

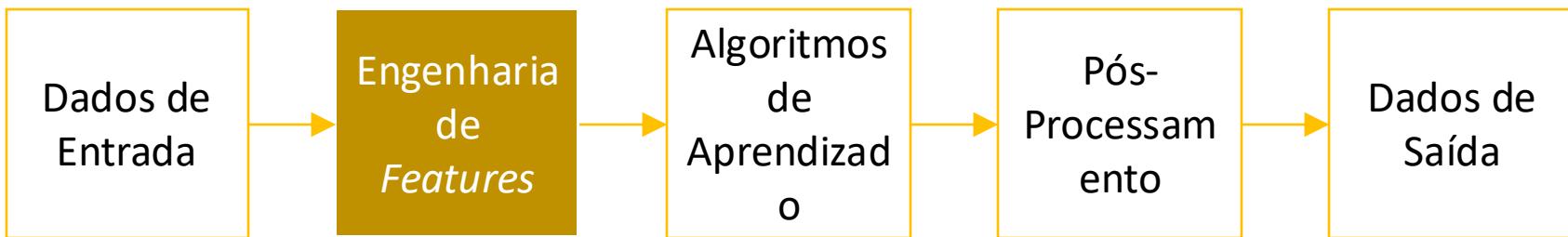
# Aprendizado de Máquina

Paradigma baseado em Otimização  
Redes Neurais IV

Prof. Me. Otávio Parraga

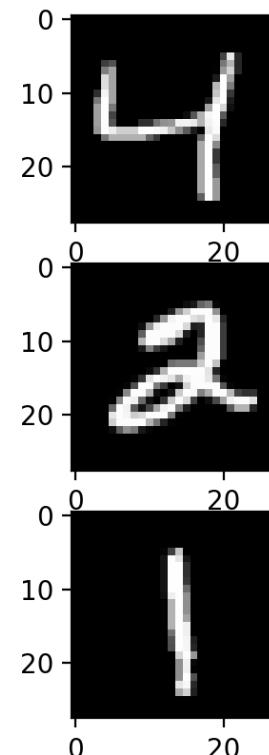
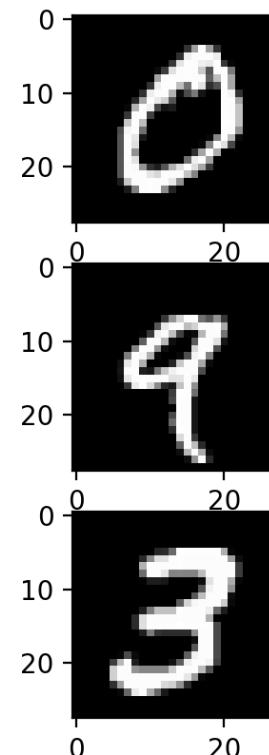
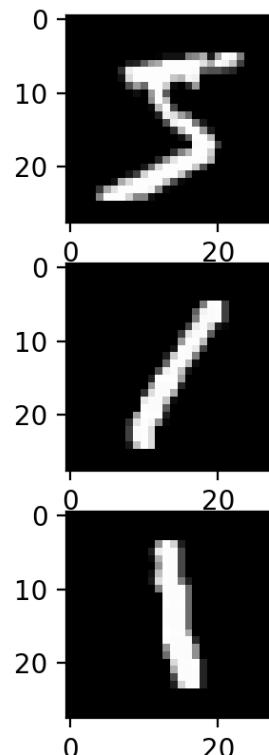
# Engenharia de Features

- O processo de Aprendizado de Máquina
- Engenharia de Features como um dos processos **mais** custosos
  - E **menos** recompensadores!!!



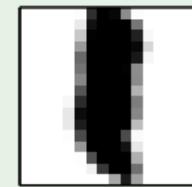
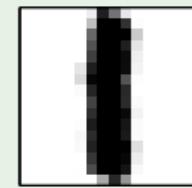
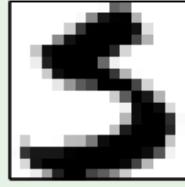
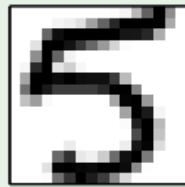
# MNIST

- Representar dados é uma tarefa difícil
- Como representar imagens de números de forma a diferenciar eles?



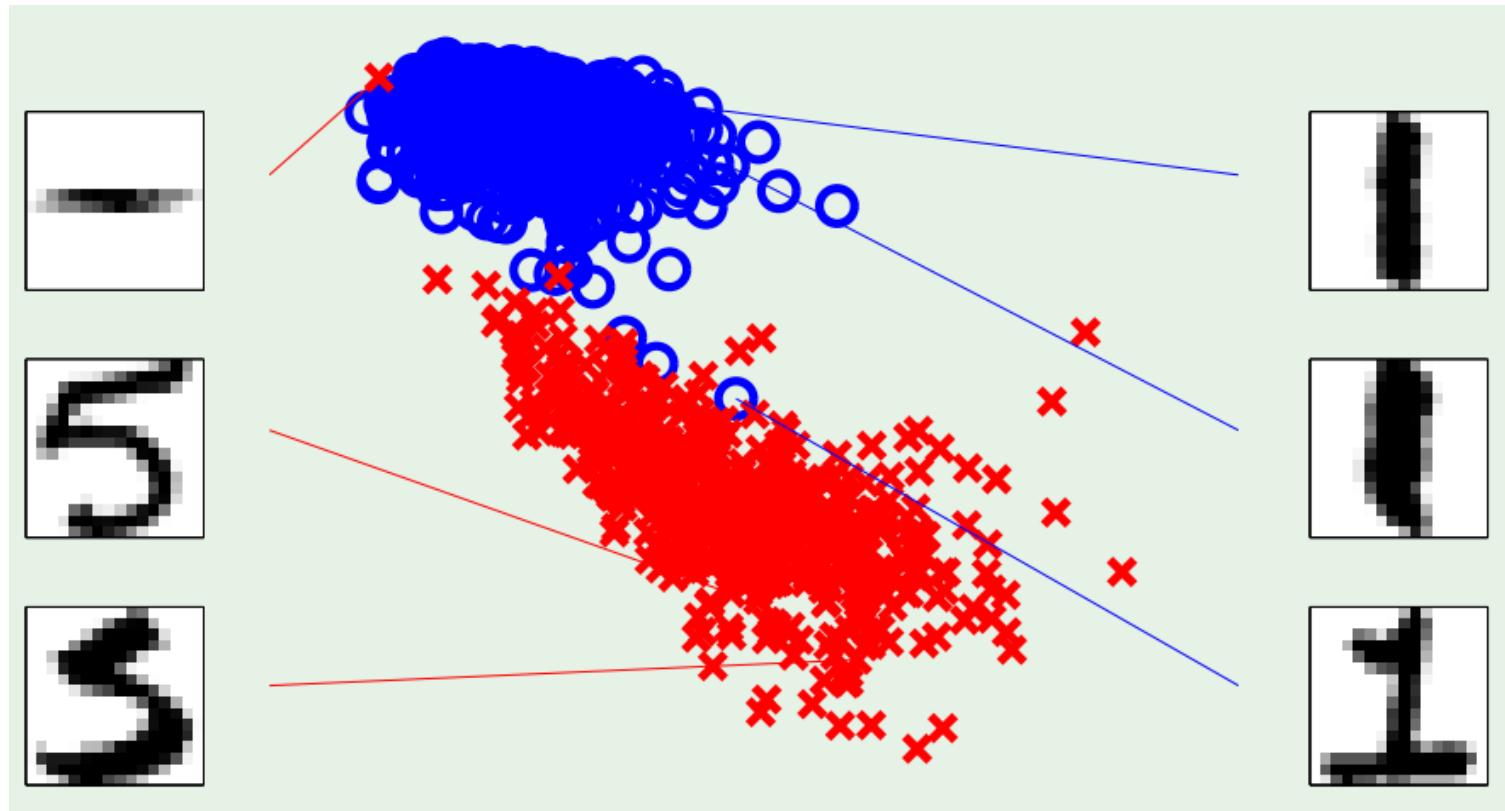
# MNIST

- Poderíamos diferenciar por duas características:
  - $x_1$ : Intensidade média dos pixels
  - $x_2$ : Simetria vertical



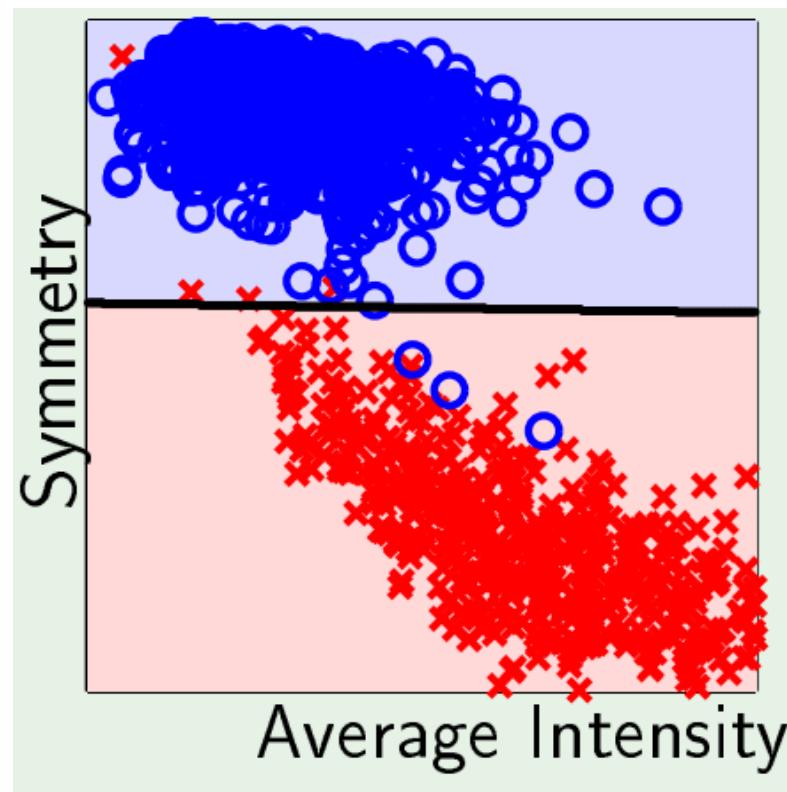
# MNIST

- Poderíamos diferenciar por duas características:
  - $x_1$ : Intensidade média dos pixels
  - $x_2$ : Simetria vertical



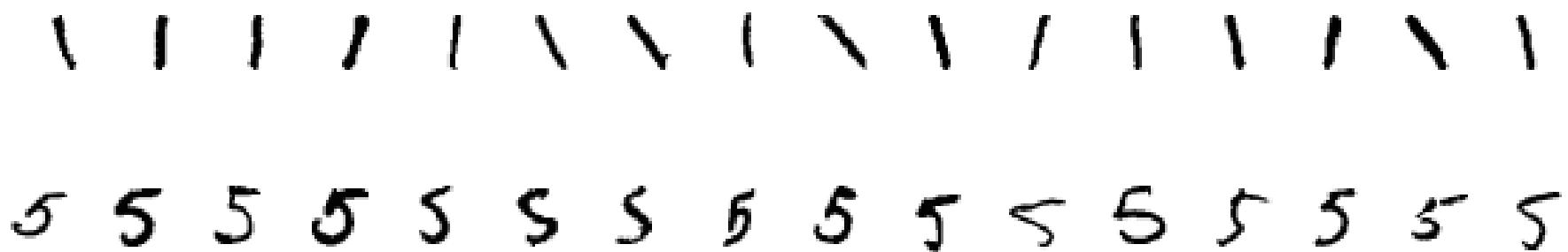
# MNIST

- Conseguimos treinar um separador linear!



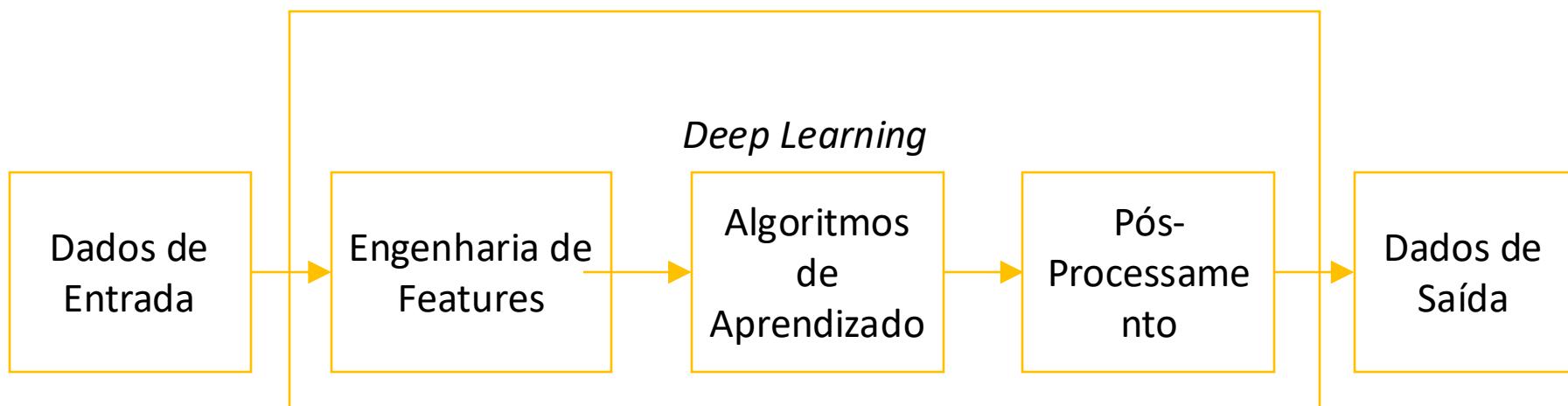
# O Problema

- Como poderíamos lidar com a variabilidade?
  - Engenharia de Features se torna uma tarefa **custosa e pouco recompensadora**
    - Especialmente com dados **não estruturados**

