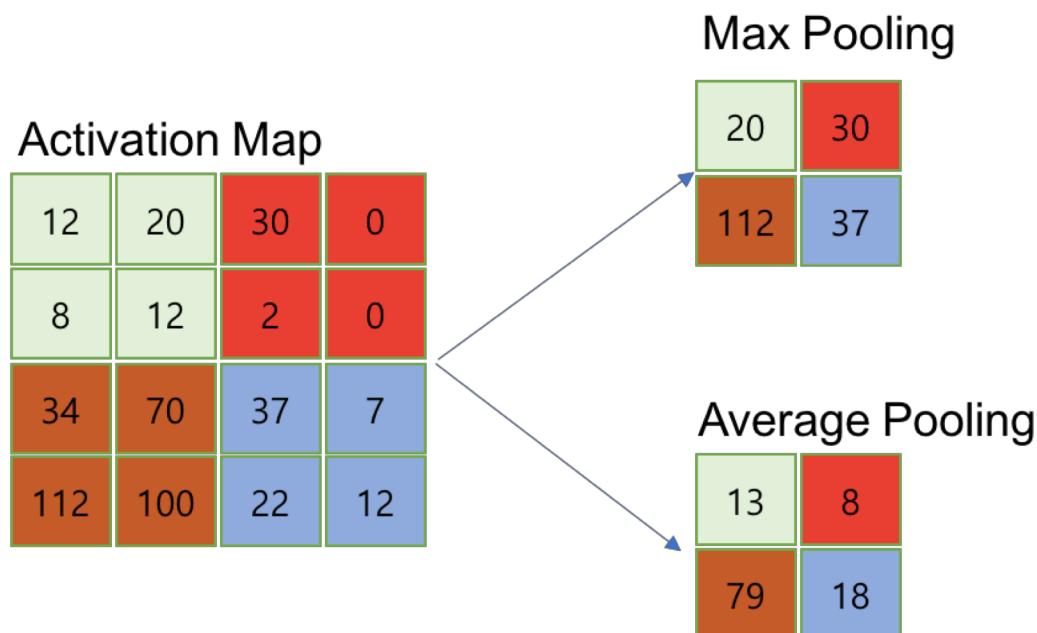
20	30
112	37

Average Pooling

13	8
79	18

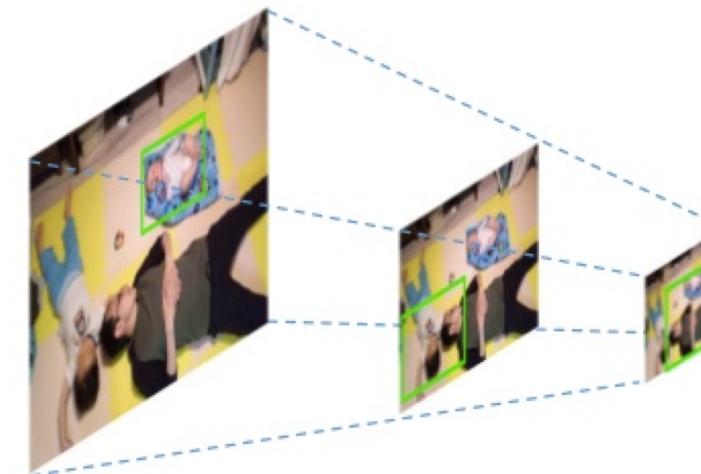
# Convolution Neural Net (CNN)

