

# Computer vision. Convolutional Neural Networks.

December 2, 2016

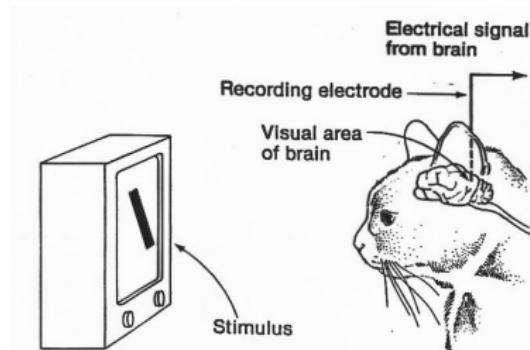
# Agenda

- Quick History of Computer Vision
- CNN Applications
- CNN Layers

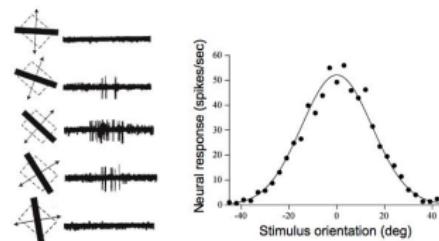
# Credits

- cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson (lectures 6, 7)

# Hubel and Wiesel. Vision research



## V1 physiology: orientation selectivity



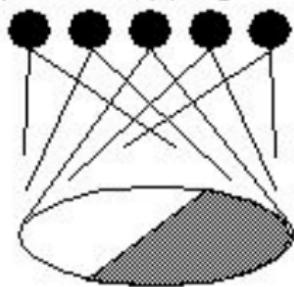
Hubel & Wiesel, 1968

Hubel & Wiesel, 1959, Receptive fields of single neurones in the cat's striate cortex

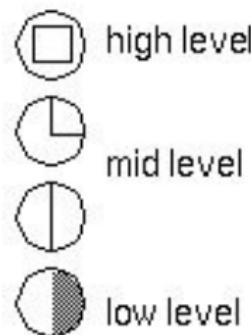
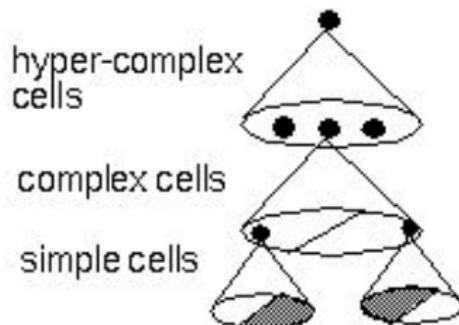
# Hubel and Wiesel. Vision research

## Hubel & Weisel

topographical mapping

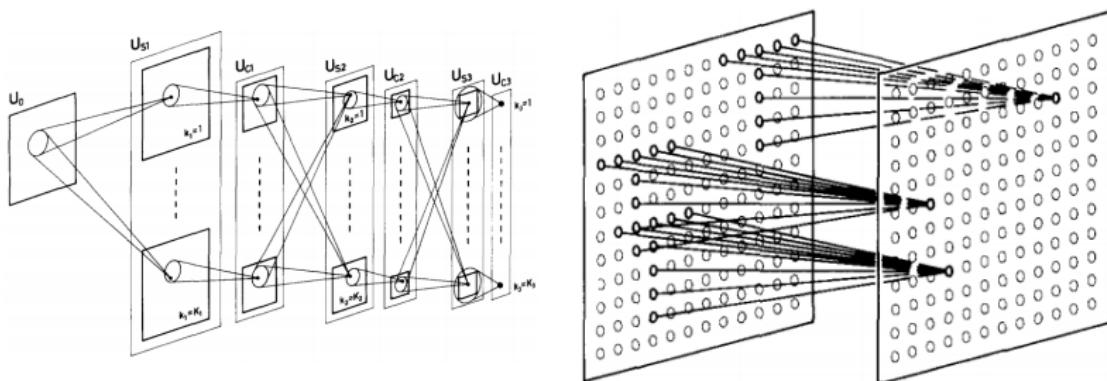


## featural hierarchy



Hubel & Wiesel, 1959, Receptive fields of single neurones in the cat's striate cortex

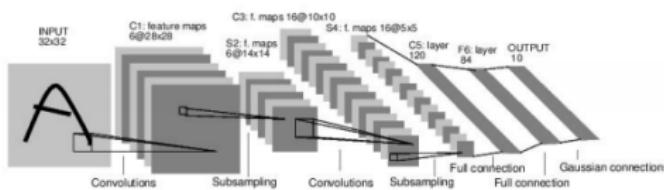
# Neocognitron. Fukushima, 1980



Neocognitron: A self-organizing neural network model

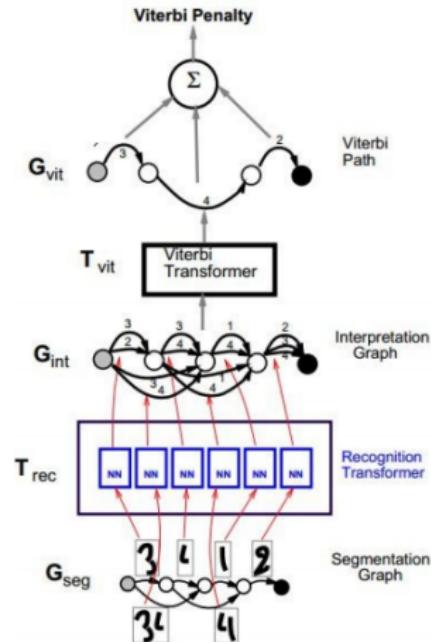
# A bit of history: Gradient-based learning applied to document recognition

[LeCun, Bottou, Bengio, Haffner  
1998]



LeNet-5

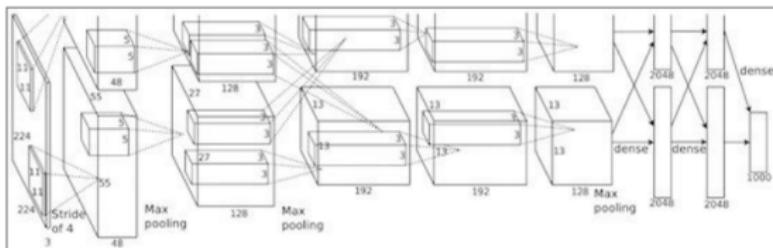
cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7



A bit of history:

## ImageNet Classification with Deep Convolutional Neural Networks

[Krizhevsky, Sutskever, Hinton, 2012]

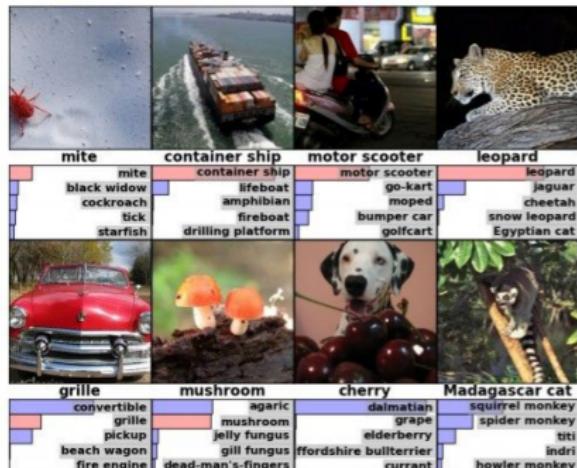


“AlexNet”