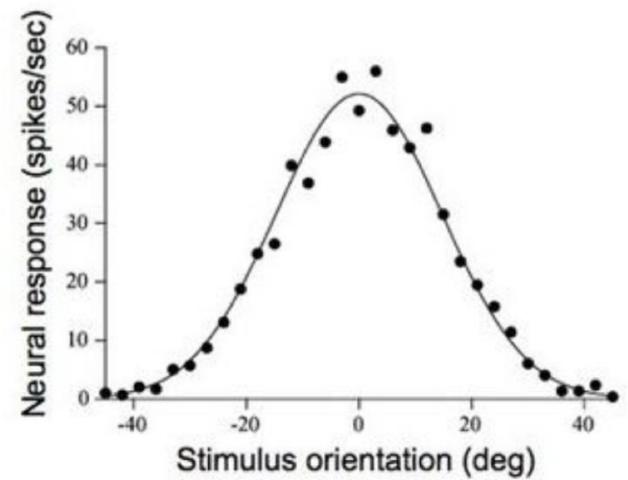
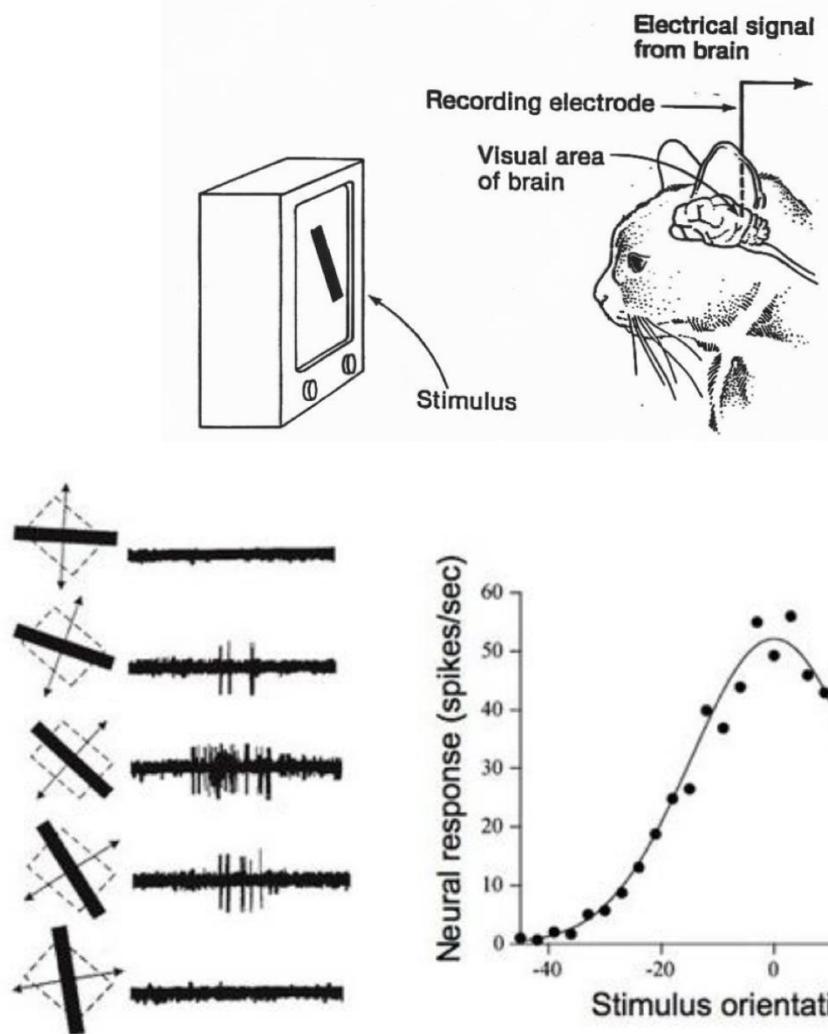
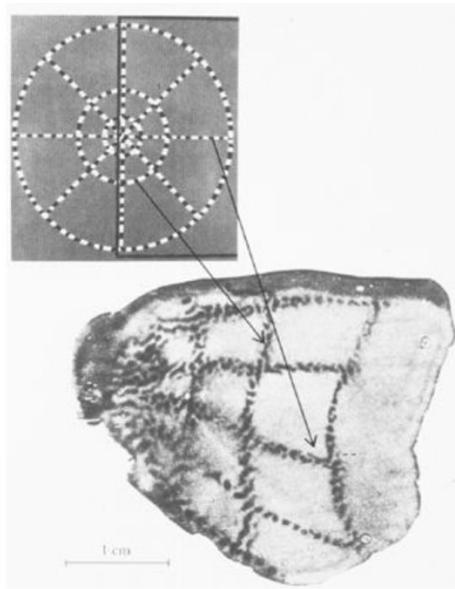


# Deep Learning

- E se delegássemos essa tarefa para o próprio algoritmo?



# Um pouco de história...

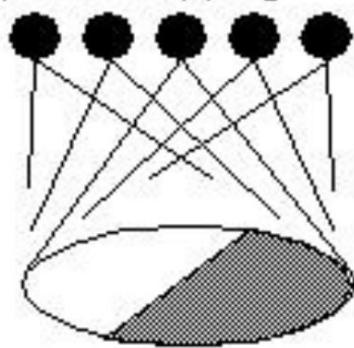


Hubel e Wiesel (1959)

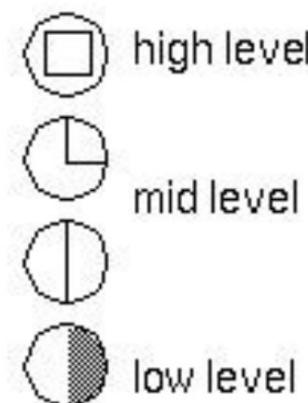
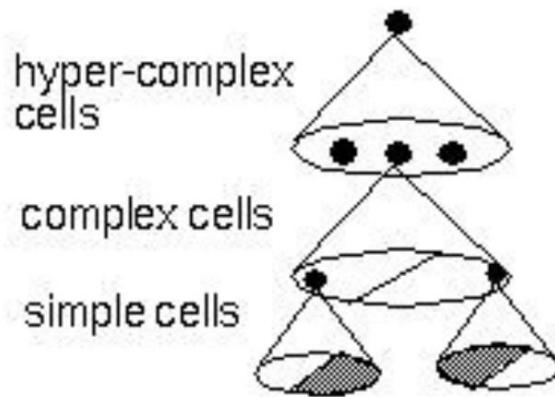
# Um pouco de história...

## Hubel & Weisel

topographical mapping



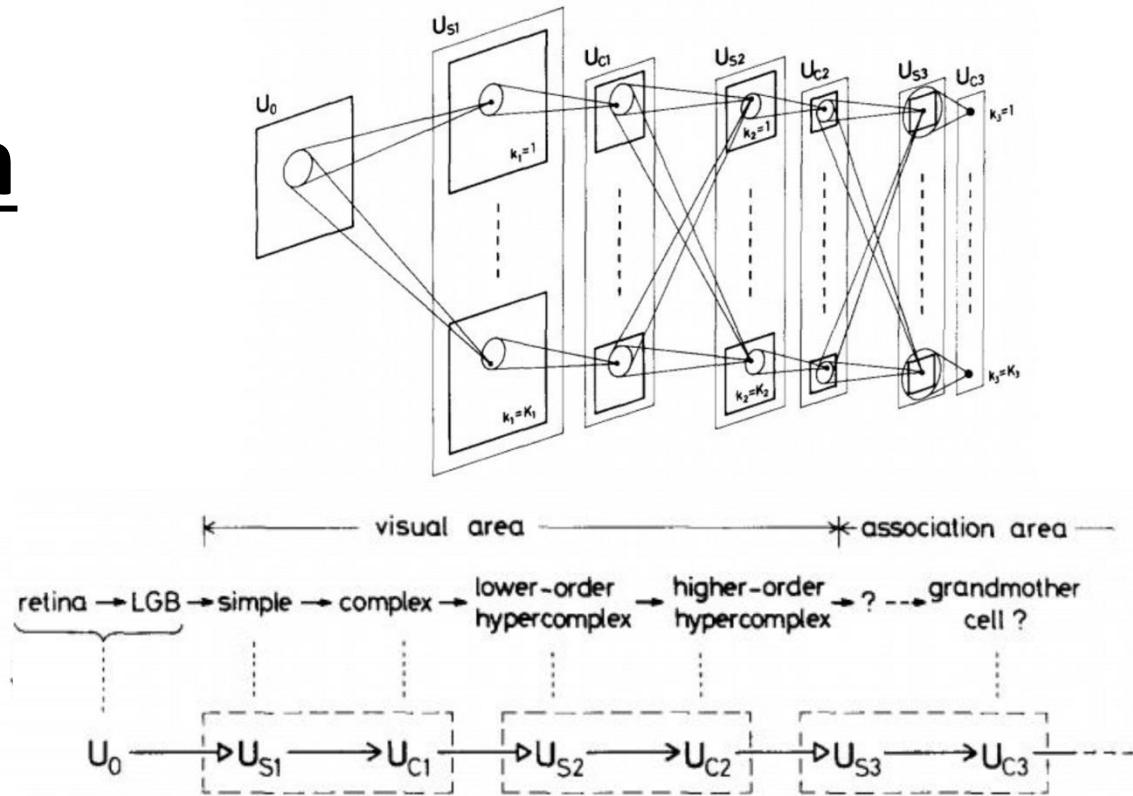
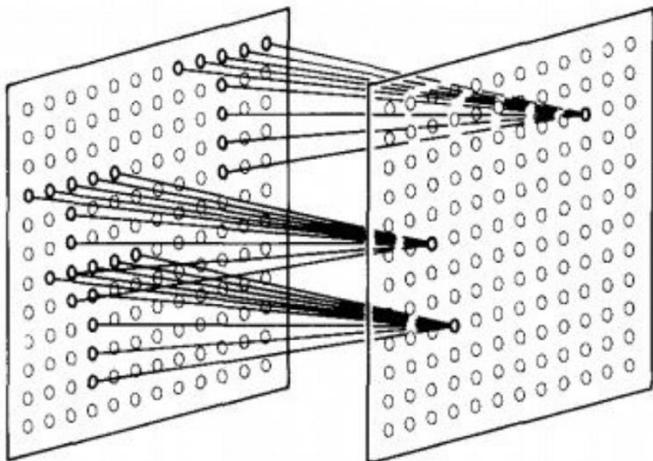
## featural hierarchy



# Um pouco de história...

(Fukushima, 1980)

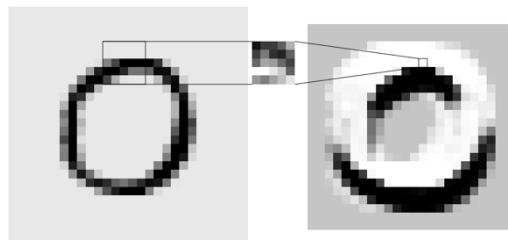
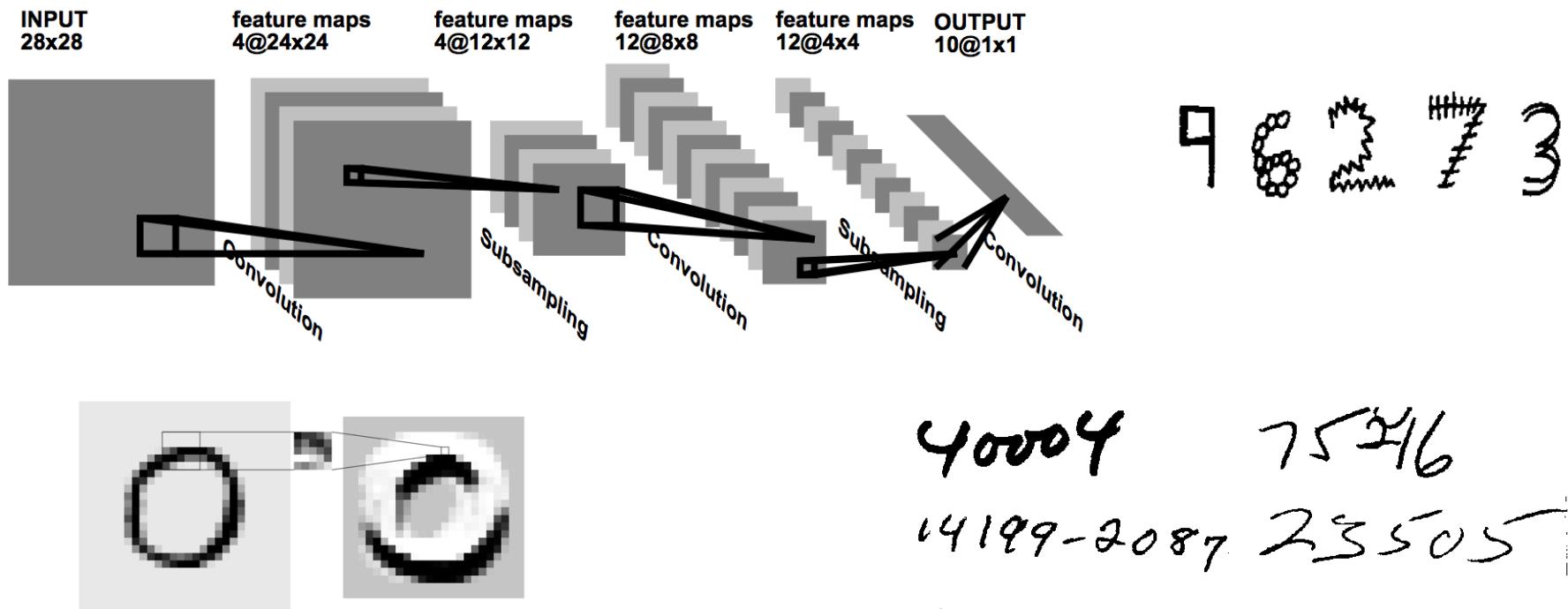
## Neurocognitron



# Um pouco de história...

1990 – A primeira convolucional! (LeNet)

98.3% de acurácia no MNIST



Y. LeCun, B. Boser, J. S. Denker, Richard E. Howard, W. Hubbard, L. D. Jackel, and D. Henderson. Handwritten digit recognition with a back-propagation network. In Advances in neural information processing systems, 1990.