## Why Pooling?



## Multi-Scaling (Pyramid Pooling)

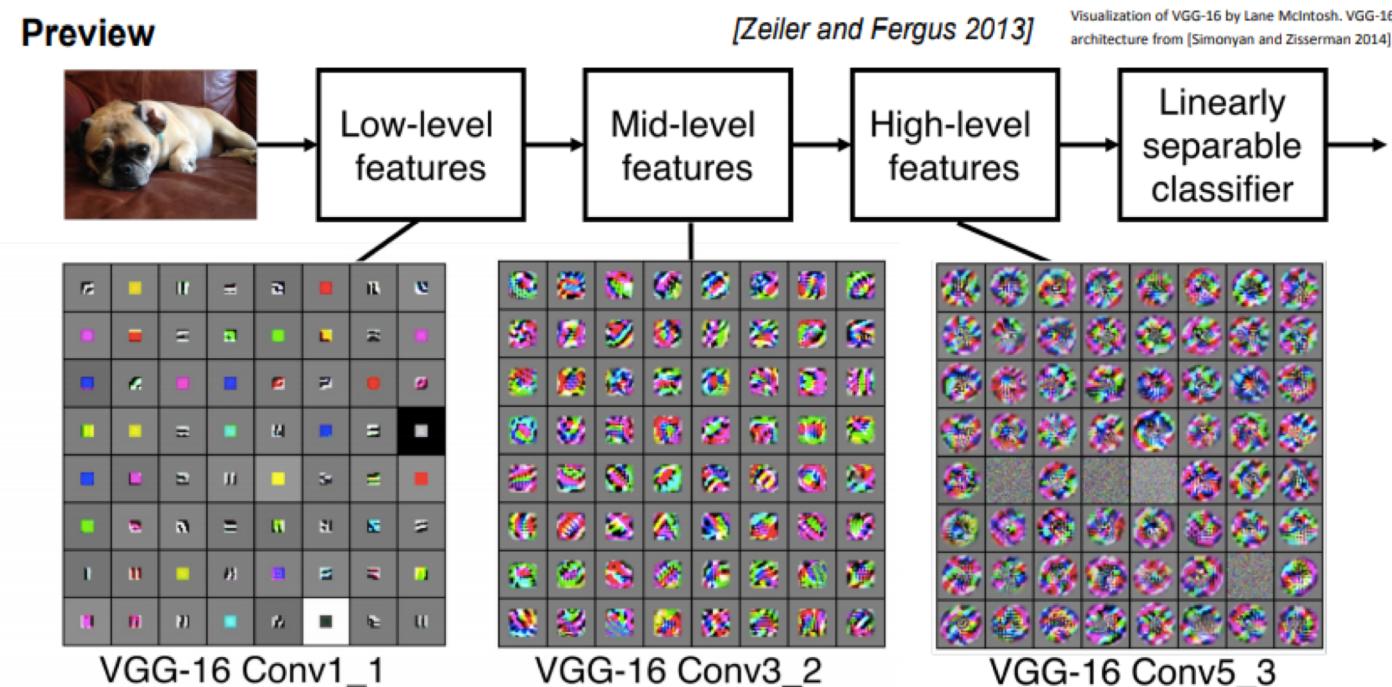
- Find objects in images at different scales
- Combined with a **sliding window**, it can find objects in various locations **with the same window size**



11

# Convolution Neural Net (CNN)

## CNN Architecture



# AlexNet

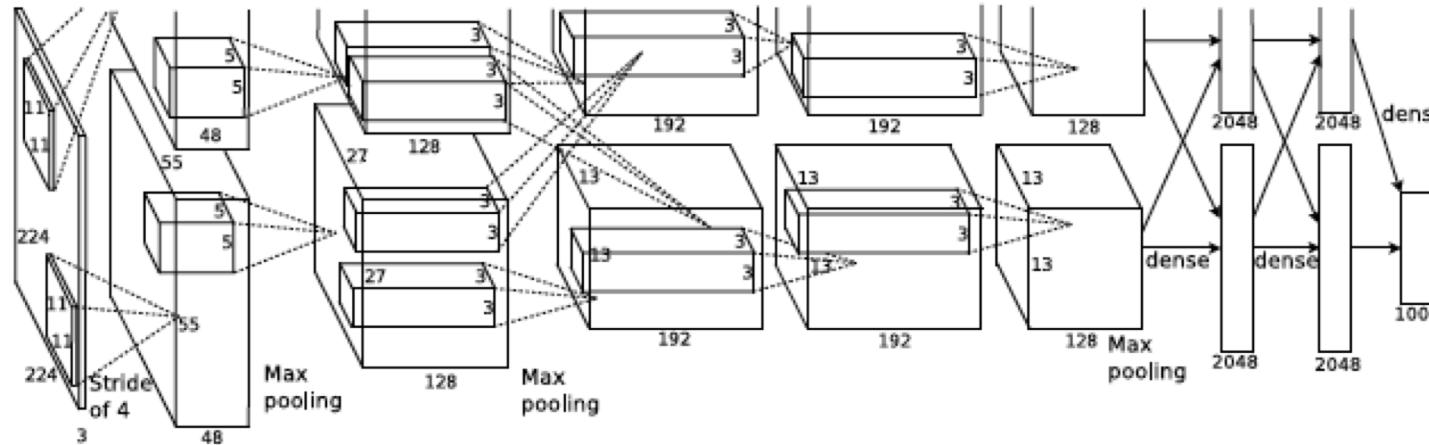


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input 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.