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

# Fast-forward to today: ConvNets are everywhere

Classification



Retrieval

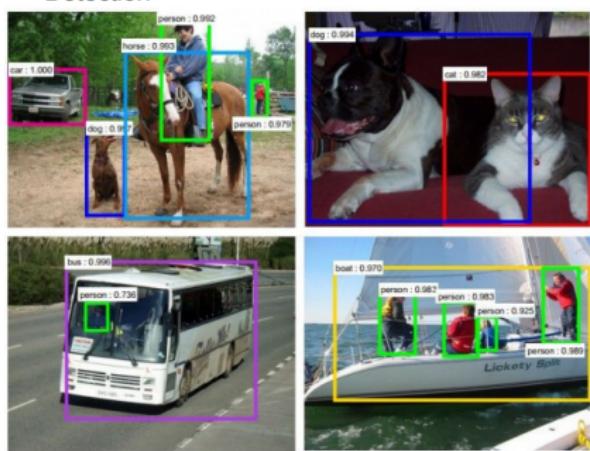


[Krizhevsky 2012]

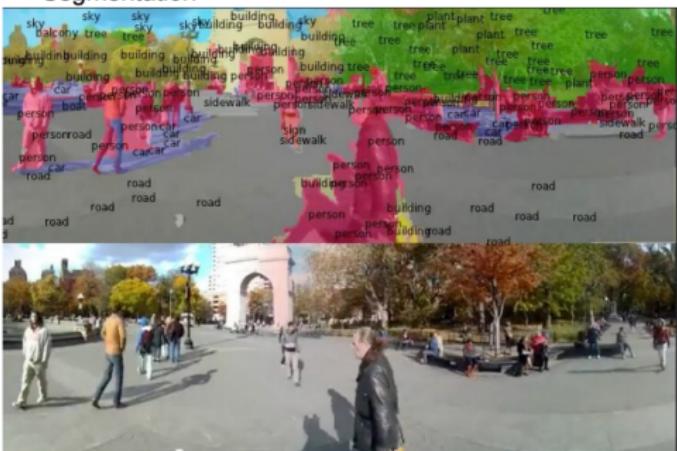
cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

# Fast-forward to today: ConvNets are everywhere

Detection



Segmentation



[Faster R-CNN: Ren, He, Girshick, Sun 2015]

[Farabet et al., 2012]

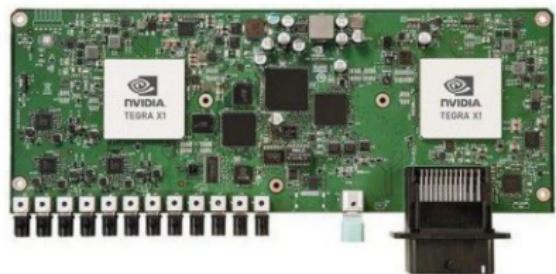
cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

# Fast-forward to today: ConvNets are everywhere



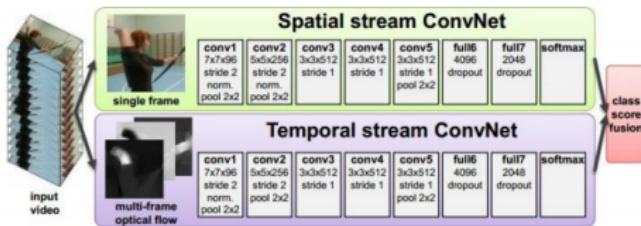
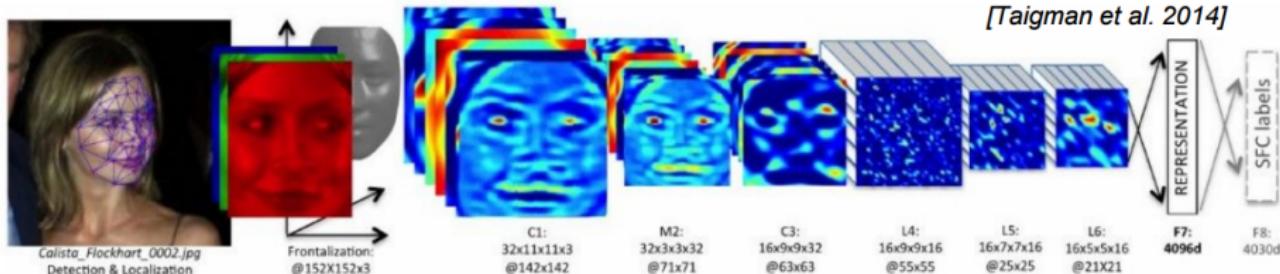
self-driving cars

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7



NVIDIA Tegra X1

# Fast-forward to today: ConvNets are everywhere



[Simonyan et al. 2014]

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7



[Goodfellow 2014]

# Fast-forward to today: ConvNets are everywhere



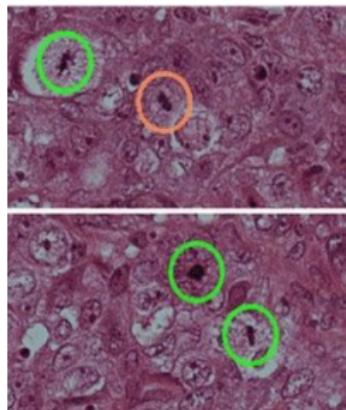
[Toshev, Szegedy 2014]



[Mnih 2013]

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

# Fast-forward to today: ConvNets are everywhere



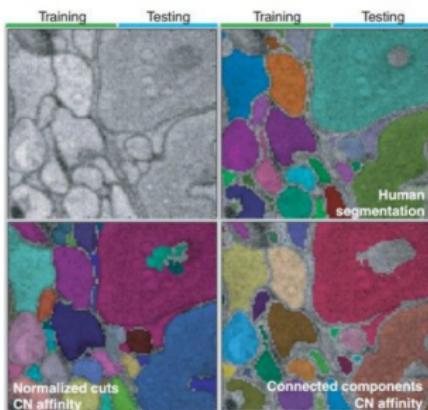
[Ciresan et al. 2013]

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

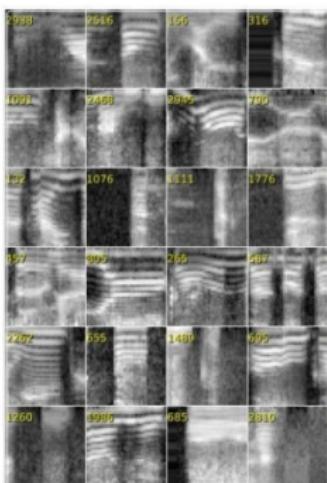


[Sermanet et al. 2011]  
[Ciresan et al.]

## Fast-forward to today: ConvNets are everywhere



[Turaga et al., 2010]



I caught this movie on the Sci-Fi channel recently. It actually turned out to be pretty decent as far as B-list horror/suspense films go. The gore gets tame after the first 1/3, but the suspense factor is still there. I would say it's a **friday night slasher** hybrid decide to play cat-and-mouse with them. Those are further complicated when they pick up a ridiculously whiney bitchslap. What makes this film unique is that the combination of comedy and terror works well in this movie, unlike so many others. The two girls are likable enough and there are some good chase/suspense scenes. Nice pacing and comic timing make this movie more than just feasible for the household here. Definitely worth checking out.