40004 75216  
14199-2087 23505  
96203 14310  
44151 05753

# Um pouco de história...

1995 – LeNet 5

99.1% de acurácia no MNIST

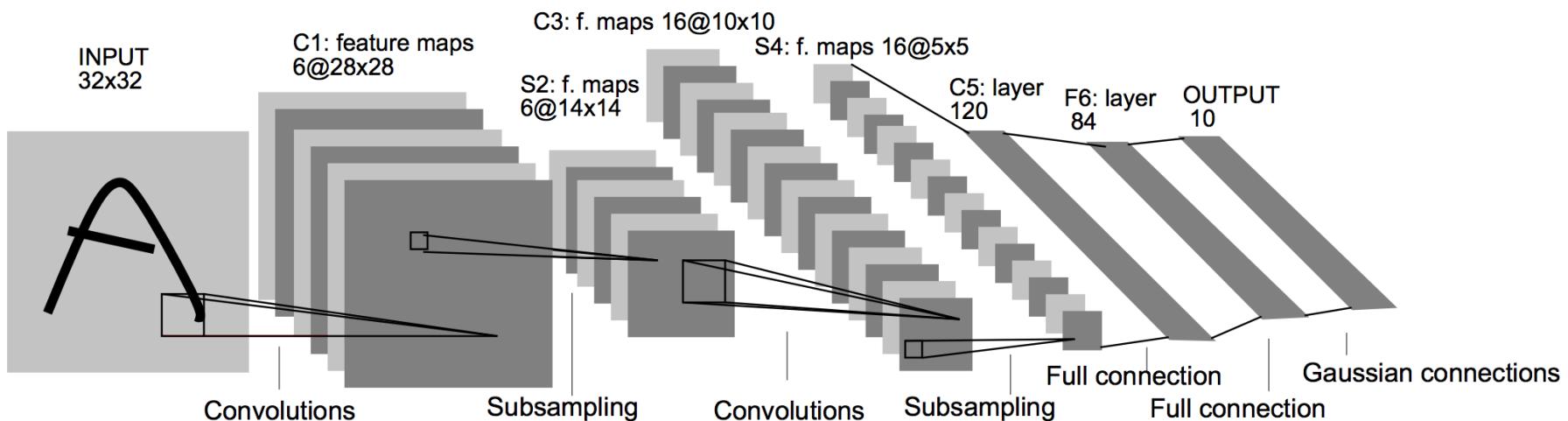
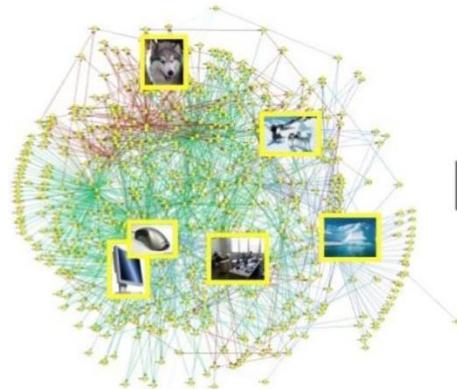


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Y. LeCun, L. D. Jackel, L. Bottou, A. Brunot, C. Cortes, J. S. Denker, H. Drucker, I. Guyon, U. A. Muller, E. Sackinger, P. Simard, and V. Vapnik. Comparison of learning algorithms for handwritten digit recognition. In F. Fogelman and P. Gallinari, editors, *International Conference on Artificial Neural Networks*, pages 53-60, Paris, 1995.

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, november 1998.

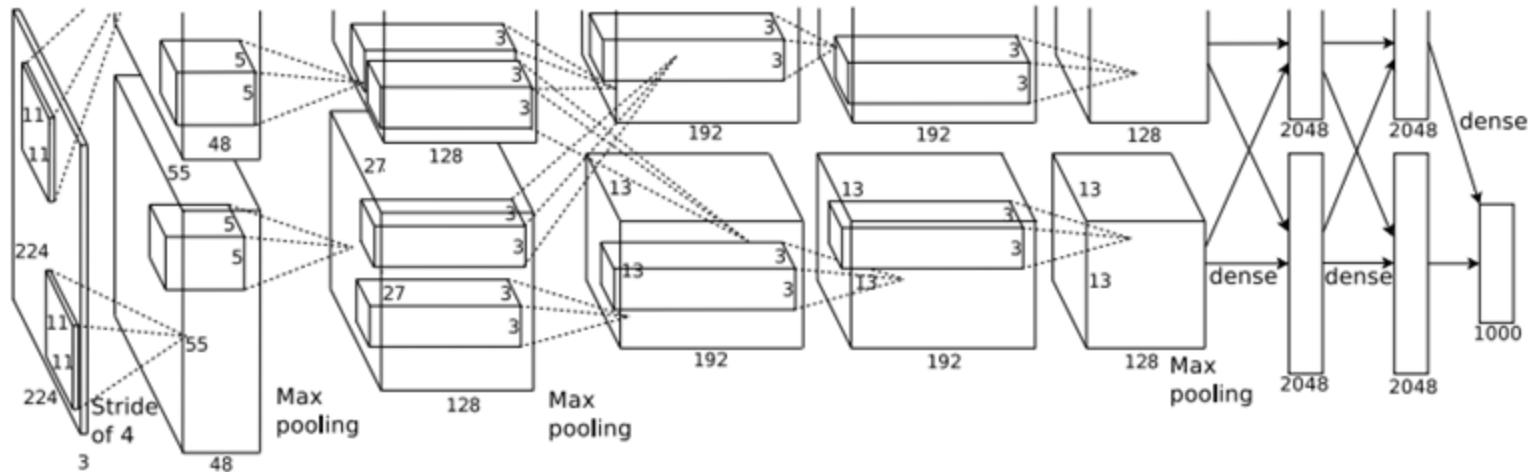
# Um pouco de história...



IMAGENET

2012 – AlexNet

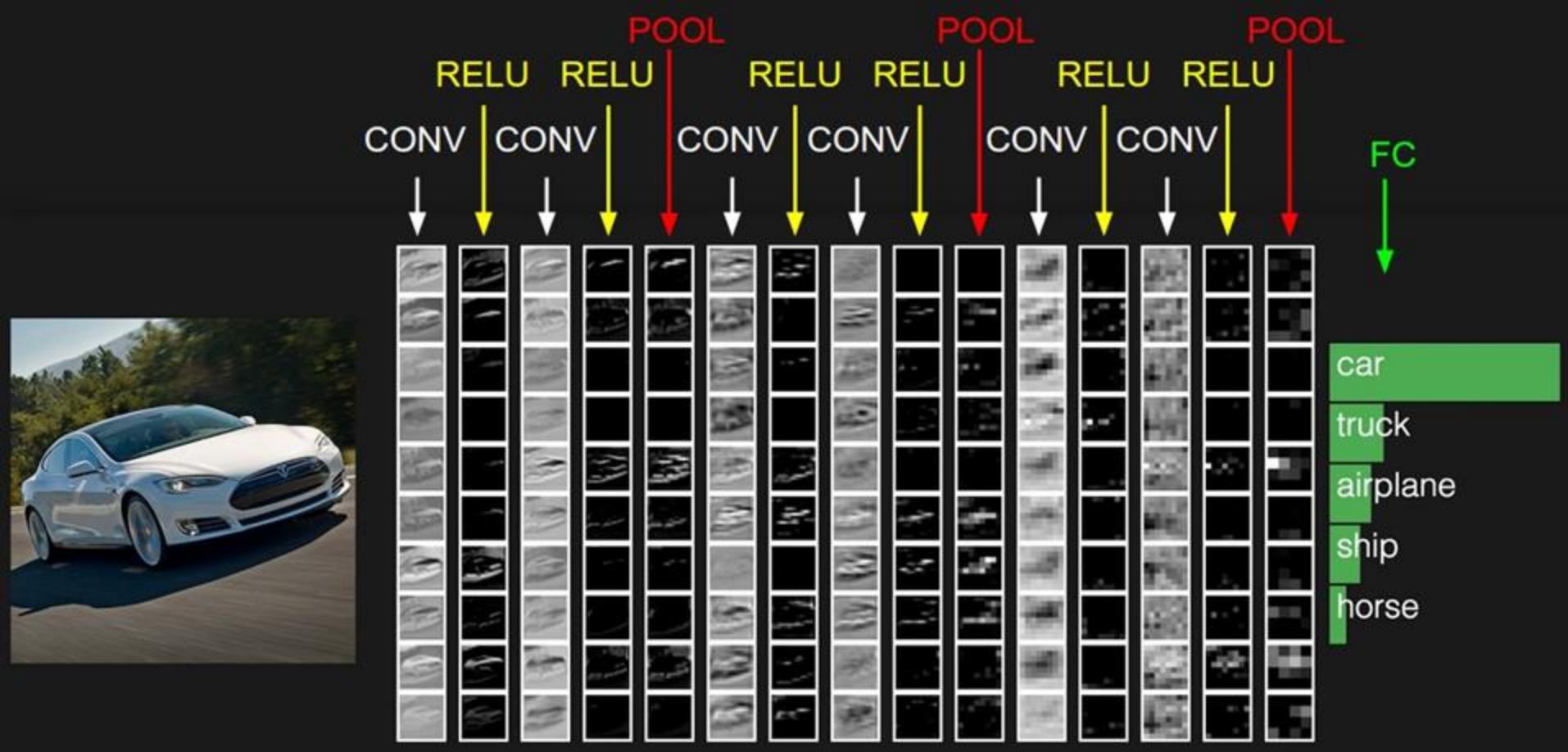
84.7% de acurácia  
no ImageNet



Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In Advances in neural information processing systems , pages 1097–1105, 2012.

# O que É uma Rede Neural Convolucional:

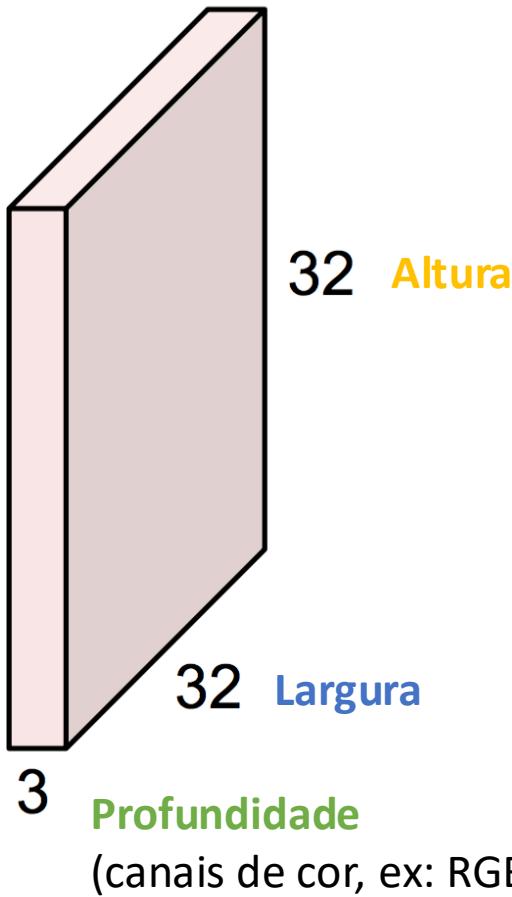
- "Rede Neural *feed-forward* cujos neurônios estão dispostos de tal forma a responder por regiões sobrepostas (**campos receptores**) que preenchem o campo visual"
- "Múltiplas camadas de **pequenas coleções de neurônios** que processam porções da imagem de entrada (campos receptores)"
- "Sequências de **camadas de convolução** intercaladas por funções de ativação"



Começaremos com a camada de convolução...