## AlexNet

- CNN Architecture
- Activation function 'Relu'
- GPU
- Preventing overfitting
  - Dropout
  - Data augmentation

# AlexNet

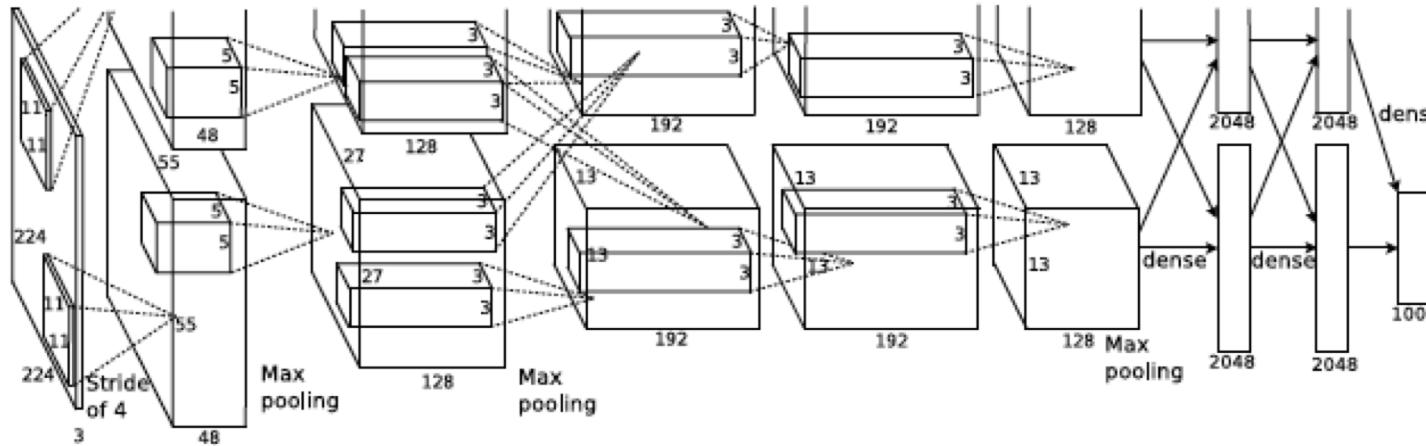


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input 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.

## AlexNet Architecture

- 5 Convolution+pooling Layer
- 3 Fully Connected Layer
- Two-GPU net  
because of hardware limitation  
at the time

# Why Deep learning

Model	Top-1 (val)	Top-5 (val)	Top-5 (test)
<i>SIFT + FVs [7]</i>	—	—	26.2%
1 CNN	40.7%	18.2%	
5 CNNs	38.1%	16.4%	<b>16.4%</b>
1 CNN*	39.0%	16.6%	—
7 CNNs*	36.7%	15.4%	<b>15.3%</b>

Table 2: Comparison of error rates on ILSVRC-2012 validation and test sets. In *italics* are best results achieved by others. Models with an asterisk\* were “pre-trained” to classify the entire ImageNet 2011 Fall release. See Section 6 for details.

Pattern recognition  
&  
machine learning

10.9% !

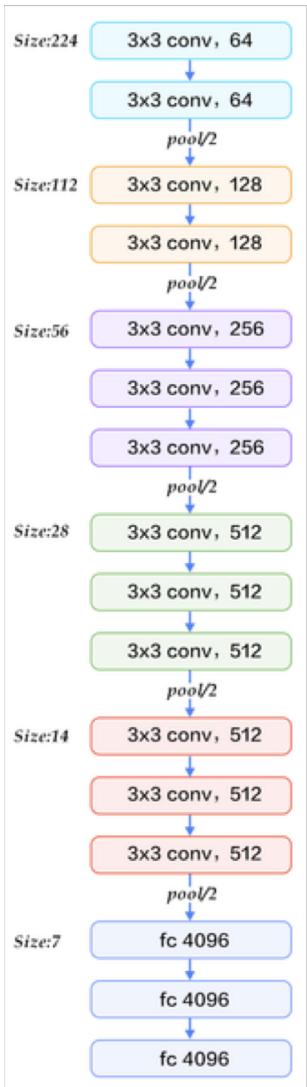
Deep learning

## Top 5 test error

The Top-5 error rate is the percentage of test examples for which the correct class was not in the top 5 predicted classes.