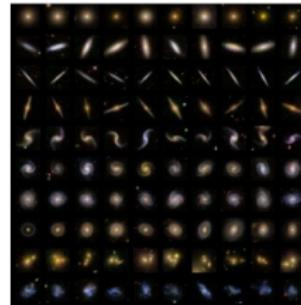
I just saw this on a local independent station in the New York City area. The show deserved attention but where I live the director, George Cukor, *Citizen Kane*, *Wings*, *Gaslight*, had enough energy, it was every bit as bad, every for his problems and stupid as ever. George Cukor's movie *Gaslight* has like a simple man's Michael Bay - with all the explosions, but accumulates promises. There is no point to the conspiracy, no burning budget crisis, the war against the people. We are left to ourselves to connect the dots from one bit of graffiti on various walls to the point to the next. Thus, the burning budget crisis, the war against the people. We are left to ourselves to connect the dots from one bit of graffiti on various walls to the point to the next. Thus, the burning budget crisis, the war against the people. In long, long, Islamic extension, the fate of social security, 47 million Americans without health care, stagnating wages, and the death of the middle class are all subsumed by the sheer tenor of graffiti). A truly, stunningly idiotic film.

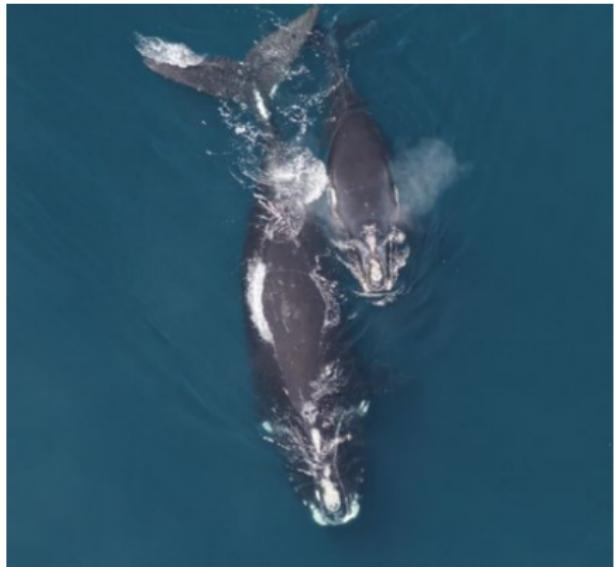
Graphics is far from the best part of the game. **THIS IS THE NUMBER ONE BEST TH GAME IN THE SERIES.** Next to Underground, it deserves strong love. It is at least **good**. There are massive levels, massive unlockable characters... it's just a massive game. **MORE THAN MONEY OR ROLLING GEAR,** this is the kind of money that **ISN'T SPENT PROPERLY.** And even though graphics suck, that doesn't make a game good. Actually, the graphics were great at the time. Today the graphics are crappo. **WHO CARES?** As they say in Canada, This is the fun part, ya! You get to go to Canada in THPS3! Well, I don't know if they say that, but they might, who knows. Well, Canadian people do. Want a remote? I'm going to rip tape. This game rocks. Buy it, play it, enjoy it, love it. AWE-SOME BRILLIANCE.

The first was good and original. I was a not bad home-theater movie. So I heard a second one was made and I had to watch it. What really makes this movie work is Jud Nelson's character and the sometimes clever script. A pretty good script for a person who wrote the Final Destination movies and the direction was okay. Sometimes there's scenes where it looks like it was filmed using a home video camera with a grainy feel. Great scale - for a TV movie. It was worth the rental and probably worth buying just to get that nice score and watching Jud Nelson's Stanley doing what he does best! I suggest new viewers to watch the first one before watching the sequel, just so you'll have an idea what Stanley is like and get a little history background.

Bennet et al. 2011

[Denil et al. 2014]





*Whale recognition, Kaggle Challenge*

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7



*Mnih and Hinton, 2010*

# Image Captioning

Describes without errors	Describes with minor errors	Somewhat related to the image	Unrelated to the image
			
			
			

[Vinyals et al., 2015]

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

# CV Challenges



08	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	04	56	62	00
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	44	04	56	62	00		
81	49	31	73	55	79	14	29	93	71	40	67	57	83	30	03	49	13	36	65		
52	70	95	23	04	60	11	42	68	44	08	56	01	32	56	71	37	02	36	91		
22	31	16	71	51	63	03	89	41	92	36	54	22	40	40	20	66	33	13	80		
24	47	15	60	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50		
32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70		
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21		
24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72		
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33	95		
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92		
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57		
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58		
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40		
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	27	98	66		
03	44	68	87	57	62	20	72	03	46	33	67	46	55	12	32	63	93	53	69		
04	42	16	73	33	75	39	11	24	94	72	18	08	46	29	32	40	62	76	36		
20	69	36	41	72	30	23	88	31	75	89	69	82	67	59	85	74	04	36	16		
20	73	35	29	78	31	90	01	74	31	49	71	46	04	41	16	23	57	05	54		
01	70	54	71	83	51	54	69	16	92	33	48	61	43	52	01	89	23	00	48		

What the computer sees

→  
image classification  
82% cat  
15% dog  
2% hat  
1% mug

<http://cs231n.github.io/classification/>

# CV Challenges

Viewpoint variation



Scale variation



Deformation



Occlusion



Illumination conditions



Background clutter

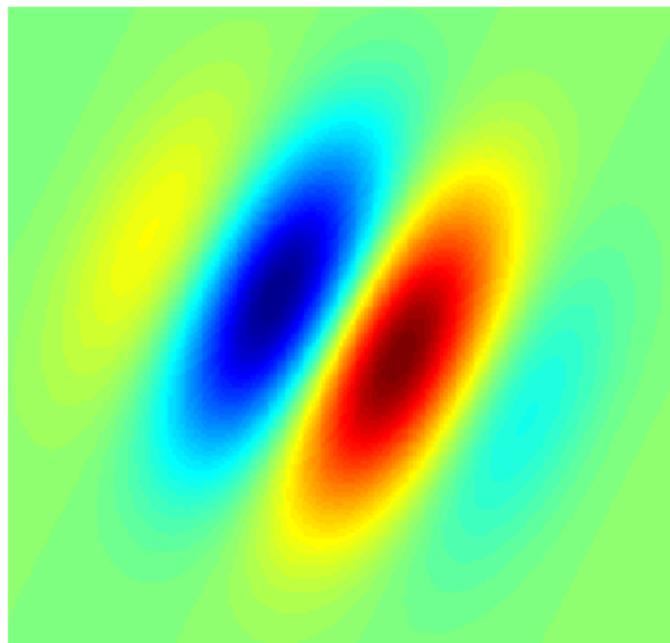


Intra-class variation



<http://cs231n.github.io/classification/>

# Gabor filters



[https://en.wikipedia.org/wiki/Gabor\\_filter](https://en.wikipedia.org/wiki/Gabor_filter)

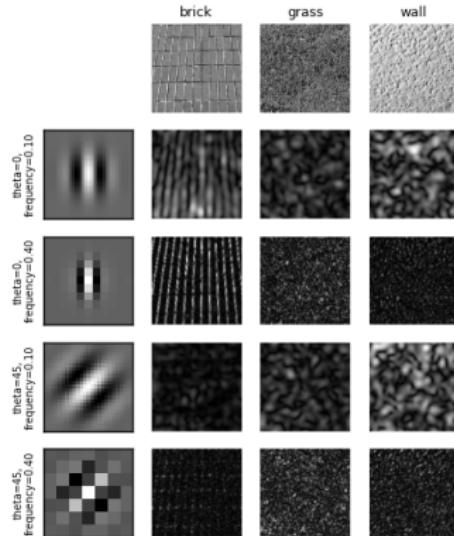
# Gabor filters



<https://cvtuts.wordpress.com/2014/04/27/gabor-filters-a-practical-overview/>

# Gabor filters

Image responses for Gabor filter kernels



[http://scikit-image.org/docs/dev/auto\\_examples/features\\_detection/plot\\_gabor.html#sphx-glr-auto-examples-features-detection-plot-gabor-py](http://scikit-image.org/docs/dev/auto_examples/features_detection/plot_gabor.html#sphx-glr-auto-examples-features-detection-plot-gabor-py)