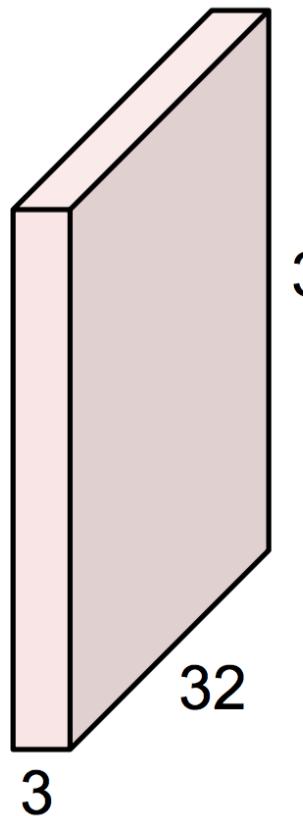
# Camada de Convolução

Imagen  $32 \times 32 \times 3$



# Camada de Convolução

Imagen  $32 \times 32 \times 3$



Filtro  $5 \times 5 \times 3$

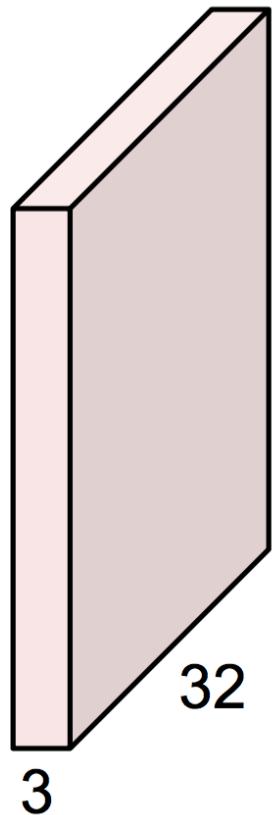


**Convoluir** o filtro com a imagem

Deslizar espacialmente o filtro  
pela imagem, computando  
produtos internos

# Camada de Convolução

Imagen  $32 \times 32 \times 3$



Filtros sempre compreendem a profundidade inteira do volume de entrada

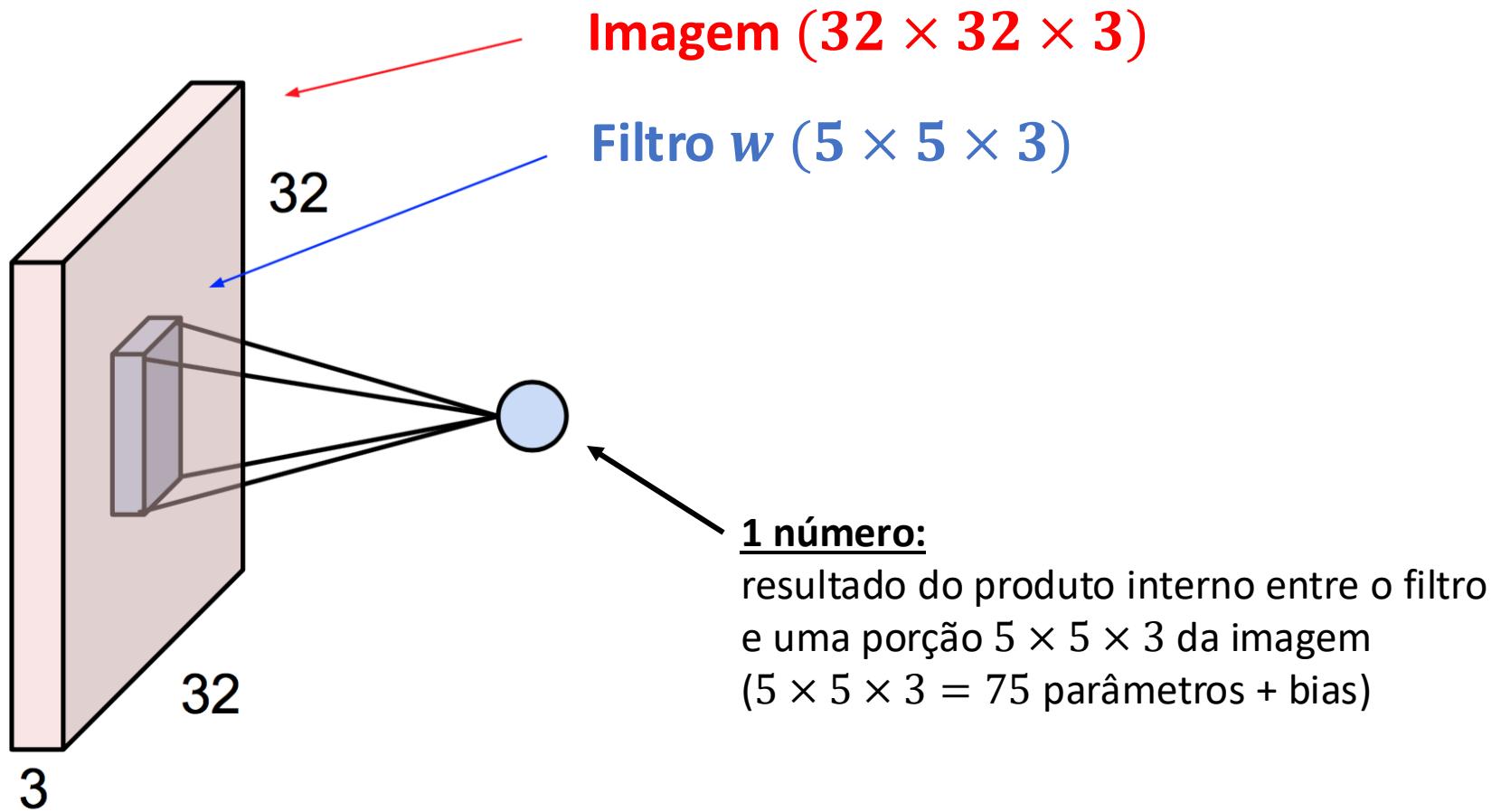
Filtro  $5 \times 5 \times 3$



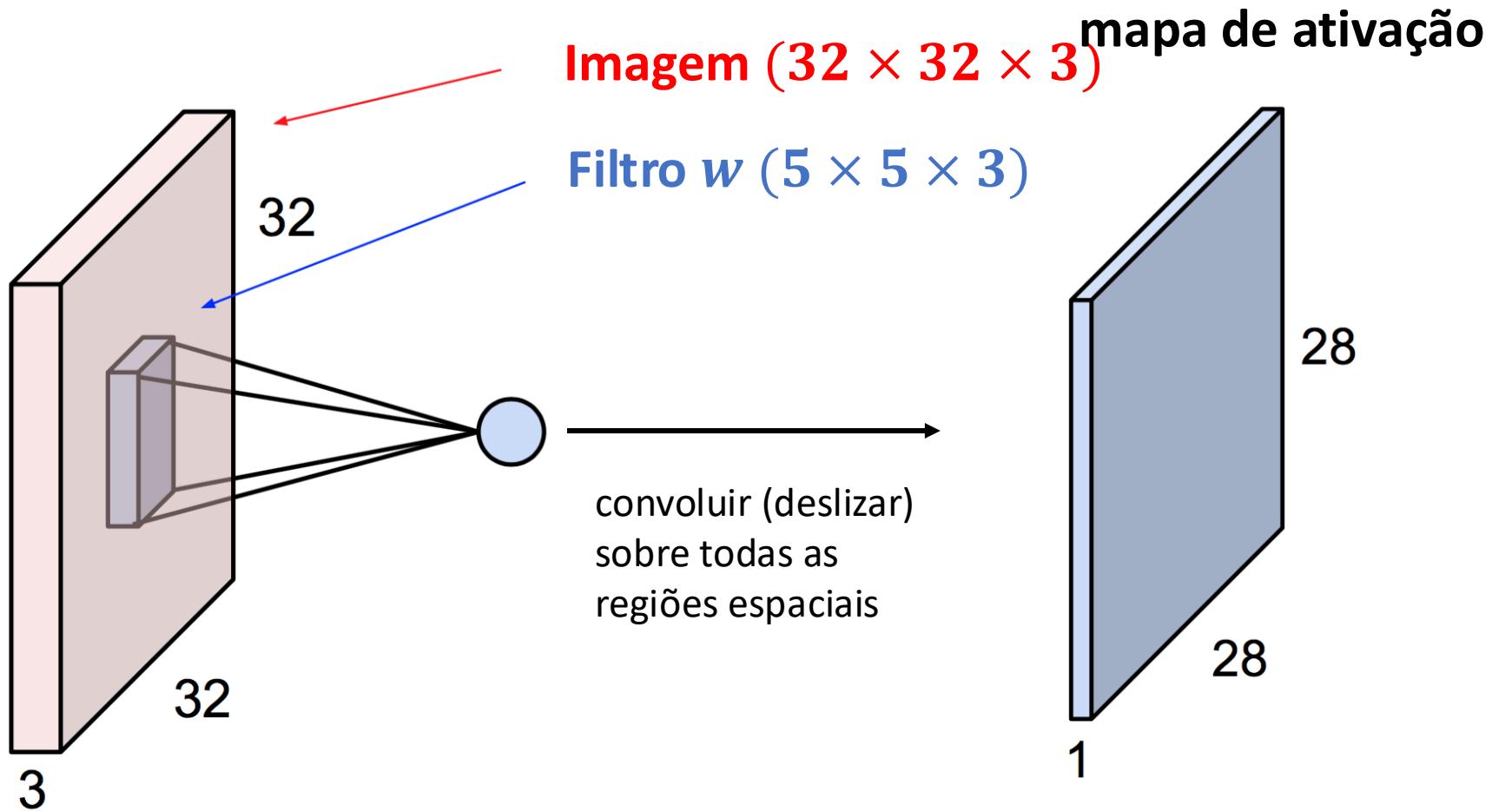
**Convoluir** o filtro com a imagem

Deslizar espacialmente o filtro pela imagem, computando produtos internos

# Camada de Convolução

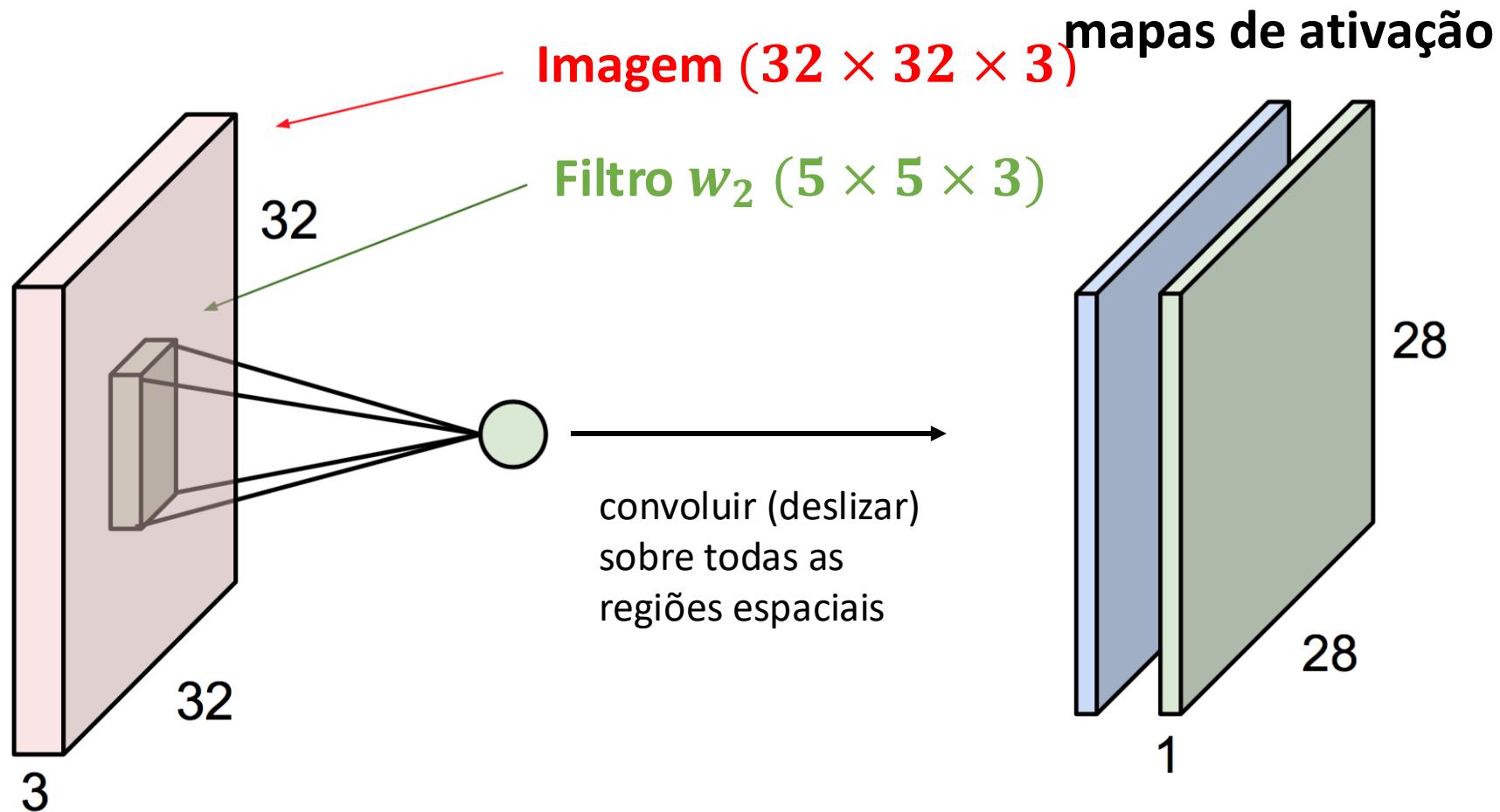


# Camada de Convolução



# Camada de Convolução

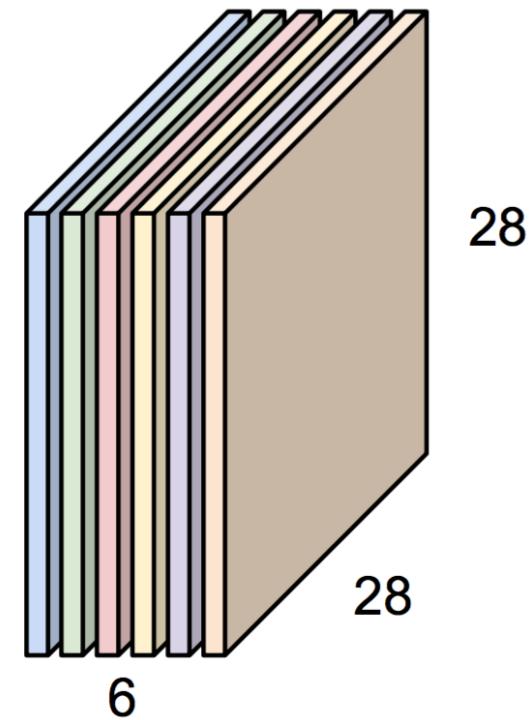
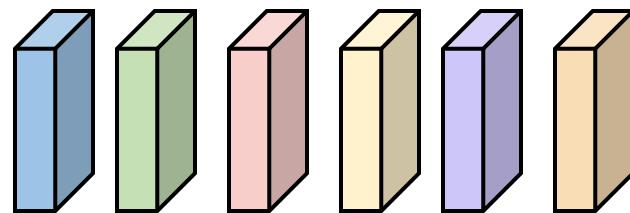
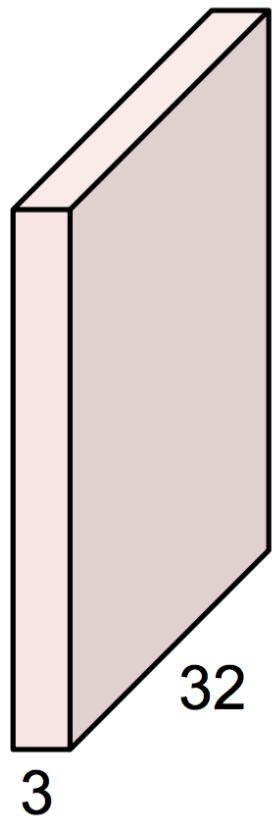
Considere a existência de um segundo filtro (verde)



# Camada de Convolução

Se tivermos 6 filtros  $5 \times 5 \times 3$ ,  
teremos 6 feature maps distintos:

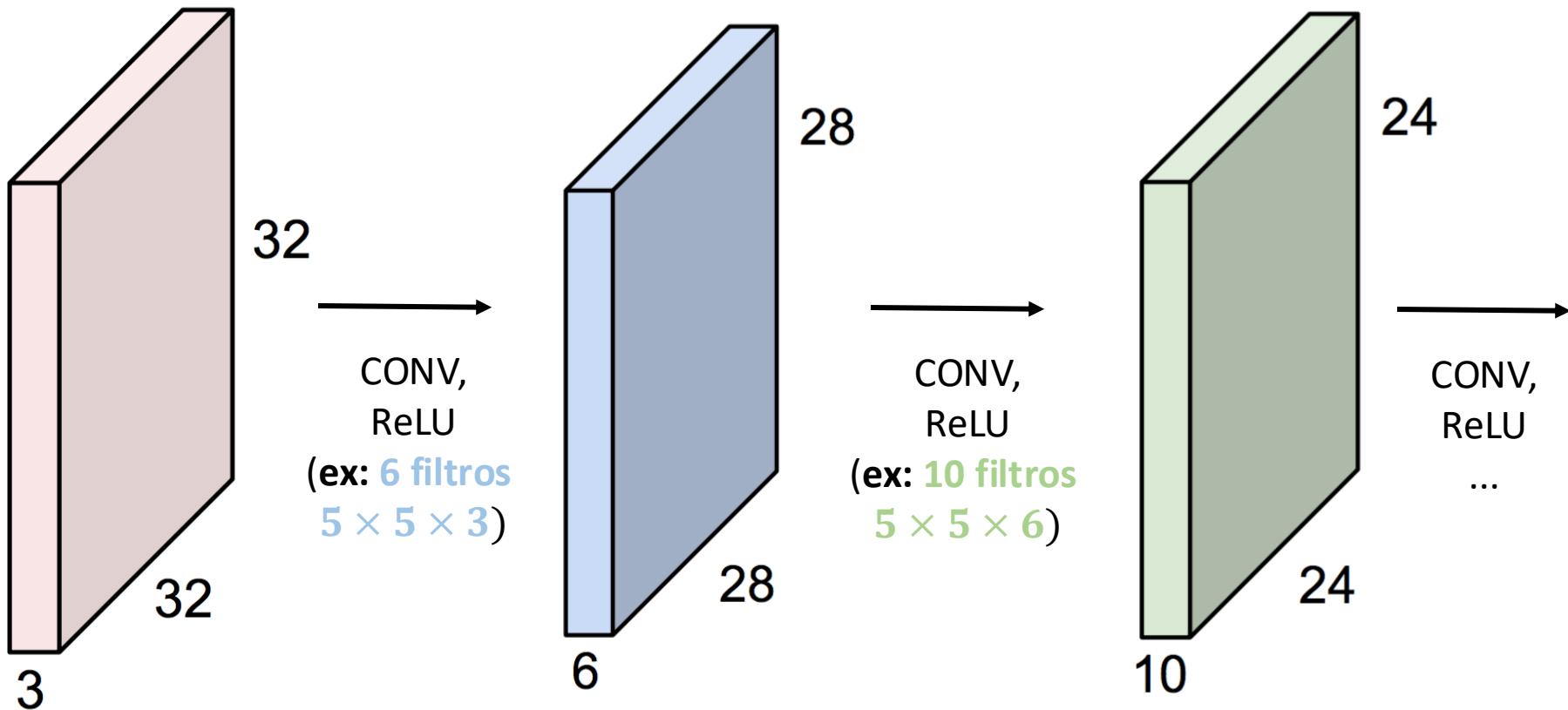
**mapas de ativação**



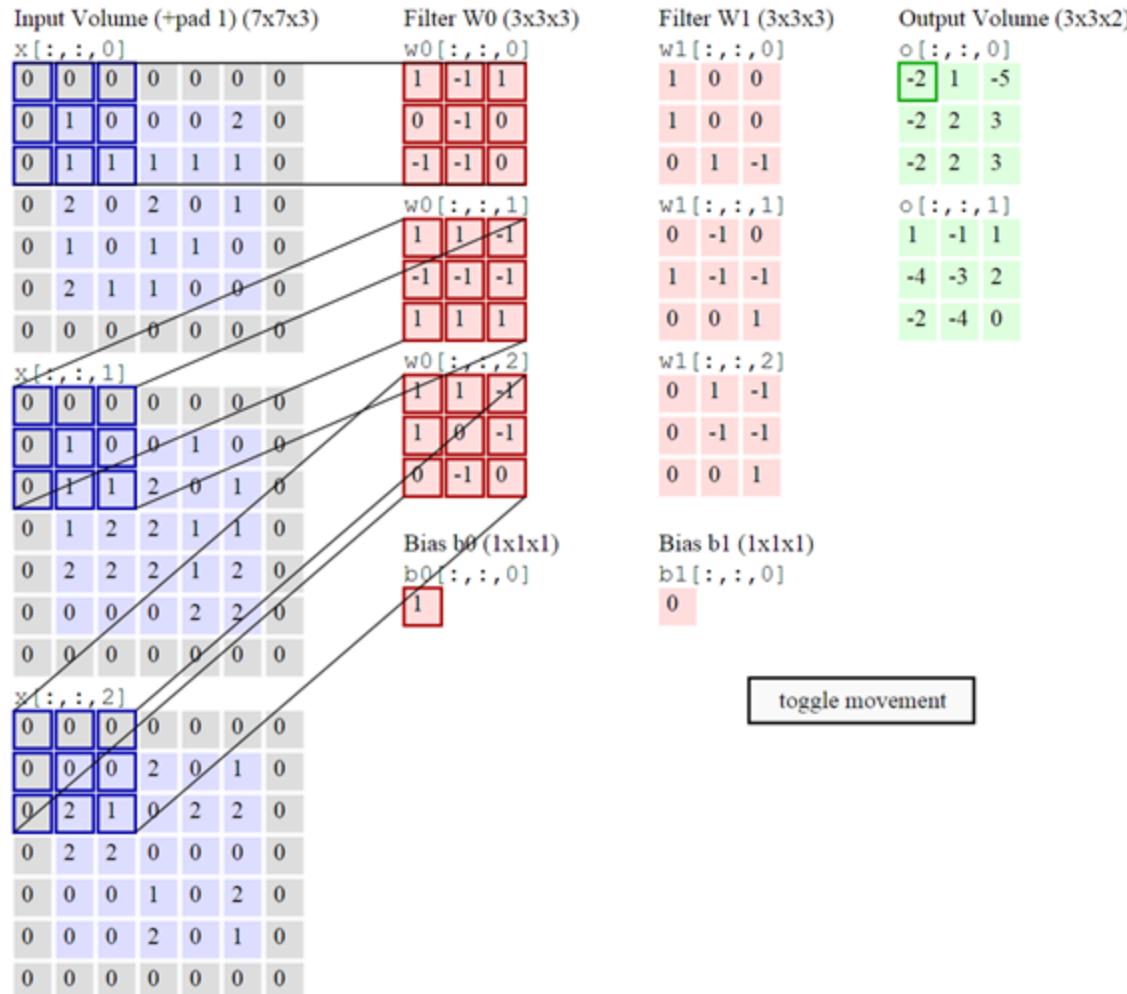
Empilhamos o *feature maps* para gerar uma "nova imagem" de tamanho  $28 \times 28 \times 6$

# Rede Convolucional

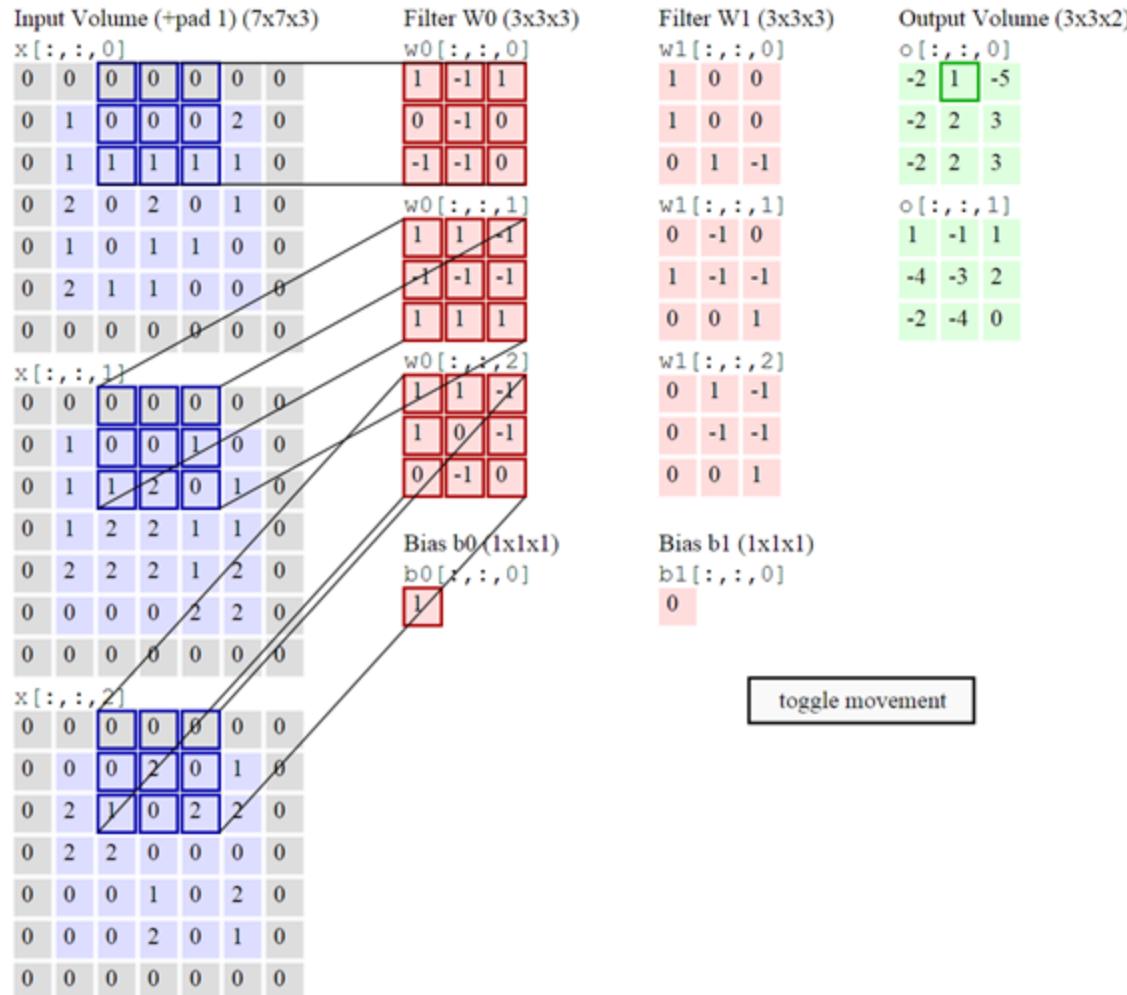
“Sequência de camadas de convolução intercaladas por funções de ativação”