# VGG

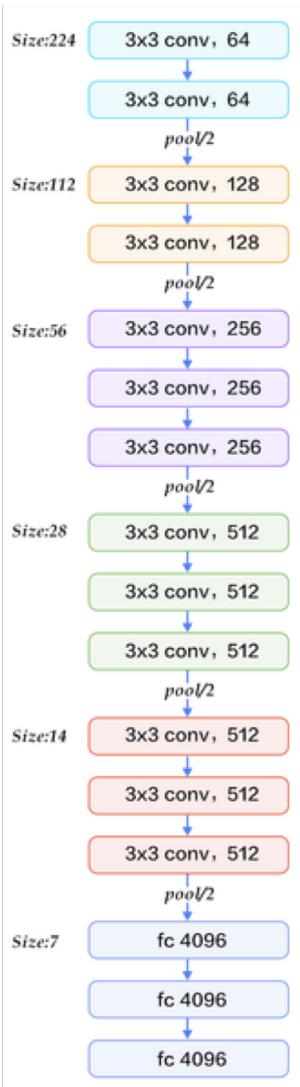
---



- *Very Deep Convolutional Networks For Large Scale Image Recognition - Karen Simonyan and Andrew Zisserman; 2015*
- ILSVRC'14 2nd in classification, 1st in localization
- Smaller convolution filter
- Deeper than previous network

# VGG

---



- **Smaller convolution filter**
  - AlexNet (2012)  
filter size = (11 x 11), stride = 4
  - ZFNet (2013)  
filter size = (7 x 7), stride = 2
- **Deeper layer**
  - AlexNet (2012)  
layer depth = 8 layers

# Architecture

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input ( $224 \times 224$ RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256	conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
softmax					

- **Convolution layer (blue)**
  - kernel size = (3x3), stride = 1
- **Max pooling layer (red)**
  - kernel size = (2x2), stride = 2
- **Convolution layer + max pooling layer**
  - get spatial resolution for the feature maps

# Architecture

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input ( $224 \times 224$ RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

- **Activation function = 'Relu'**
  - for non-linearity
- **Padding = 1 pixel (convolution layer)**
- **No LRN (Local Response Normalization)**
  - not improve the performance
  - increase memory consumption and computation time
- **Parameter**

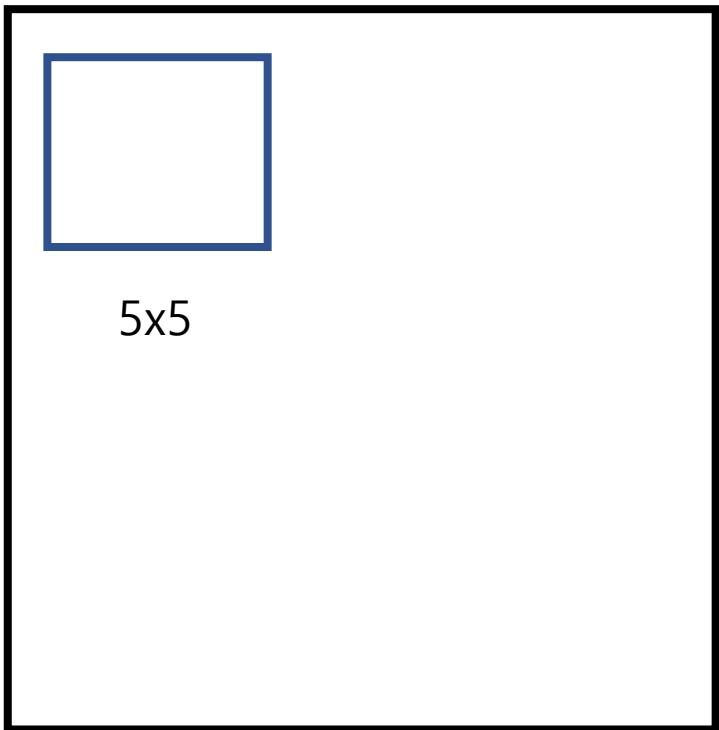
Table 2: Number of parameters (in millions).