# HOG



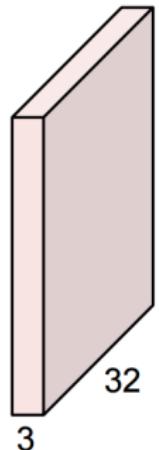
Histogram of Oriented Gradients



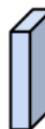
[http://scikit-image.org/docs/dev/auto\\_examples/features\\_detection/plot\\_hog.html#sphx-glr-auto-examples-features-detection-plot-hog-py](http://scikit-image.org/docs/dev/auto_examples/features_detection/plot_hog.html#sphx-glr-auto-examples-features-detection-plot-hog-py)

# Convolution Layer

32x32x3 image



5x5x3 filter

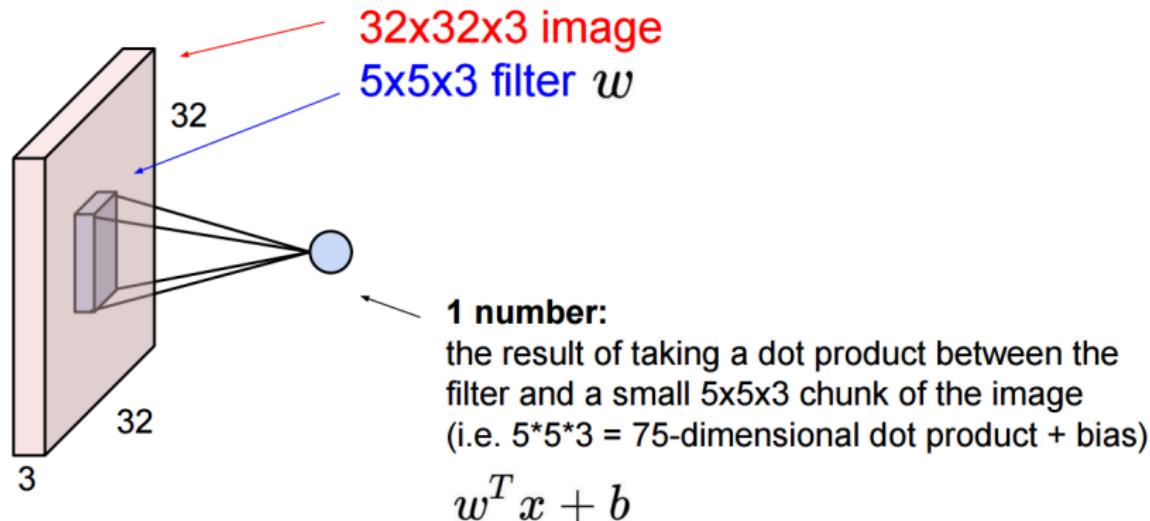


Filters always extend the full depth of the input volume

**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

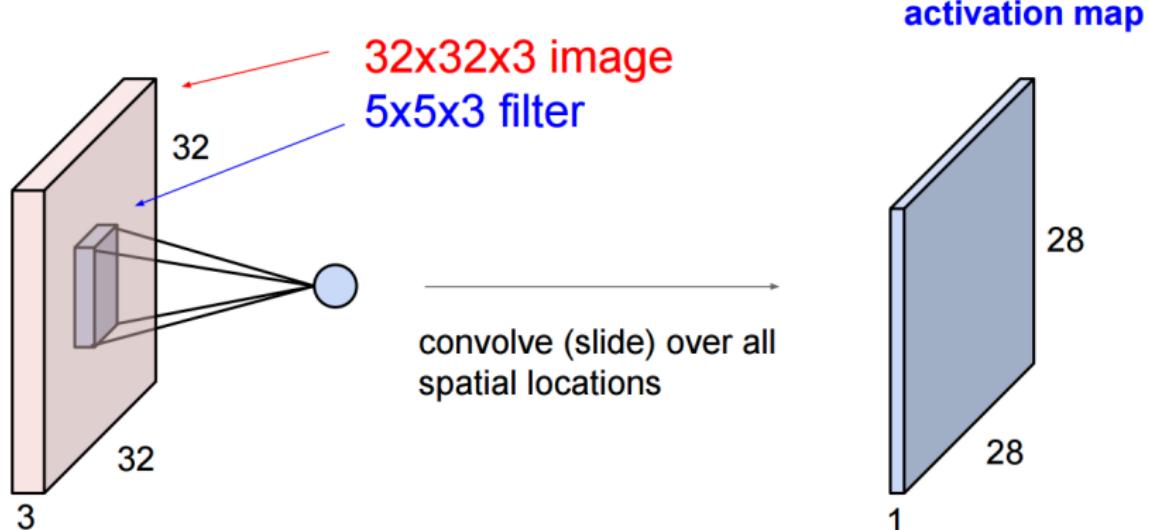
cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