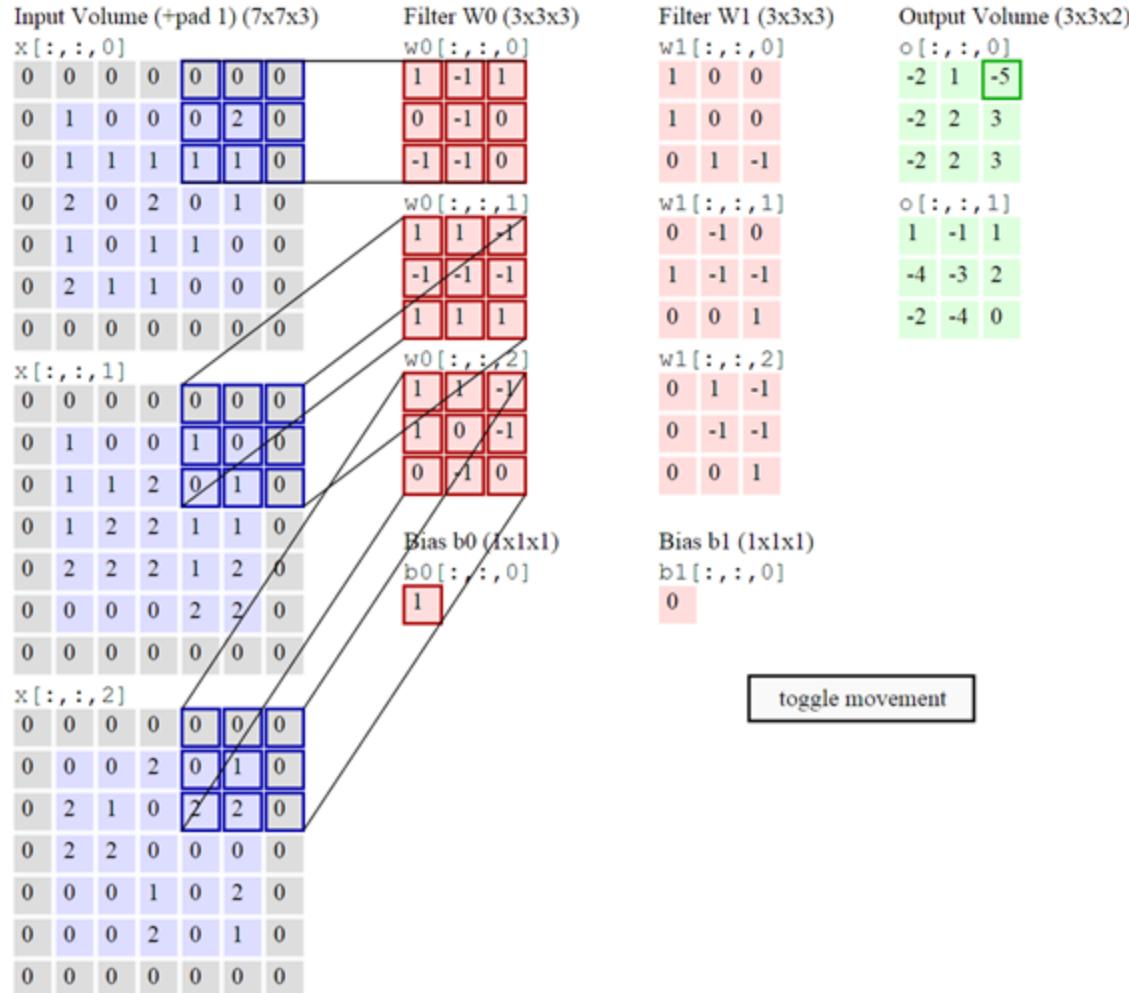
# Redes Neurais Convolucionais



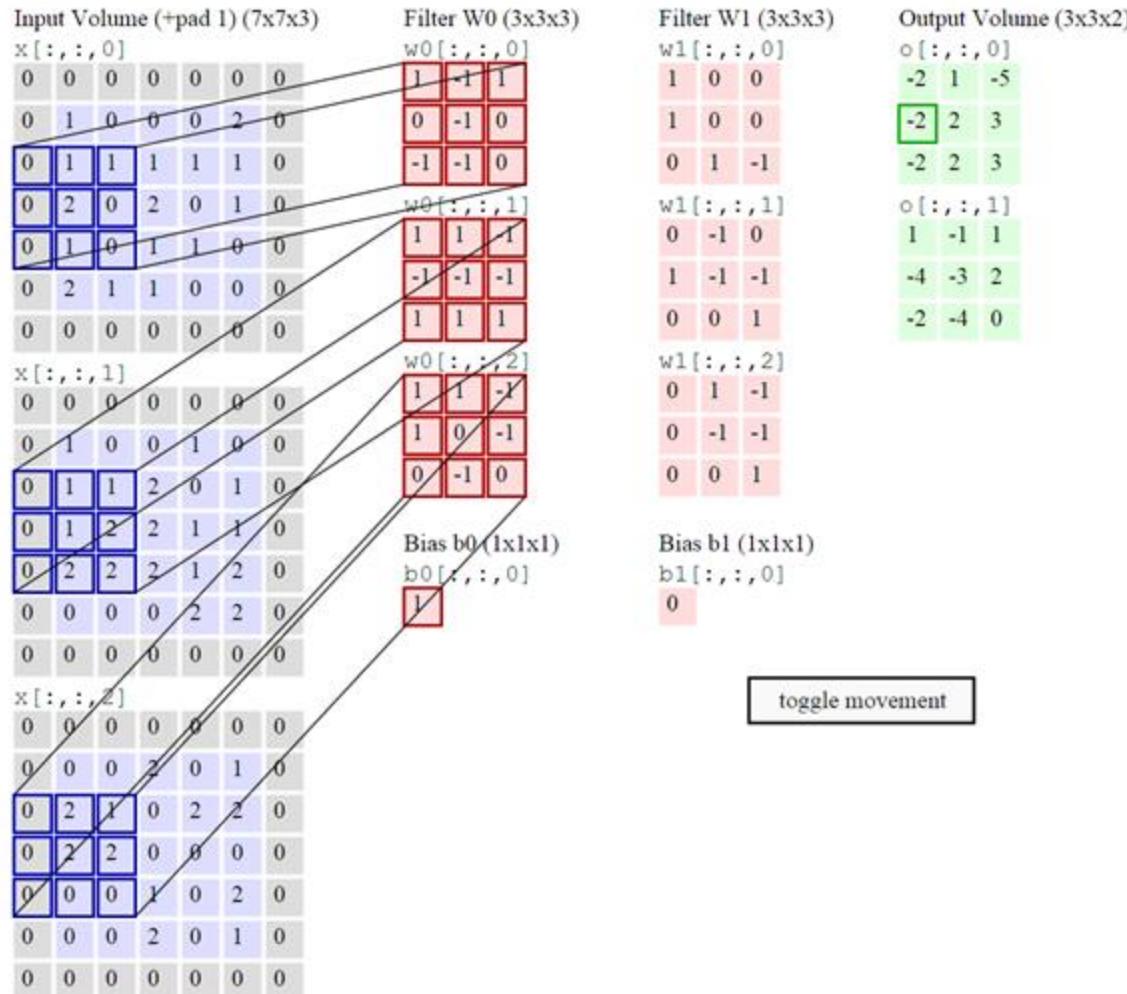
# Redes Neurais Convolucionais



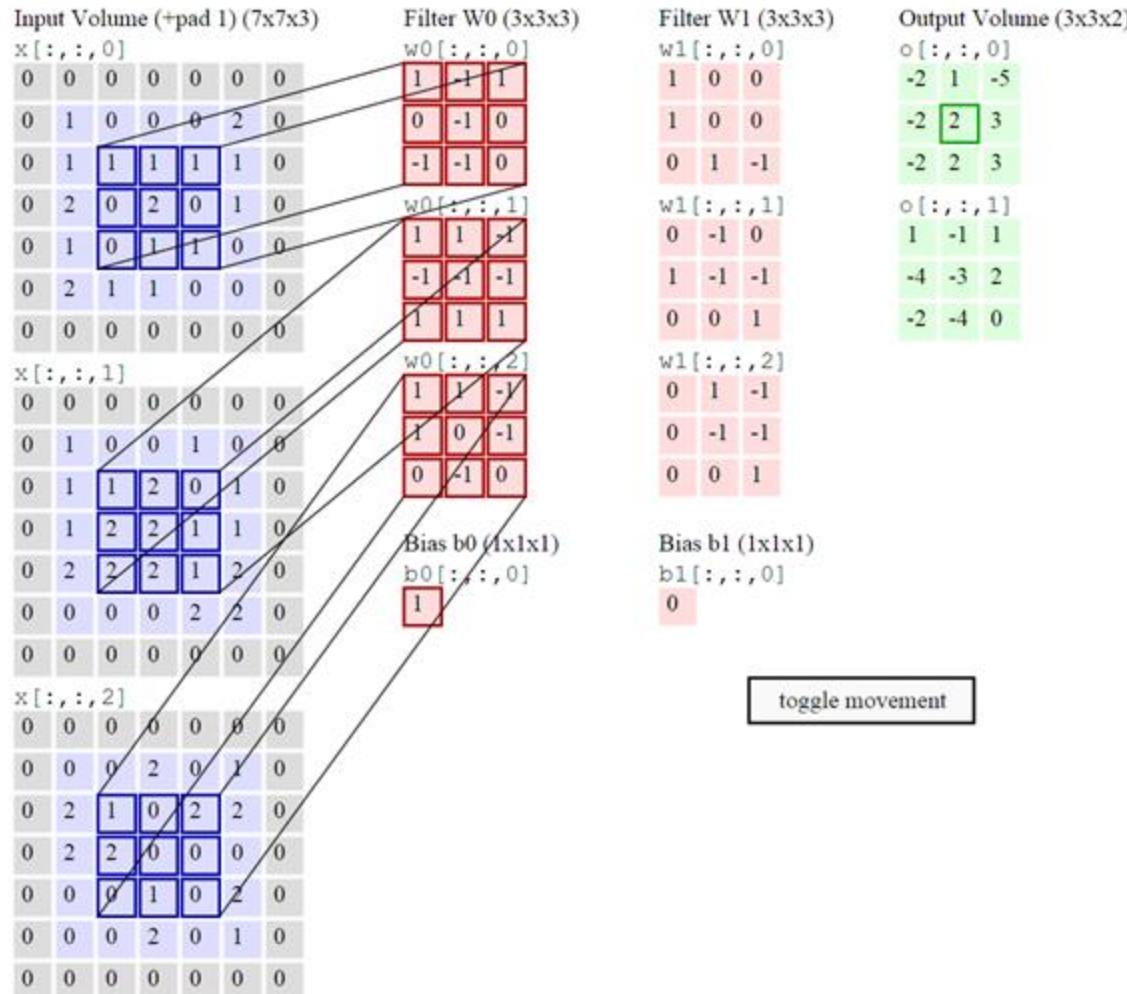
# Redes Neurais Convolucionais



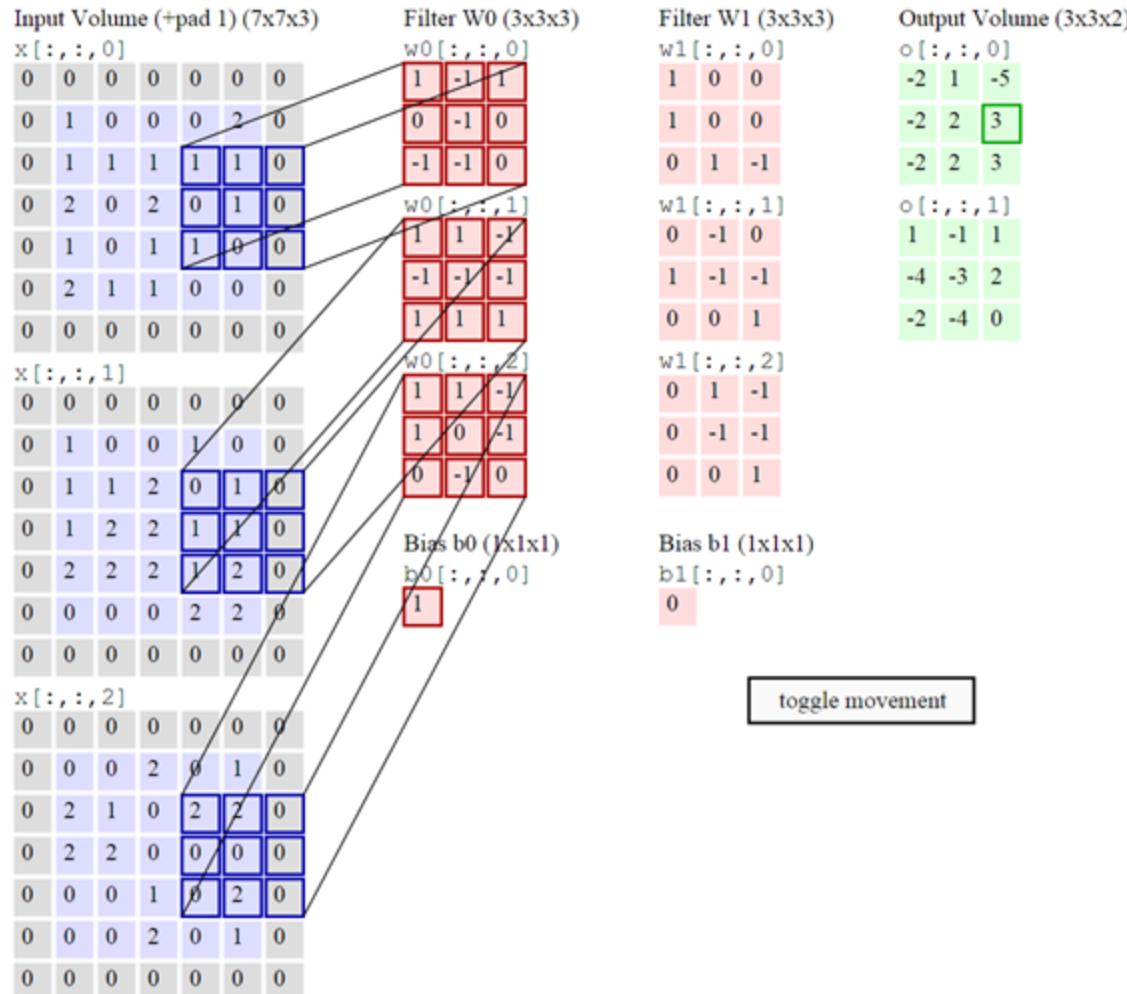
# Redes Neurais Convolucionais



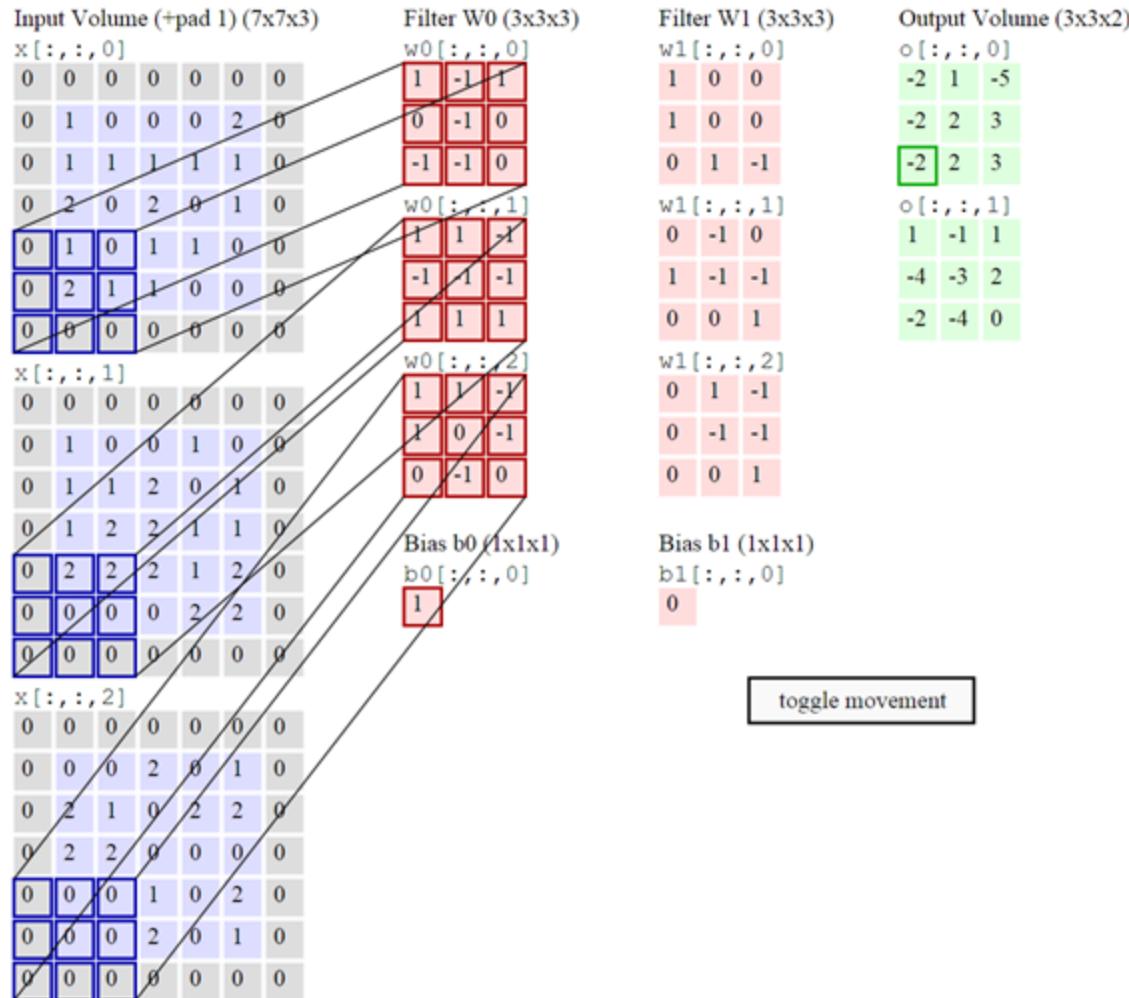