Network	A,A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144

## Smaller Convolution filter

---

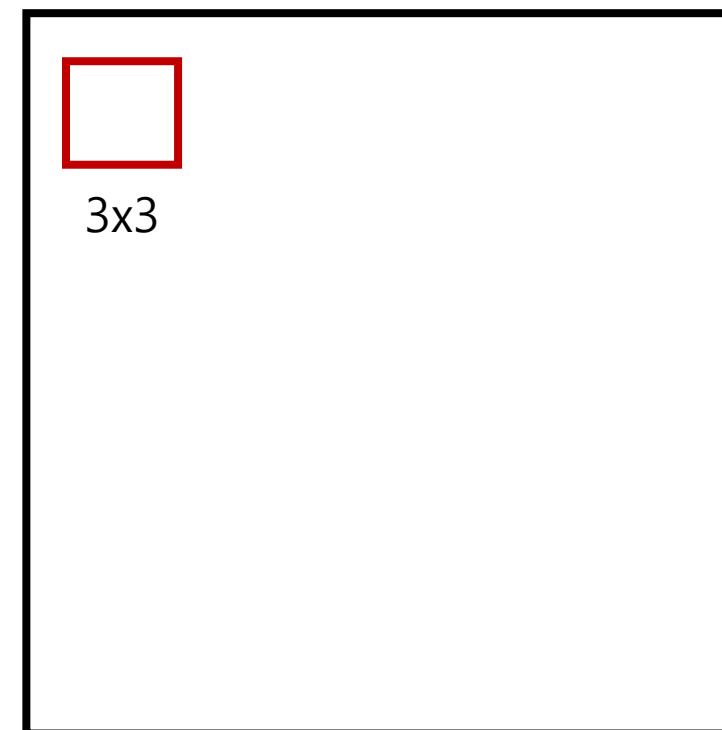
No padding, stride = 1



$28 \times 28$

Output = (23x23)

No padding, stride = 1



$28 \times 28$

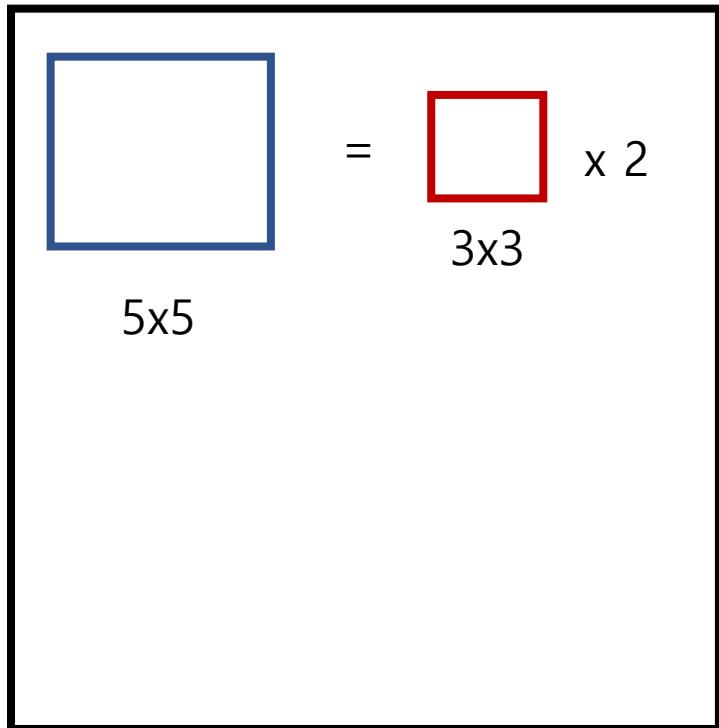
Output = (25x25)

$3 \times 3$

# Smaller Convolution filter

---

No padding, stride = 1



## Why use smaller filter instead of bigger one

- Makes the decision function more discriminative (increase non-linearity)
- Decrease the number of parameter

The diagram compares the number of parameters between different kernel sizes. On the left, a blue-bordered square represents a 5x5 kernel, labeled "5x5". To its right is the equation  $= 5 \times 5(25)$ . To the right of that is a red-bordered square representing a 3x3 kernel, labeled "3x3". To its right is the equation  $= 2 \times (3 \times 3) (18)$ .