# Convolution Layer



cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

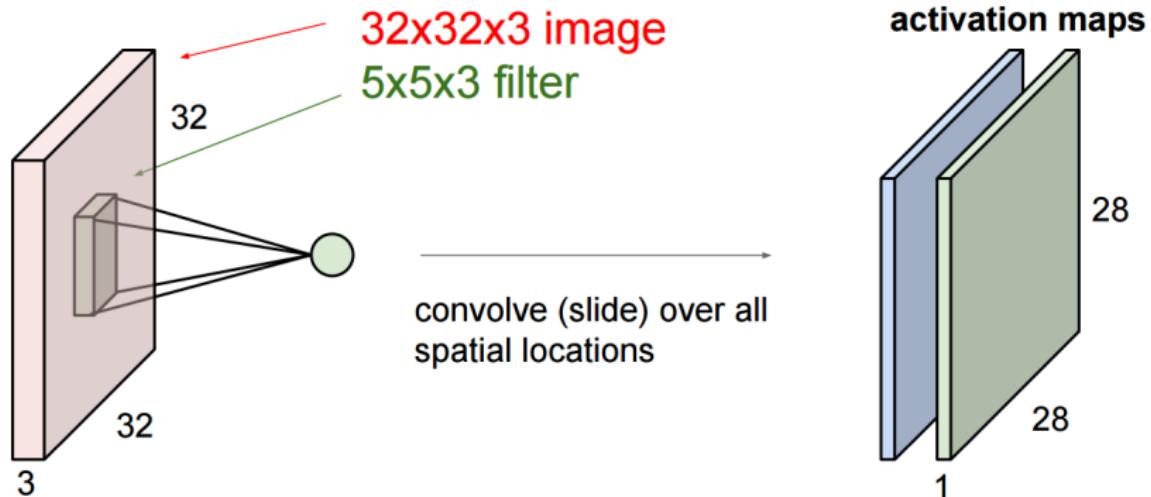
# Convolution Layer



cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

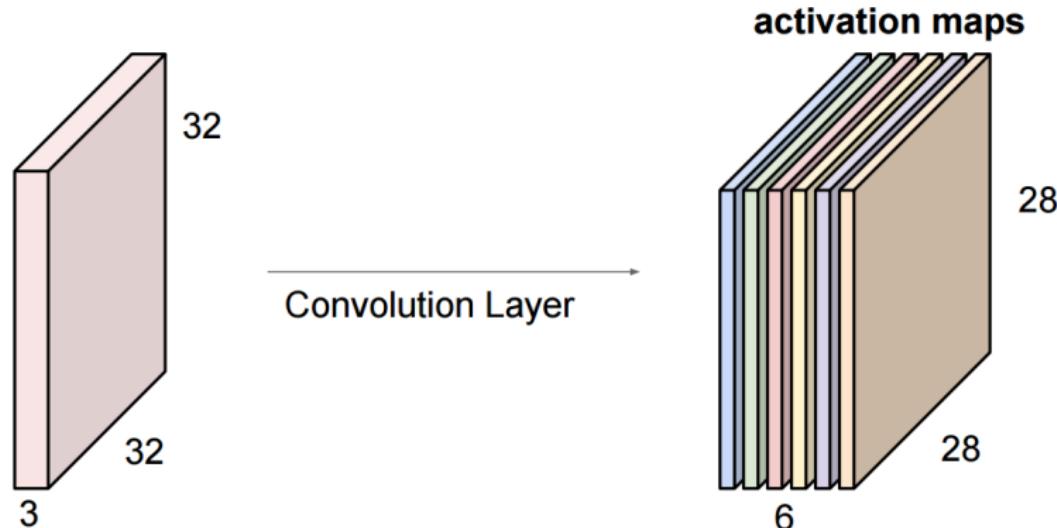
# Convolution Layer

consider a second, green filter



cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

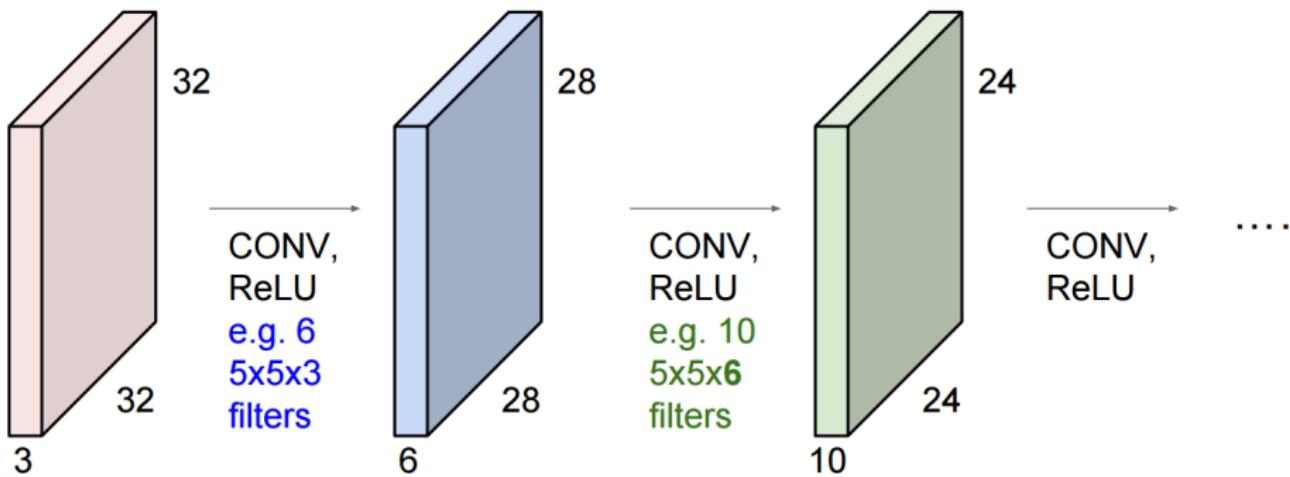
For example, if we had 6  $5 \times 5$  filters, we'll get 6 separate activation maps:



We stack these up to get a “new image” of size  $28 \times 28 \times 6$ !

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions

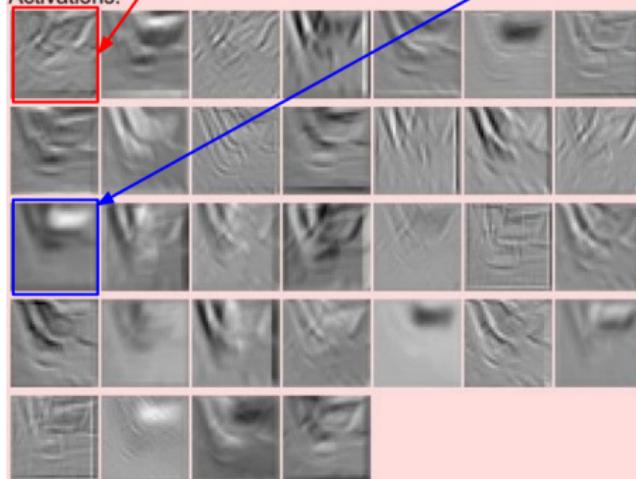


cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7



one filter =>  
one activation map

Activations:



example 5x5 filters  
(32 total)

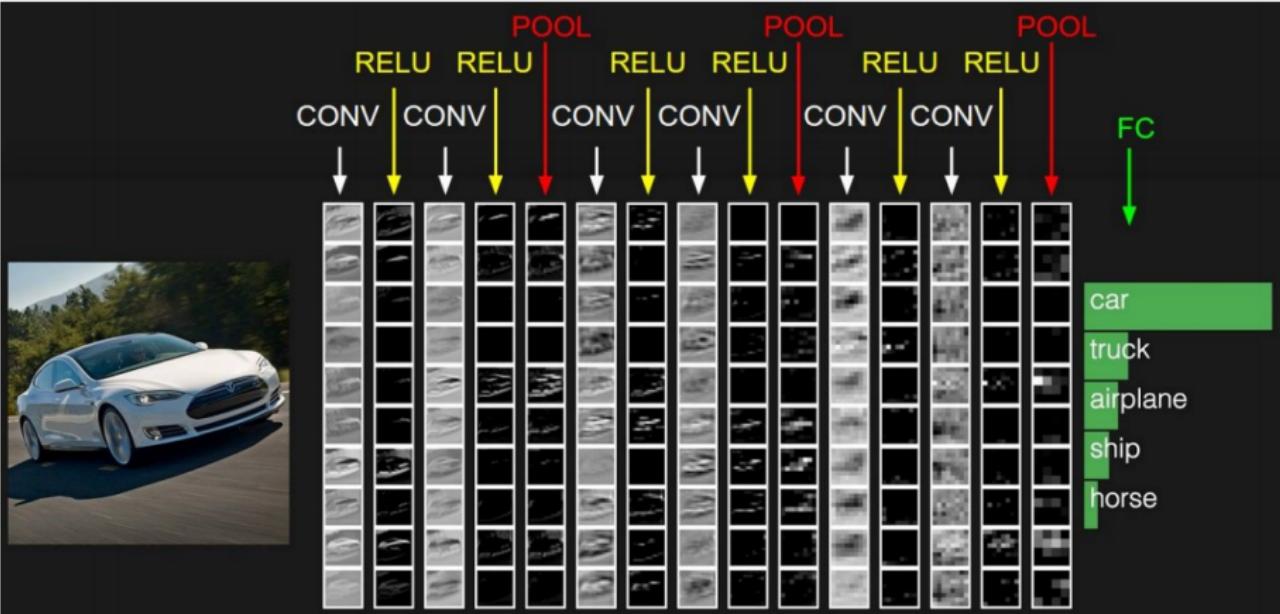
We call the layer convolutional  
because it is related to convolution  
of two signals:

$$f[x,y] * g[x,y] = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1, n_2] \cdot g[x-n_1, y-n_2]$$