# Redes Neurais Convolucionais



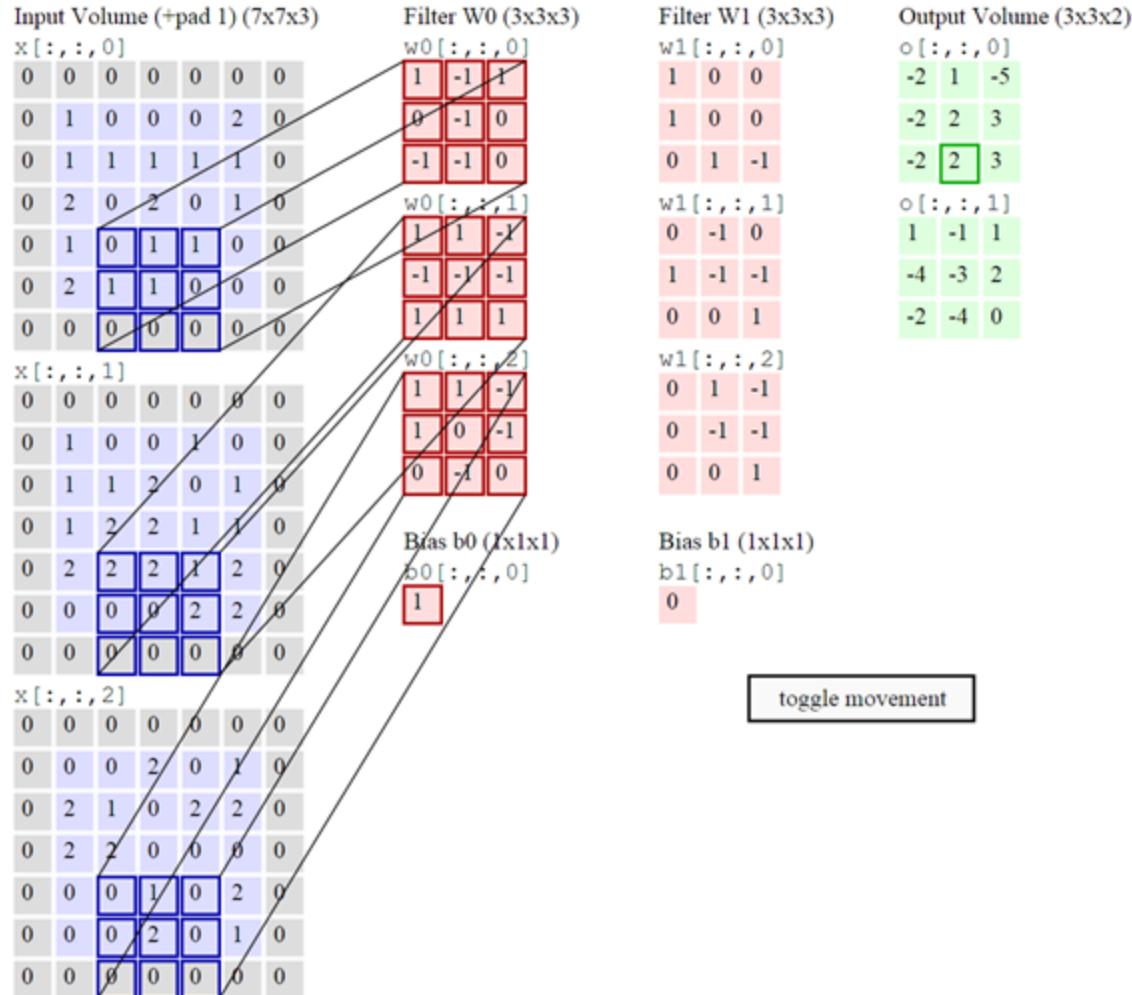
# Redes Neurais Convolucionais



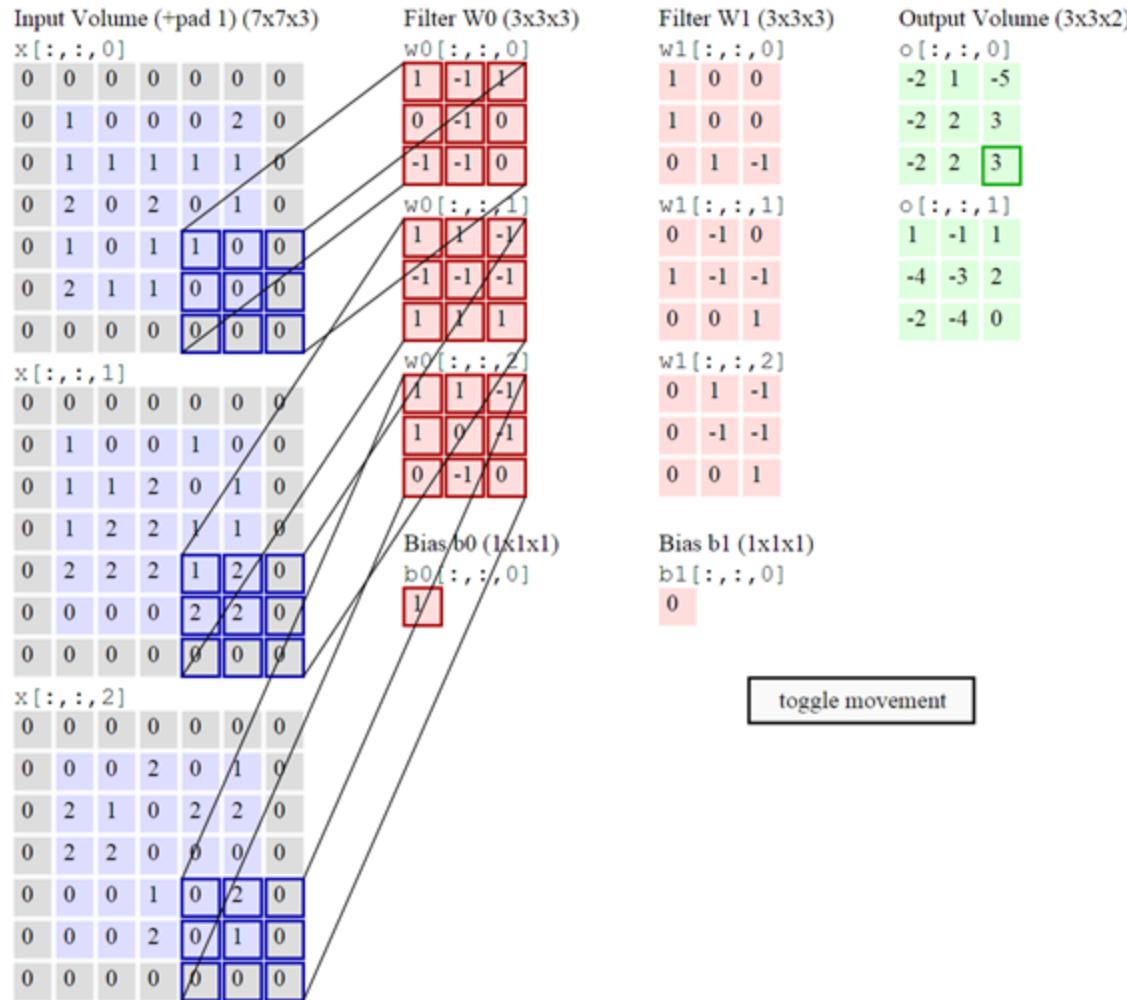
# Redes Neurais Convolucionais



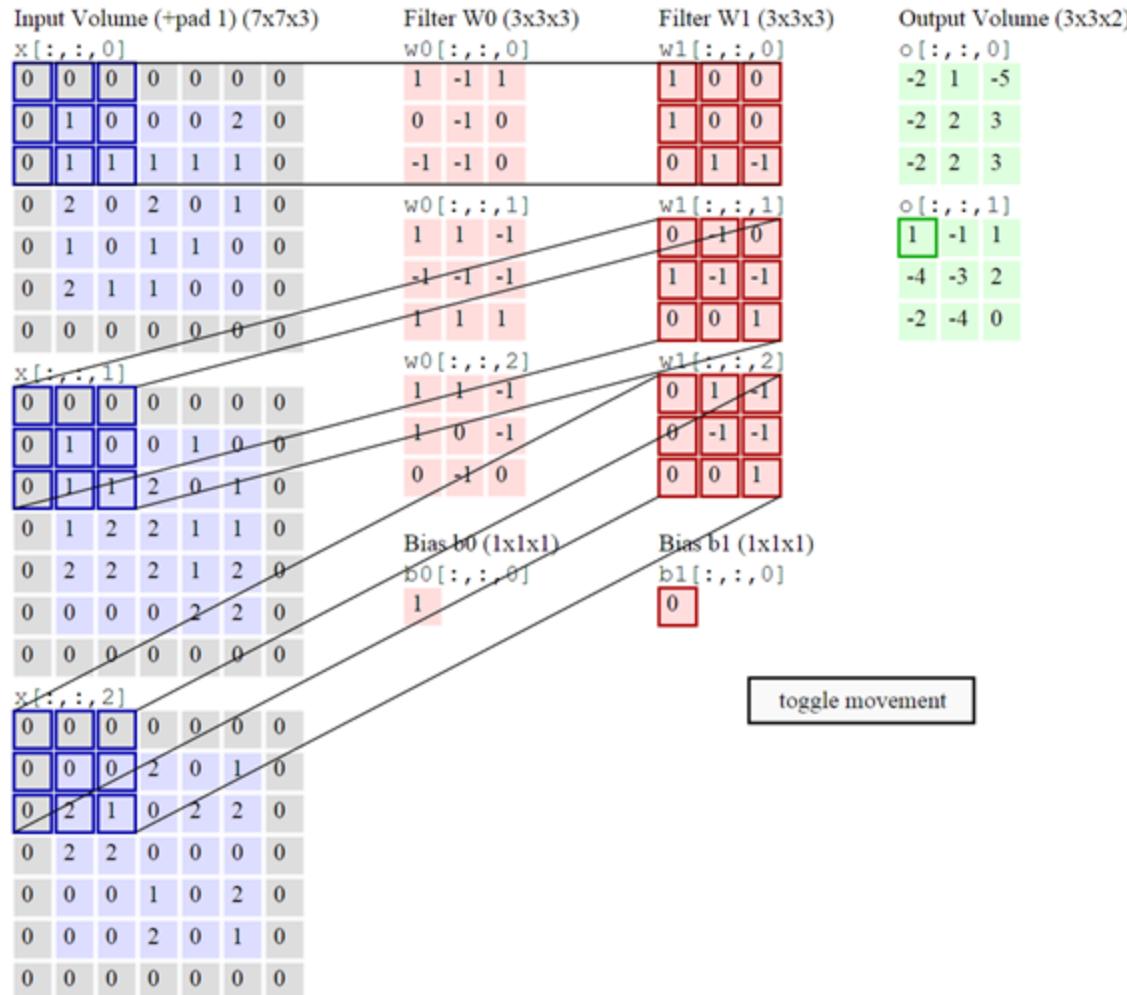
# Redes Neurais Convolucionais



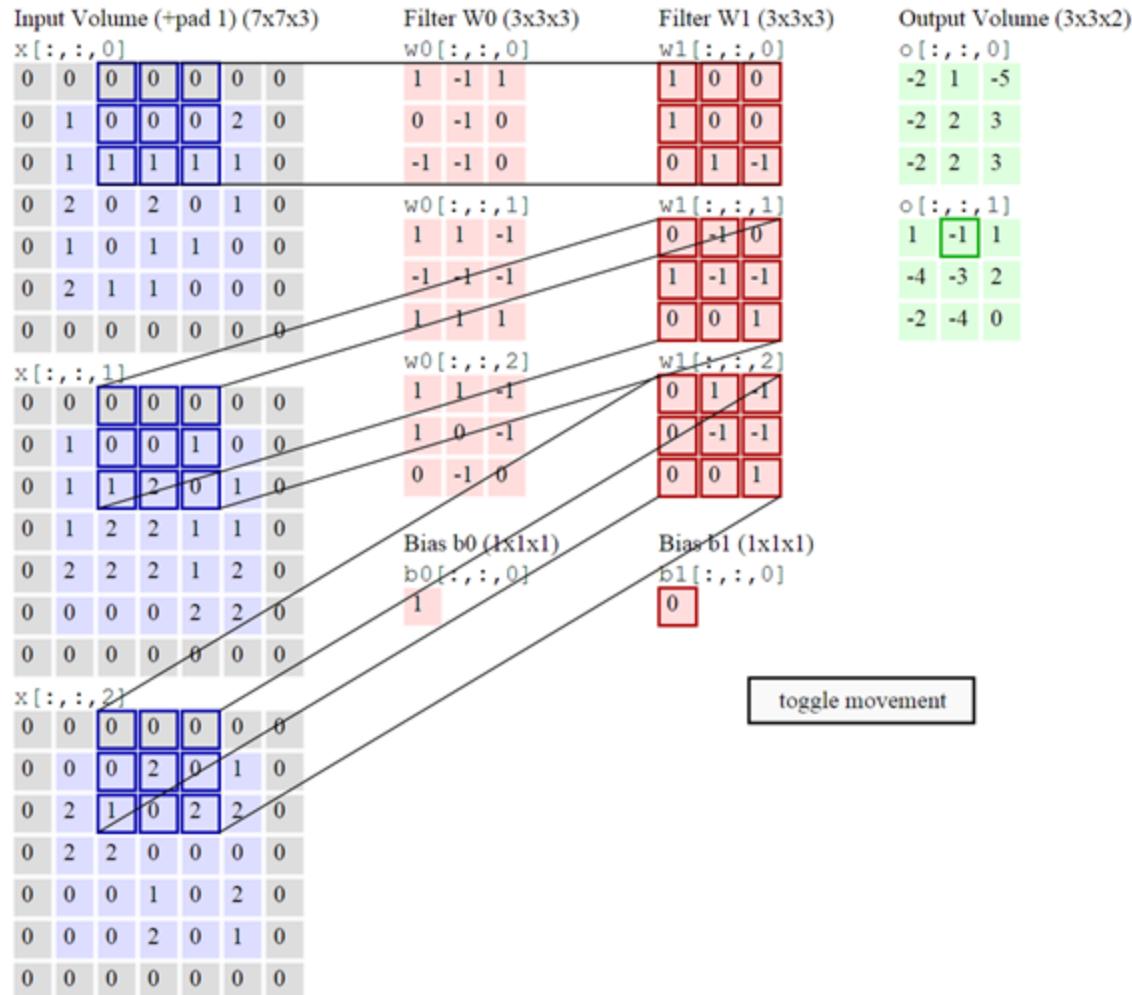
# Redes Neurais Convolucionais



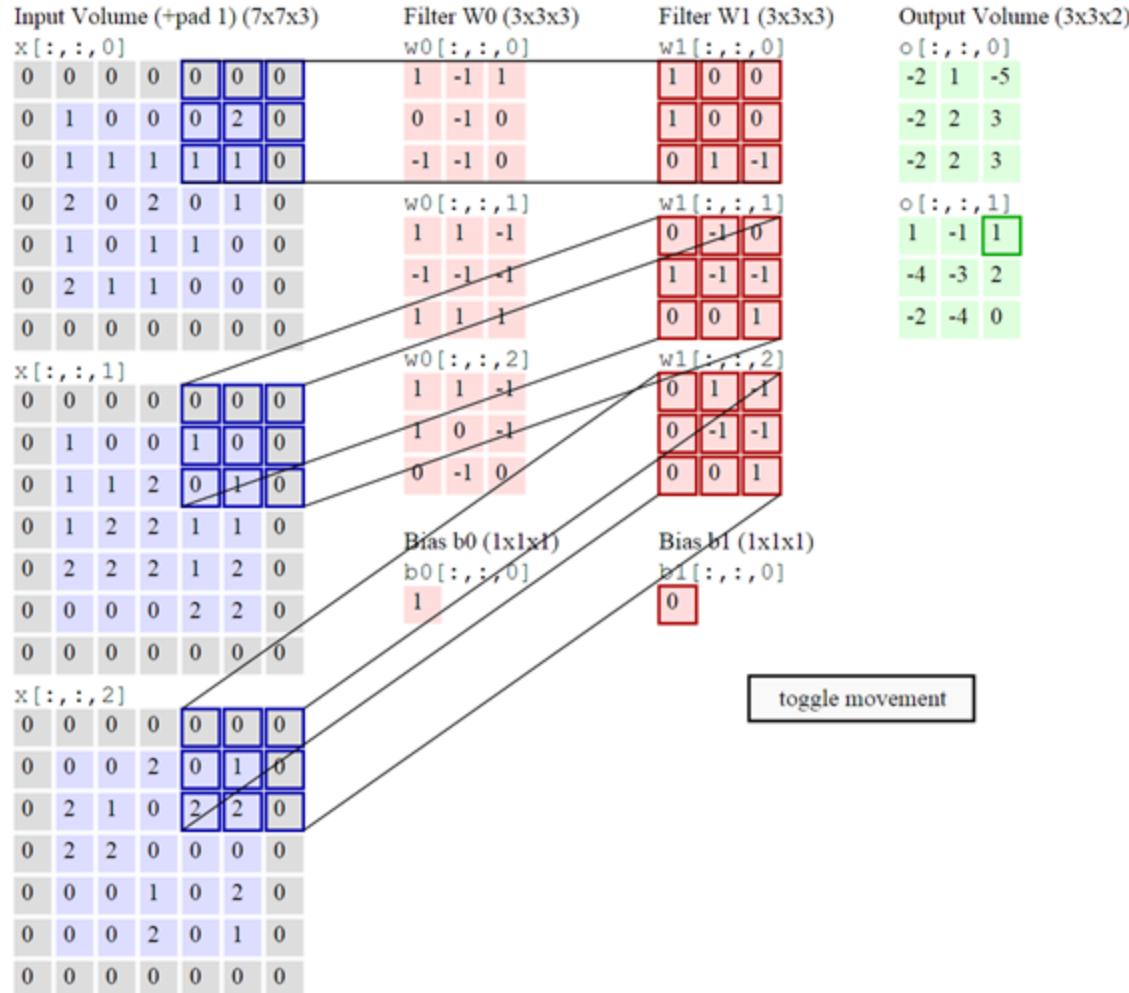