# Deeper layer

## Vanishing gradient (Fine tuning)

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

## Use pretrained 'A' VGG model

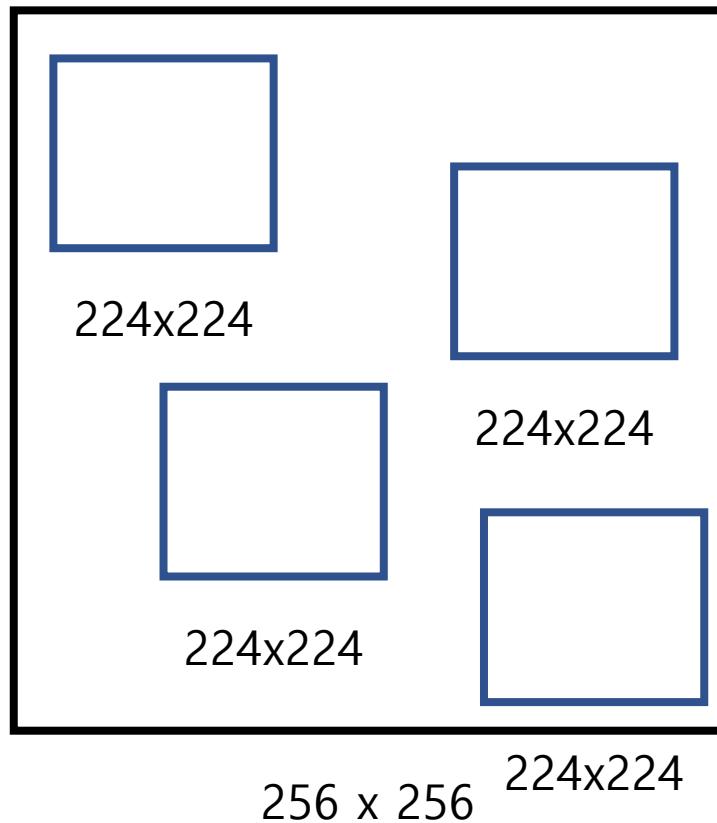
- Net 'A' is simple model than other networks.
- initialized weights in first 4 convolution layers and last 3 FC layers using pretrained net 'A'

# Train

---

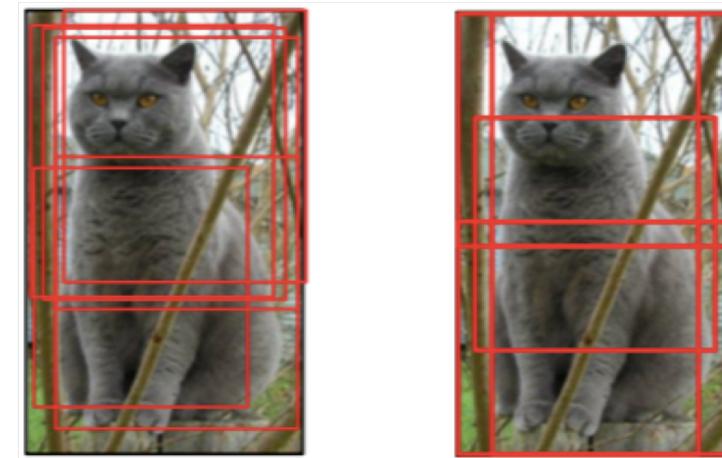
## Data augmentation

Random Cropping!



