

Cheat Sheet – Famous CNNs

AlexNet – 2012

Why: AlexNet was born out of the need to improve the results of the ImageNet challenge.

What: The network consists of 5 Convolutional (CONV) layers and 3 Fully Connected (FC) layers. The activation used is the Rectified Linear Unit (ReLU).

How: Data augmentation is carried out to reduce over-fitting, Uses Local response localization.

AlexNet Network - Structural Details													
Input	Output	Layer	Stride	Pad	Kernel size	in	out	# of Param					
227 227 3	55 55 96	conv1	4	0	11 11	3	96	34944					
55 55 96	27 27 96	maxpool1	2	0	3 3	96	96	0					
27 27 96	27 27 256	conv2	1	2	5 5	96	256	614656					
27 27 256	13 13 256	maxpool2	2	0	3 3	256	256	0					
13 13 256	13 13 384	conv3	1	1	3 3	256	384	885120					
13 13 384	13 13 384	conv4	1	1	3 3	384	384	1327488					
13 13 384	13 13 256	conv5	1	1	3 3	384	256	884992					
13 13 256	6 6 256	maxpool5	2	0	3 3	256	256	0					
									fc6	1	9216	4096	37752832
									fc7	1	4096	4096	16781312
									fc8	1	4096	1000	4097000
Total													62,378,344

VGGNet – 2014

Why: VGGNet was born out of the need to reduce the # of parameters in the CONV layers and improve on training time

What: There are multiple variants of VGGNet (VGG16, VGG19, etc.)

How: The important point to note here is that all the conv kernels are of size 3x3 and maxpool kernels are of size 2x2 with a stride of two.

VGG16 - Structural Details												
#	Input Image		output	Layer	Stride	Kernel	in	out	Param			
1	224	224	3	conv3-64	1	3 3	3	64	1792			
2	224	224	64	conv3-64	1	3 3	64	64	36928			
3	112	112	64	maxpool	2	2 2	64	64	0			
4	112	112	128	conv3-128	1	3 3	64	128	73856			
5	112	112	128	conv3-128	1	3 3	128	128	147584			
6	56	56	128	maxpool	2	2 2	128	128	65664			
7	56	56	256	conv3-256	1	3 3	128	256	295168			
8	56	56	256	conv3-256	1	3 3	256	256	590080			
9	28	28	256	maxpool	2	2 2	256	256	590080			
10	28	28	512	conv3-512	1	3 3	256	512	1180160			
11	28	28	512	conv3-512	1	3 3	512	512	2359808			
12	14	14	512	conv3-512	1	3 3	512	512	2359808			
13	14	14	512	maxpool	2	2 2	512	512	2359808			
14	14	14	512	conv3-512	1	3 3	512	512	2359808			
15	14	14	512	conv3-512	1	3 3	512	512	2359808			
16	1	1	4096	maxpool	2	2 2	512	512	2359808			
17	1	1	25088	1 4096	fc	1	1	25088	4096	10764544		
18	1	1	4096	1 4096	fc	1	1	4096	1000	16781312		
19	1	1	4096	1 1000	fc	1	1	4096	1000	4097000		
Total										138,423,208		

ResNet – 2015

Why: Neural Networks are notorious for not being able to find a simpler mapping when it exists. ResNet solves that.

What: There are multiple versions of ResNetXX architectures where 'XX' denotes the number of layers. The most used ones are ResNet50 and ResNet101. Since the vanishing gradient problem was taken care of (more about it in the How part), CNN started to get deeper and deeper

How: ResNet architecture makes use of shortcut connections do solve the vanishing gradient problem. The basic building block of ResNet is a Residual block that is repeated throughout the network.

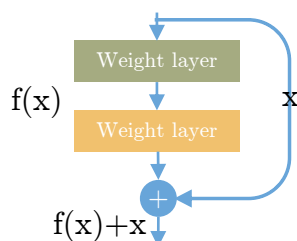


Figure 1 ResNet Block

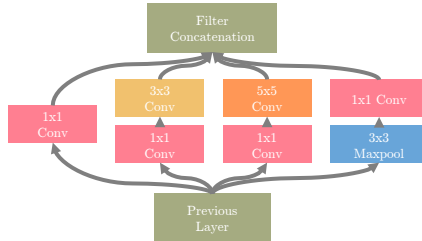


Figure 2 Inception Block

Inception – 2014

Why: Larger kernels are preferred for more global features, on the other hand, smaller kernels provide good results in detecting area-specific features. For effective recognition of such a variable-sized feature, we need kernels of different sizes. That is what Inception does.

What: The Inception network architecture consists of several inception modules of the following structure. Each inception module consists of four operations in parallel, 1x1 conv layer, 3x3 conv layer, 5x5 conv layer, max pooling

How: Inception increases the network space from which the best network is to be chosen via training. Each inception module can capture salient features at different levels.