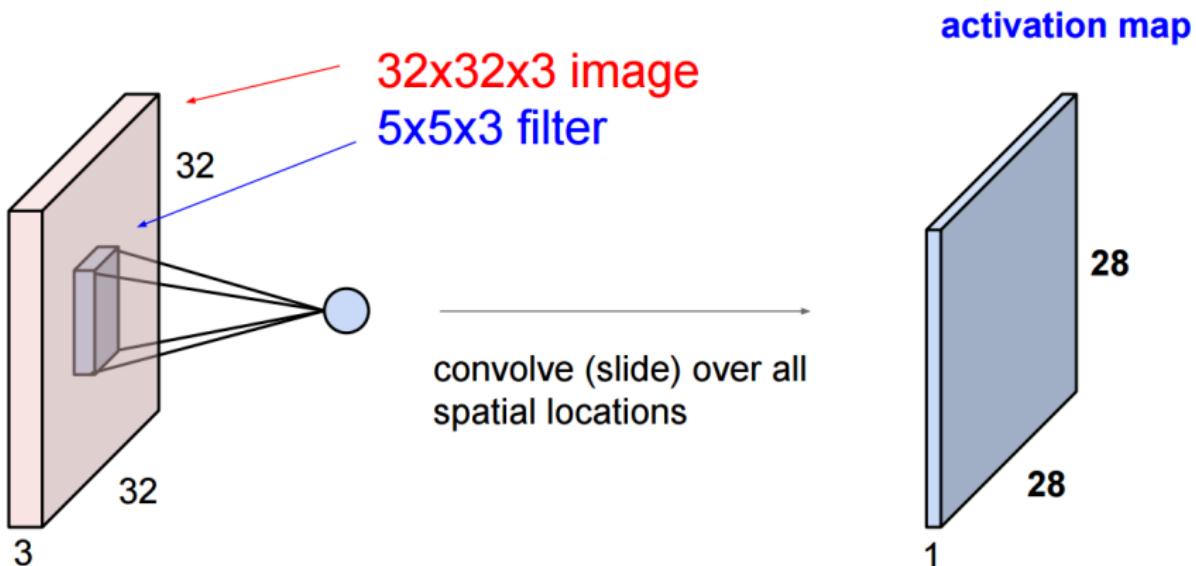
elementwise multiplication and sum of  
a filter and the signal (image)

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7



cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

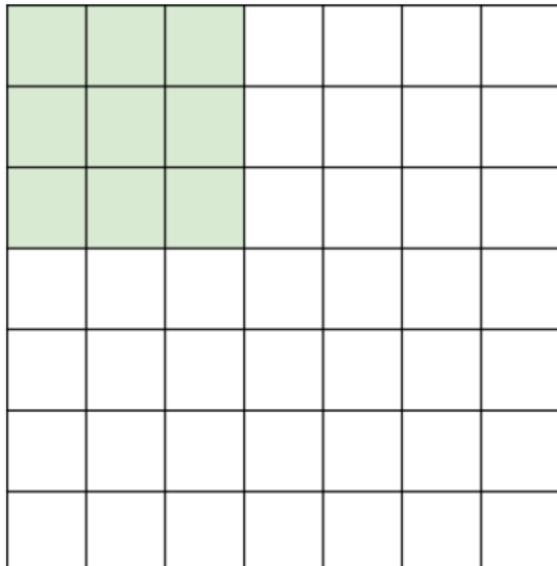
A closer look at spatial dimensions:



cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

## A closer look at spatial dimensions:

7

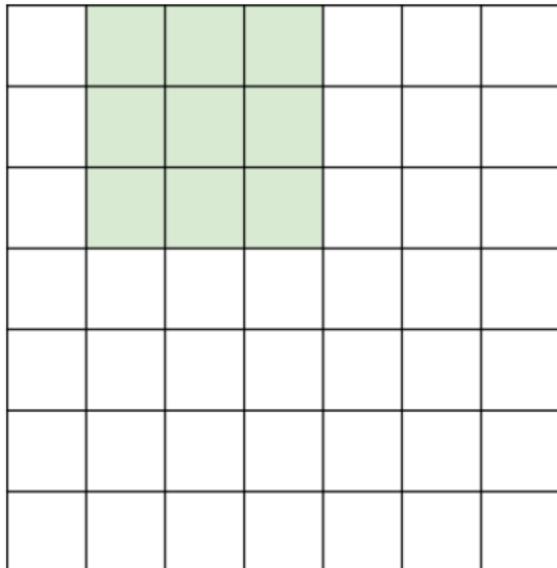


7x7 input (spatially)  
assume 3x3 filter

7

## A closer look at spatial dimensions:

7

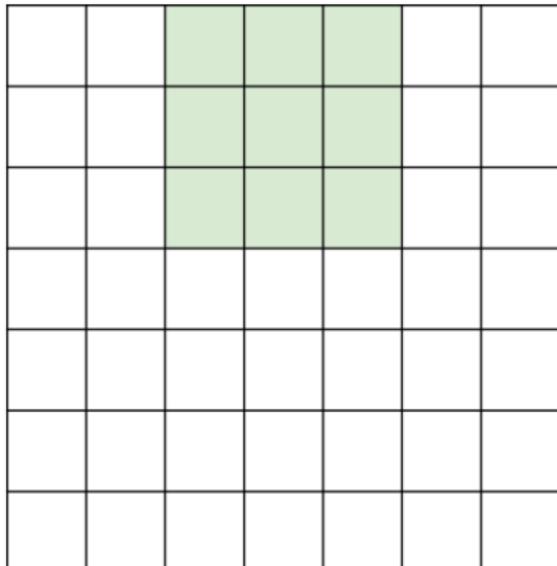


7x7 input (spatially)  
assume 3x3 filter

7

## A closer look at spatial dimensions:

7

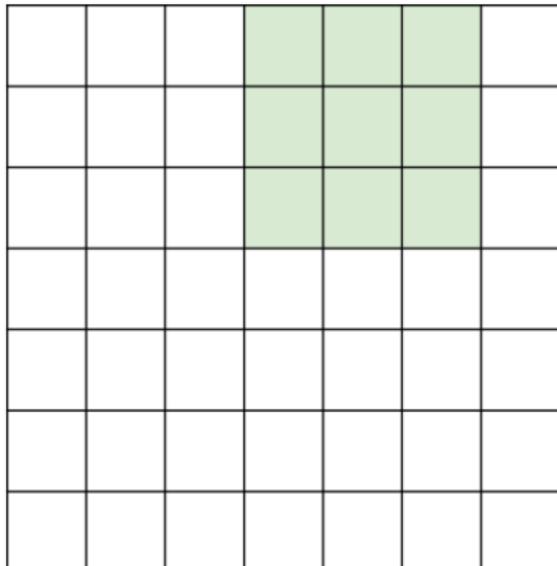


7x7 input (spatially)  
assume 3x3 filter

7

## A closer look at spatial dimensions:

7

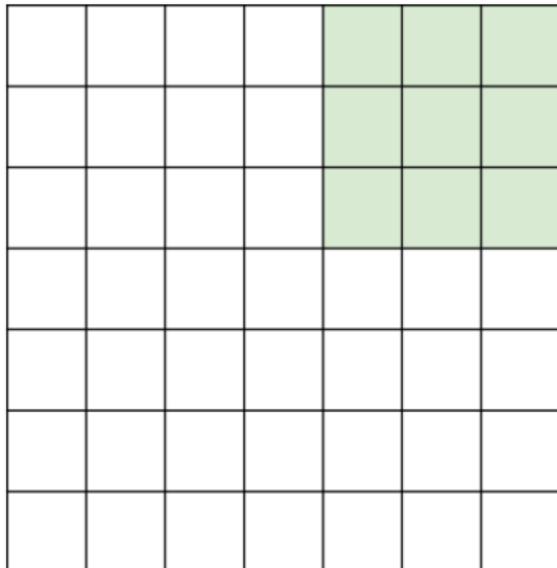


7x7 input (spatially)  
assume 3x3 filter

7

## A closer look at spatial dimensions:

7

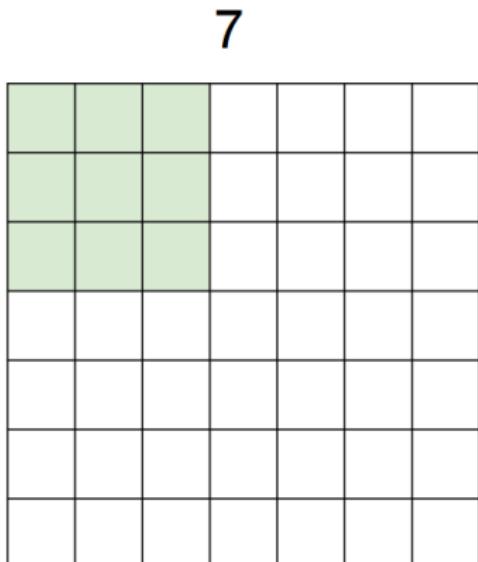


7x7 input (spatially)  
assume 3x3 filter

7

=> 5x5 output

## A closer look at spatial dimensions:

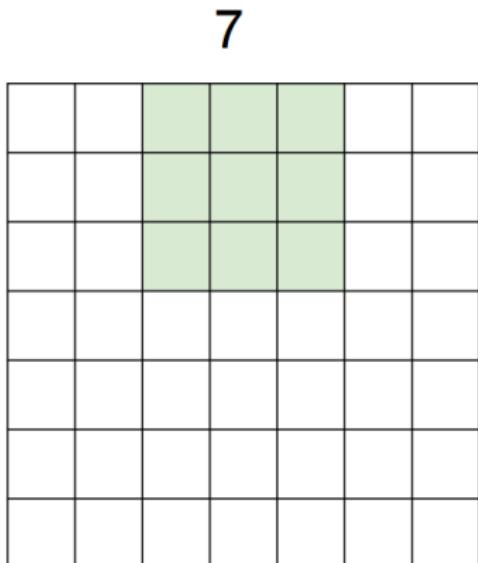


7

7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

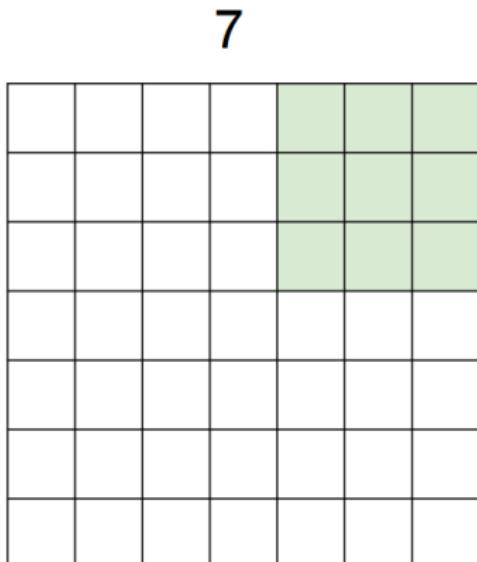
## A closer look at spatial dimensions:



7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

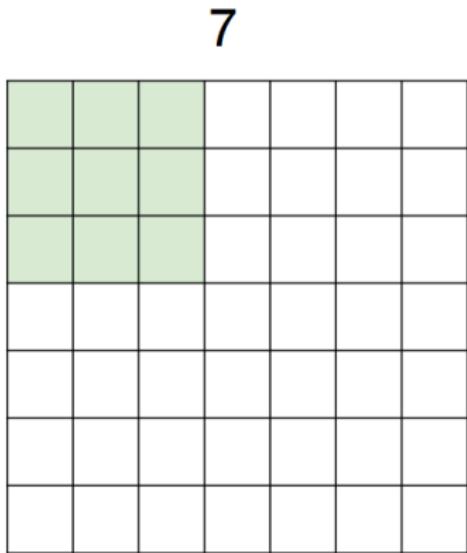
## A closer look at spatial dimensions:



7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**  
**=> 3x3 output!**

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

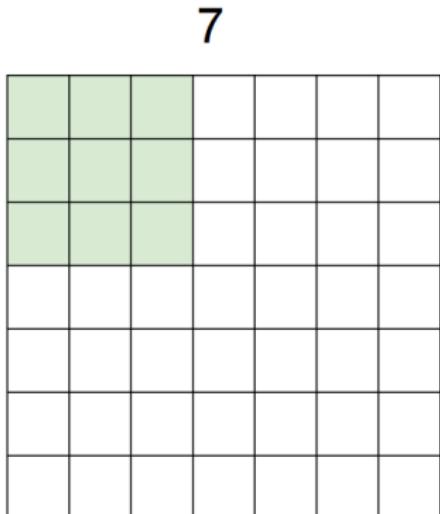
## A closer look at spatial dimensions:



7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 3?**

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

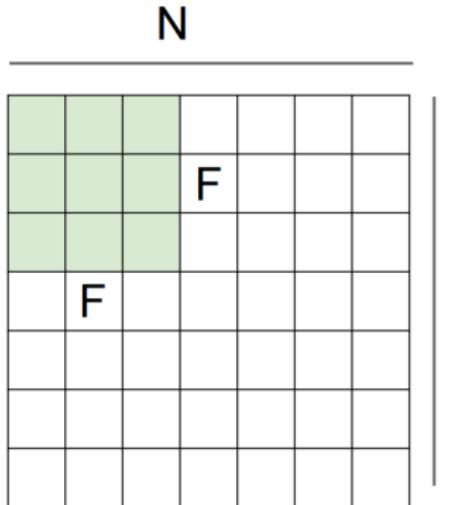
A closer look at spatial dimensions:



7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 3?**

**doesn't fit!**  
cannot apply 3x3 filter on  
7x7 input with stride 3.

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7



Output size:  
 $(N - F) / \text{stride} + 1$

e.g.  $N = 7$ ,  $F = 3$ :

$$\text{stride } 1 \Rightarrow (7 - 3)/1 + 1 = 5$$

$$\text{stride } 2 \Rightarrow (7 - 3)/2 + 1 = 3$$

$$\text{stride } 3 \Rightarrow (7 - 3)/3 + 1 = 2.33 \therefore$$

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

# In practice: Common to zero pad the border

0	0	0	0	0	0		
0							
0							
0							
0							

e.g. input 7x7

**3x3 filter, applied with stride 1**

**pad with 1 pixel border => what is the output?**

(recall:)  
$$(N - F) / \text{stride} + 1$$

# In practice: Common to zero pad the border

0	0	0	0	0	0		
0							
0							
0							
0							

e.g. input 7x7

3x3 filter, applied with **stride 1**

**pad with 1 pixel border => what is the output?**

**7x7 output!**

# In practice: Common to zero pad the border

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input 7x7

**3x3 filter, applied with stride 1**

**pad with 1 pixel border => what is the output?**

**7x7 output!**

in general, common to see CONV layers with  
stride 1, filters of size FxF, and zero-padding with  
 $(F-1)/2$ . (will preserve size spatially)