**Random horizontal flipping & RGB color shift  
(Same as Alexnet's augmentation strategy)**

- extracting random  $224 \times 224$  patches



# Train

---

## Data augmentation

### RGB color shift

(Same as Alexnet's augmentation strategy)

- change RGB color !
- PCA on the set of RGB pixel values
- P1, p2, p3 = eigenvectors from PCA  
lambda1, 2, 3 = eigenvalues from PCA  
a1, a2, a3 = Gaussian noise (0, 0.1)

$$I_{xy} = [I_{xy}^R, I_{xy}^G, I_{xy}^B]^T \quad [\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3][\alpha_1 \lambda_1, \alpha_2 \lambda_2, \alpha_3 \lambda_3]^T$$
$$\alpha_i \sim N(0, 0.1)$$

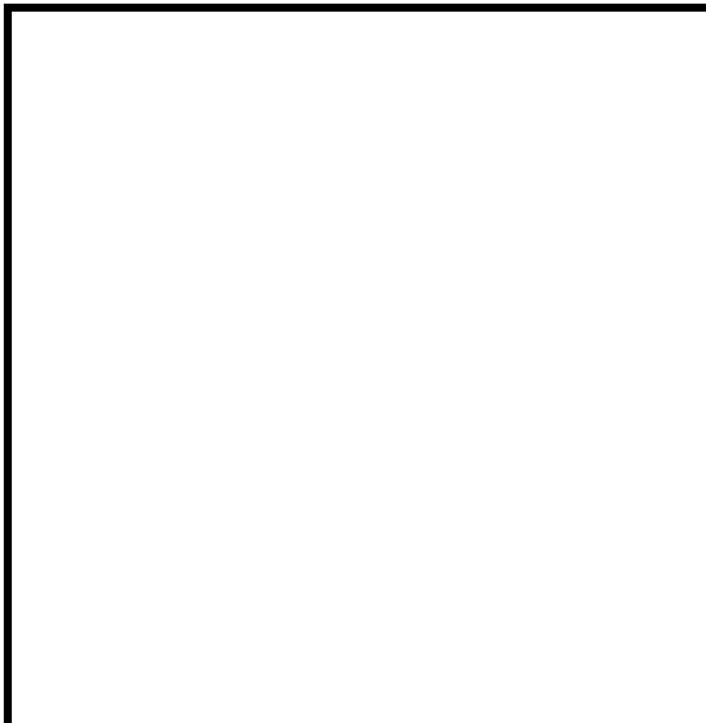
# Train

---

## Training scale



Cropping size is determined by 'S'



256 x 256

- 'S' is 'train scale' value
- ' $S \geq 224$ '
  - If  $S == '224'$  :  
Crop whole image
  - Elif  $S > '224'$  :  
Crop small part of image

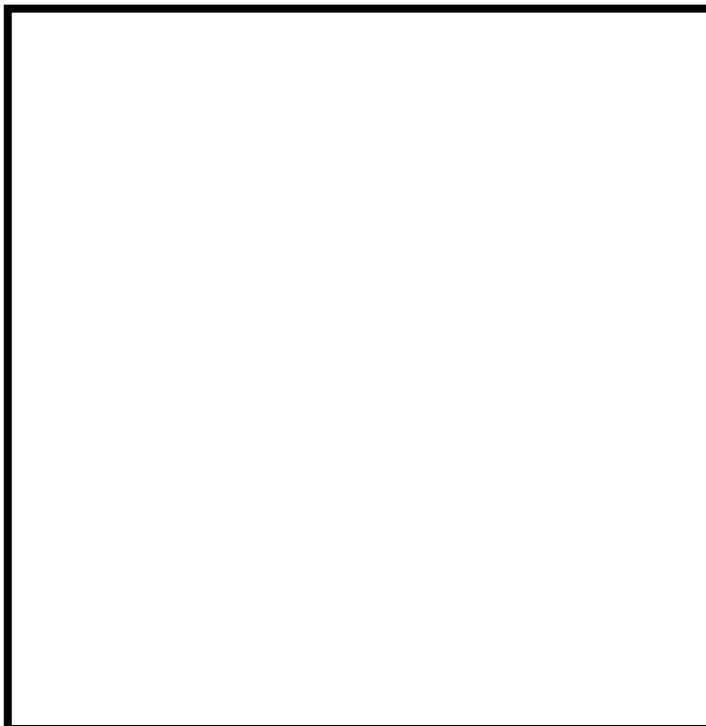
# Train

---

## Training scale



Cropping size is determined by 'S'



256 x 256

- 'S' is 'train scale' value
- ' $S \geq 224$ '
  - If  $S == '224'$  :  
Crop whole image
  - Elif  $S > '224'$  :  
Crop small part of image