Comparison					
Network	Year	Salient Feature	top5 accuracy	Parameters	FLOP
AlexNet	2012	Deeper	84.70%	62M	1.5B
VGGNet	2014	Fixed-size kernels	92.30%	138M	19.6B
Inception	2014	Wider - Parallel kernels	93.30%	6.4M	2B
ResNet-152	2015	Shortcut connections	95.51%	60.3M	11B

GoogleNet - Structural Details												
	Input Image	output	Layer	Stride	Pad	Kernel	in	out	Param			
	224 224 3	112 112 64	conv1	2	0.5	3 3	64	64	0			
	112 112 64	56 56 64	maxpool	1	0	1 1	64	64	0			
	56 56 64	56 56 64	conv1	1	0	1 1	64	64	4160			
	56 56 64	56 56 192	conv2	1	0	1 1	64	192	110784			
	56 56 192	28 28 192	maxpool2	2	0.5	3 3	192	192	0			
inception (3a)	28 28 192	28 28 96	conv3a	1	0	1 1	192	96	18528			
	28 28 96	28 28 96	conv3b	1	0	1 1	192	96	3688			
	28 28 192	28 28 64	conv4a	1	0	1 1	192	64	12352			
	28 28 96	28 28 128	conv4b	1	0	1 1	192	128	10720			
	28 28 128	28 28 32	conv5a	1	2	5 5	16	32	12832			
inception (3b)	28 28 192	28 28 32	conv5b	1	0	1 1	192	32	6176			
	28 28 32	28 28 256	depth-concat	1	0	1 1	256	128	32896			
	28 28 128	28 28 32	conv5c	1	0	1 1	256	32	8224			
	28 28 128	28 28 256	maxpool4	1	1	3 3	256	256	0			
	28 28 128	28 28 128	conv5d	1	0	1 1	256	128	32896			
inception (3c)	28 28 128	28 28 128	conv5e	1	0	1 1	256	128	32896			
	28 28 128	28 28 128	conv5f	1	2	5 5	32	96	76896			
	28 28 192	28 28 64	conv5g	1	0	1 1	256	64	18448			
	28 28 128	28 28 480	depth-concat	1	0	1 1	480	480	0			
	28 28 480	14 14 480	maxpool5	2	0.5	3 3	480	480	0			
inception (4a)	14 14 480	14 14 96	conv6a	1	0	1 1	480	96	46176			
	14 14 480	14 14 16	conv6b	1	0	1 1	480	16	7696			
	14 14 480	14 14 192	conv6c	1	0	1 1	480	192	82304			
	14 14 192	14 14 208	conv6d	1	1	3 3	96	208	179920			
	14 14 16	14 14 48	conv6e	1	2	5 5	16	48	19248			
inception (4b)	14 14 192	14 14 64	conv6f	1	0	1 1	480	64	30784			
	14 14 192	14 14 512	depth-concat	1	0	1 1	512	512	57456			
	14 14 512	14 14 24	conv7a	1	0	1 1	512	24	1560			
	14 14 512	14 14 64	conv7b	1	0	1 1	512	64	2800			
	14 14 512	14 14 112	conv7c	1	0	1 1	512	112	22616			
inception (4c)	14 14 112	14 14 160	conv7d	1	0	1 1	64	160	10400			
	14 14 16	14 14 224	conv7e	1	2	5 5	24	224	226016			
	14 14 16	14 14 64	conv7f	1	0	1 1	512	64	34464			
	14 14 160	14 14 64	conv7g	1	0	1 1	64	64	4160			
	14 14 160	14 14 512	depth-concat	1	0	1 1	512	512	65664			
inception (4d)	14 14 512	14 14 24	conv8a	1	0	1 1	512	24	2080			
	14 14 512	14 14 64	conv8b	1	0	1 1	512	64	2800			
	14 14 512	14 14 112	conv8c	1	0	1 1	512	112	22616			
	14 14 112	14 14 160	conv8d	1	0	1 1	64	160	10400			
	14 14 16	14 14 224	conv8e	1	2	5 5	24	224	226016			
inception (4e)	14 14 16	14 14 64	conv8f	1	0	1 1	512	64	34464			
	14 14 160	14 14 64	conv8g	1	0	1 1	64	64	4160			
	14 14 160	14 14 512	depth-concat	1	0	1 1	512	512	65664			
	14 14 512	14 14 24	conv9a	1	0	1 1	512	24	2080			
	14 14 512	14 14 64	conv9b	1	0	1 1	512	64	2800			
inception (5a)	14 14 512	14 14 112	conv9c	1	0	1 1	512	112	22616			
	14 14 112	14 14 160	conv9d	1	0	1 1	64	160	10400			
	14 14 16	14 14 224	conv9e	1	2	5 5	24	224	226016			
	14 14 16	14 14 64	conv9f	1	0	1 1	512	64	34464			
	14 14 160	14 14 64	conv9g	1	0	1 1	64	64	4160			
inception (5b)	14 14 160	14 14 512	depth-concat	1	0	1 1	512	512	65664			
	14 14 512	7 7 832	conv10a	2	0.5	3 3	832	832	0			
	7 7 832	7 7 160	conv10b	1	0	1 1	832	160	133280			
	7 7 832	7 7 832	maxpool6	1	0	1 1	832	832	26560			
	7 7 832	7 7 832	conv10c	1	0	1 1	832	832	212384			
inception (5c)	7 7 96	7 7 160	conv10d	1	0	1 1	160	160	28256			
	7 7 16	7 7 128	conv10e	1	2	5 5	32	128	102528			
	7 7 256	7 7 139	conv10f	1	0	1 1	832	128	106624			
	7 7 139	7 7 7	depth-concat	1	0	1 1	7	7	0			
	7 7 832	7 7 192	depth-concat	1	0	1 1	832	192	159936			
inception (5d)	7 7 832	7 7 48	conv10g	1	0	1 1	832	48	39984			
	7 7 832	7 7 832	maxpool6	1	1	3 3	832	832	0			
	7 7 832	7 7 832	conv10h	1	0	1 1	832	832	311920			
	7 7 96	7 7 384	conv10i	1	1	3 3	96	384	663936			
	7 7 16	7 7 128	conv10j	1	2	5 5	48	128	153728			
	7 7 384	7 7 128	conv10k	1	0	1 1	128	128	16512			
	7 7 1024	1 1 1024	depth-concat	1	0	1 1	1024	1024	0			
	1 1 1024	1 1 1000	fc	1	0	1 1	1024	1000	4025000			