e.g.  $F = 3 \Rightarrow$  zero pad with 1

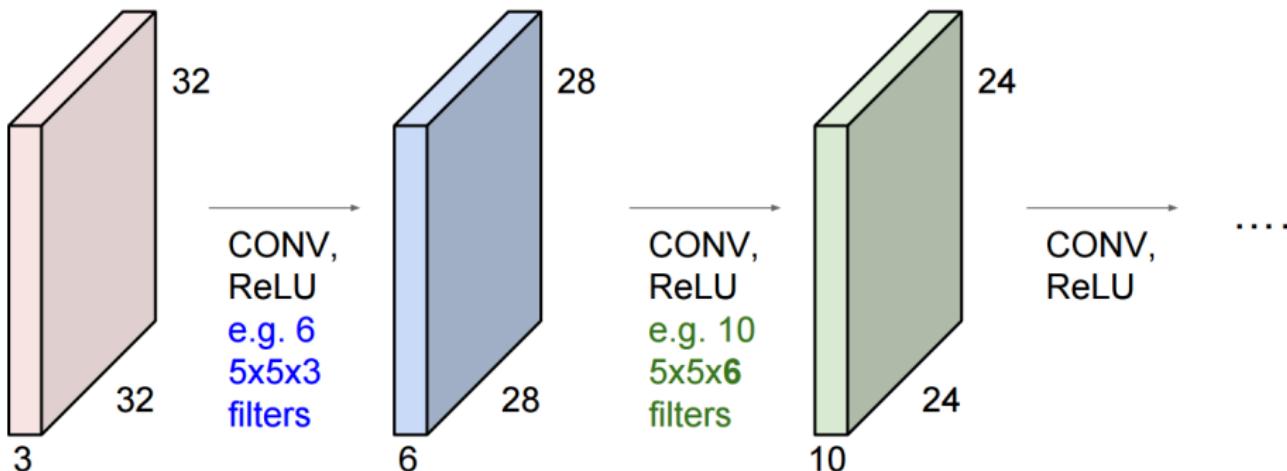
$F = 5 \Rightarrow$  zero pad with 2

$F = 7 \Rightarrow$  zero pad with 3

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

## Remember back to...

E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially!  
(32 -> 28 -> 24 ...). Shrinking too fast is not good, doesn't work well.



cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

**Summary.** To summarize, the Conv Layer:

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
  - Number of filters  $K$ ,
  - their spatial extent  $F$ ,
  - the stride  $S$ ,
  - the amount of zero padding  $P$ .
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F + 2P)/S + 1$
  - $H_2 = (H_1 - F + 2P)/S + 1$  (i.e. width and height are computed equally by symmetry)
  - $D_2 = K$
- With parameter sharing, it introduces  $F \cdot F \cdot D_1$  weights per filter, for a total of  $(F \cdot F \cdot D_1) \cdot K$  weights and  $K$  biases.
- In the output volume, the  $d$ -th depth slice (of size  $W_2 \times H_2$ ) is the result of performing a valid convolution of the  $d$ -th filter over the input volume with a stride of  $S$ , and then offset by  $d$ -th bias.

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

**Summary.** To summarize, the Conv Layer:

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
  - Number of filters  $K$ ,
  - their spatial extent  $F$ ,
  - the stride  $S$ ,
  - the amount of zero padding  $P$ .
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F + 2P)/S + 1$
  - $H_2 = (H_1 - F + 2P)/S + 1$  (i.e. width and height are computed equally by symmetry)
  - $D_2 = K$
- With parameter sharing, it introduces  $F \cdot F \cdot D_1$  weights per filter, for a total of  $(F \cdot F \cdot D_1) \cdot K$  weights and  $K$  biases.
- In the output volume, the  $d$ -th depth slice (of size  $W_2 \times H_2$ ) is the result of performing a valid convolution of the  $d$ -th filter over the input volume with a stride of  $S$ , and then offset by  $d$ -th bias.

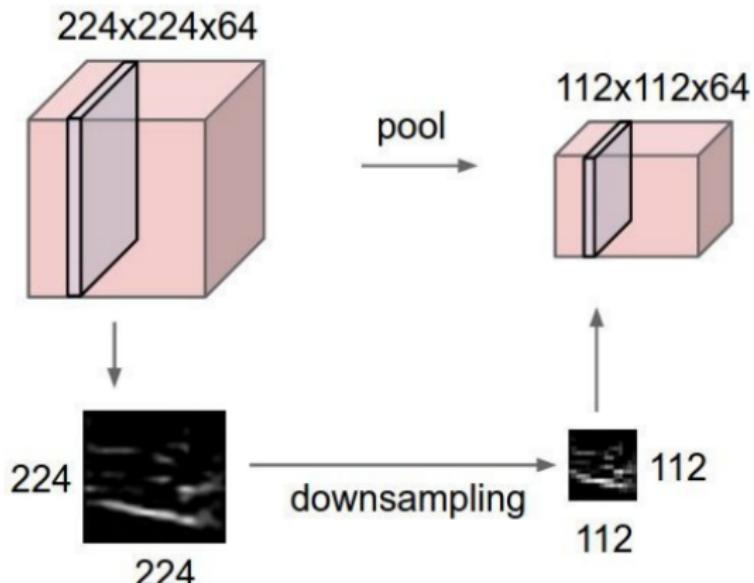
### Common settings:

- $K = (\text{powers of 2, e.g. } 32, 64, 128, 512)$
- $F = 3, S = 1, P = 1$
  - $F = 5, S = 1, P = 2$
  - $F = 5, S = 2, P = ?$  (whatever fits)
  - $F = 1, S = 1, P = 0$

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

# Pooling layer

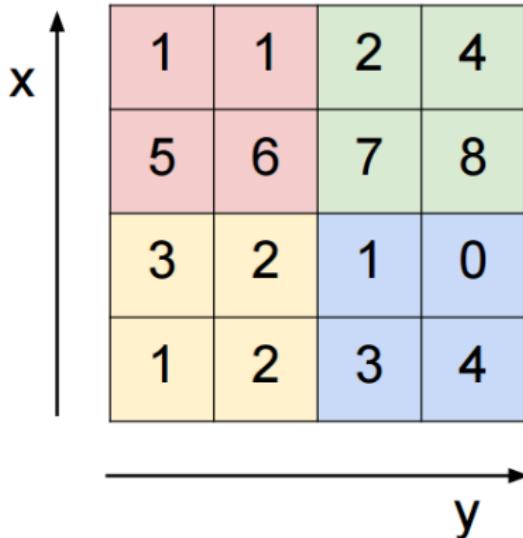
- makes the representations smaller and more manageable
- operates over each activation map independently:



cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

# MAX POOLING

Single depth slice



max pool with 2x2 filters  
and stride 2

A 2x2 grid representing the output of max pooling. The cells are colored in a repeating pattern: top-left (red, red), top-right (green, green), bottom-left (yellow, yellow), and bottom-right (blue, blue). The output values are:

6	8
3	4

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires two hyperparameters:
  - their spatial extent  $F$ ,
  - the stride  $S$ ,
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F)/S + 1$
  - $H_2 = (H_1 - F)/S + 1$
  - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

cs231n - Fei-Fei Li, Andrej Karpathy, Justin Johnson, Lecture 7