# Redes Neurais Convolucionais



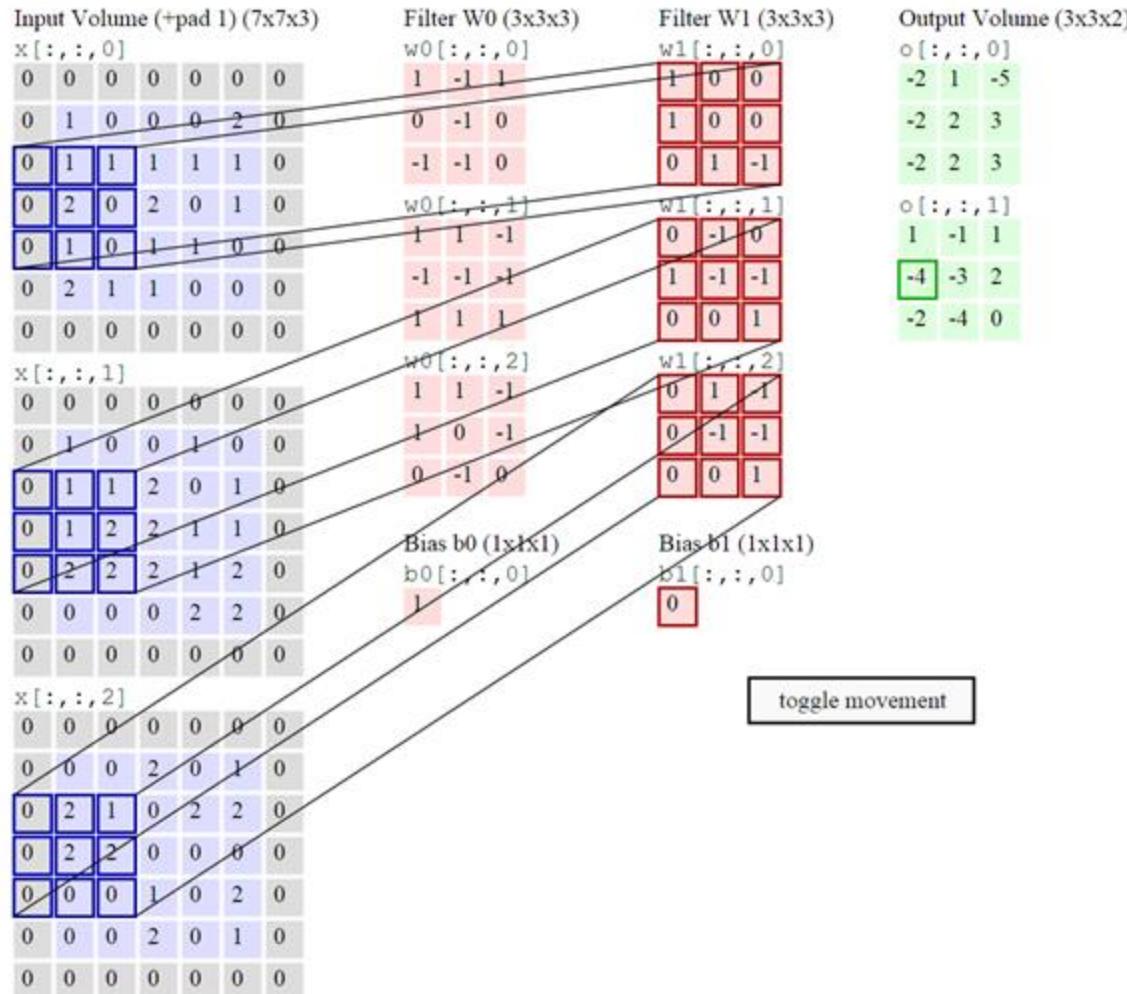
# Redes Neurais Convolucionais



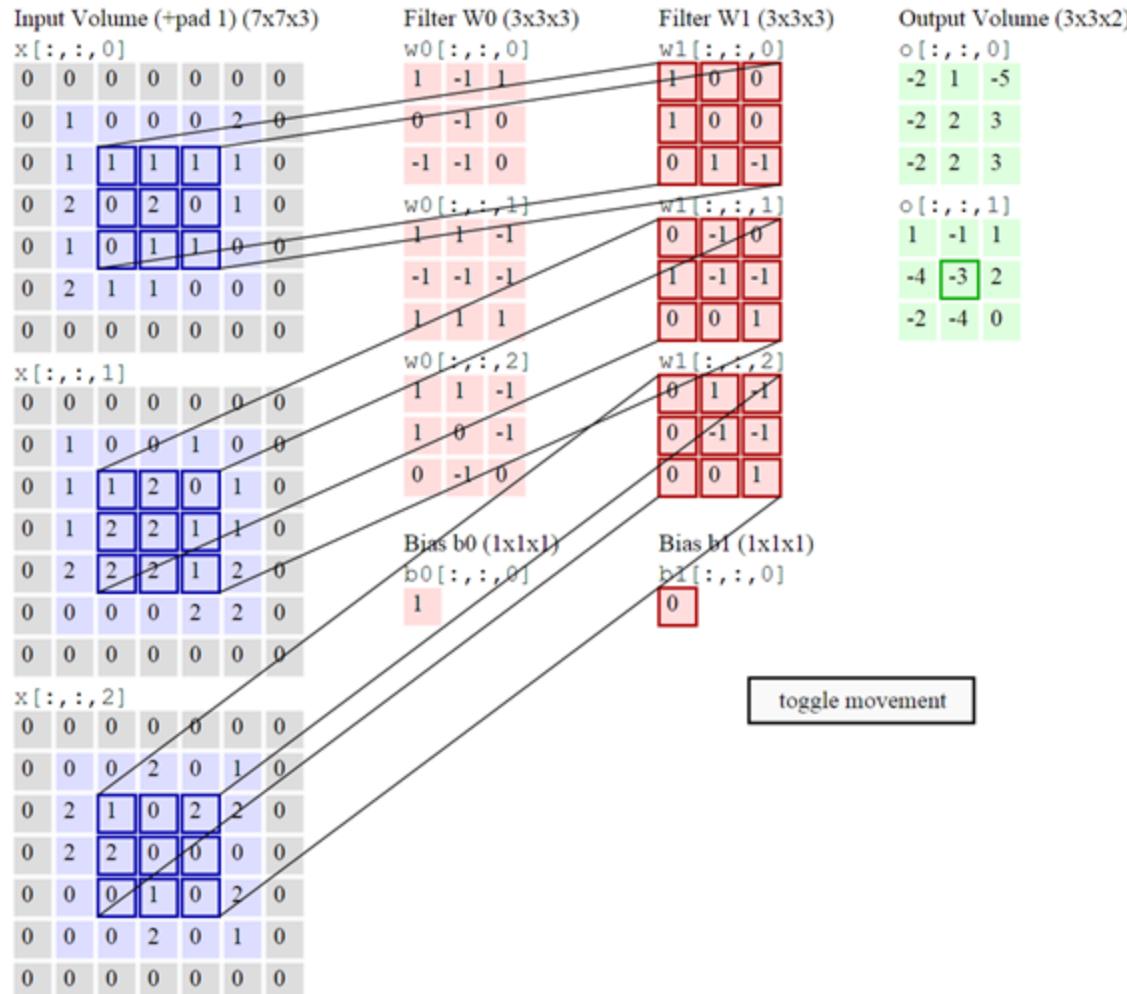
# Redes Neurais Convolucionais



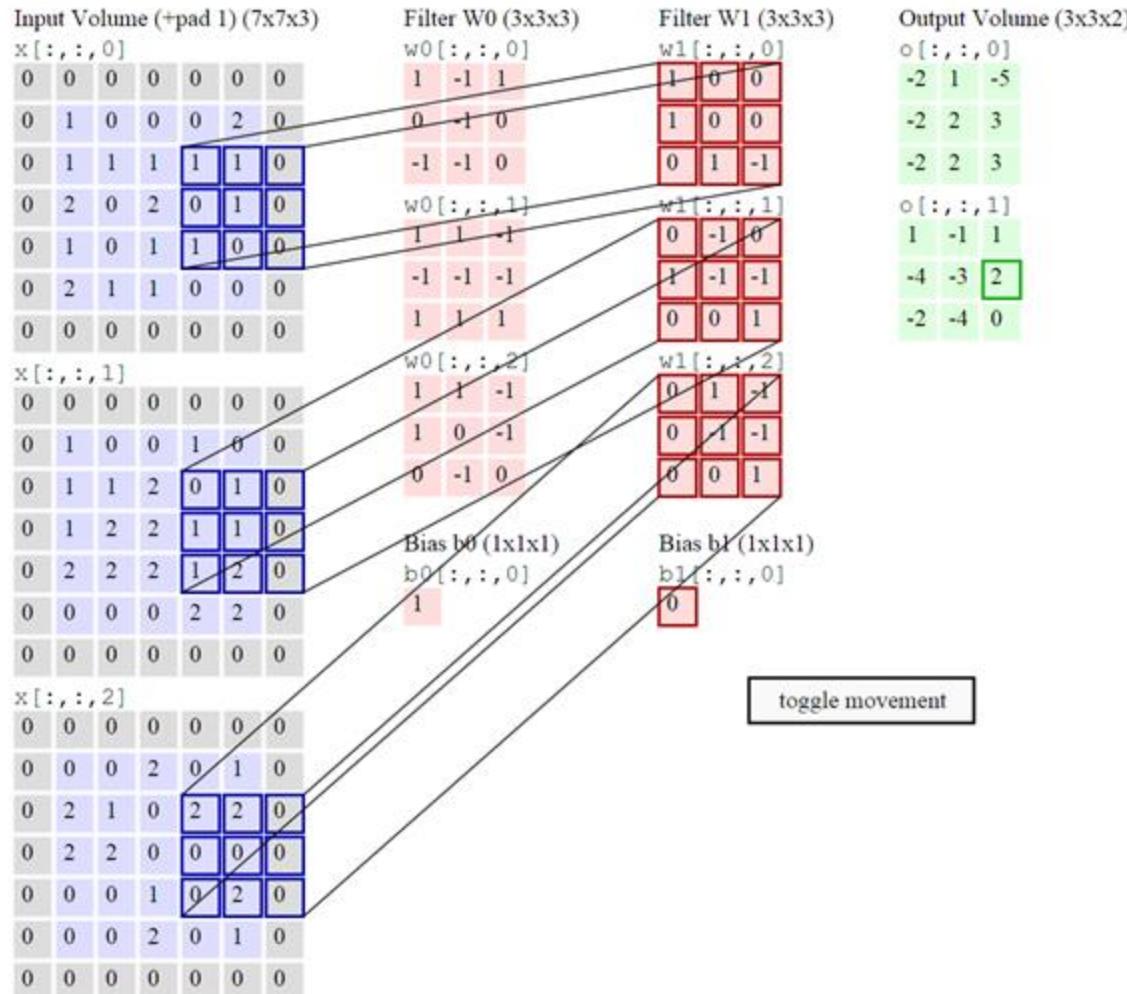
# Redes Neurais Convolucionais



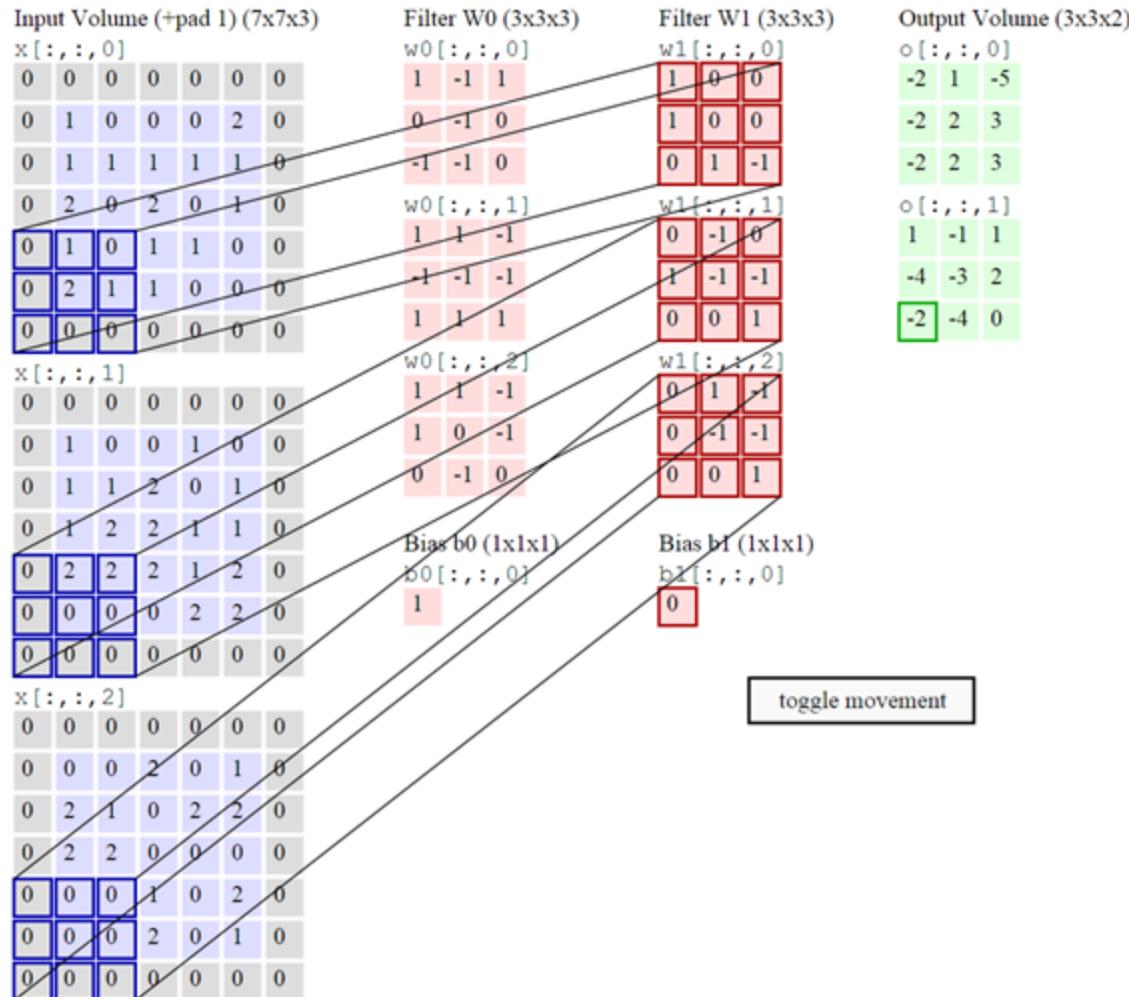
# Redes Neurais Convolucionais



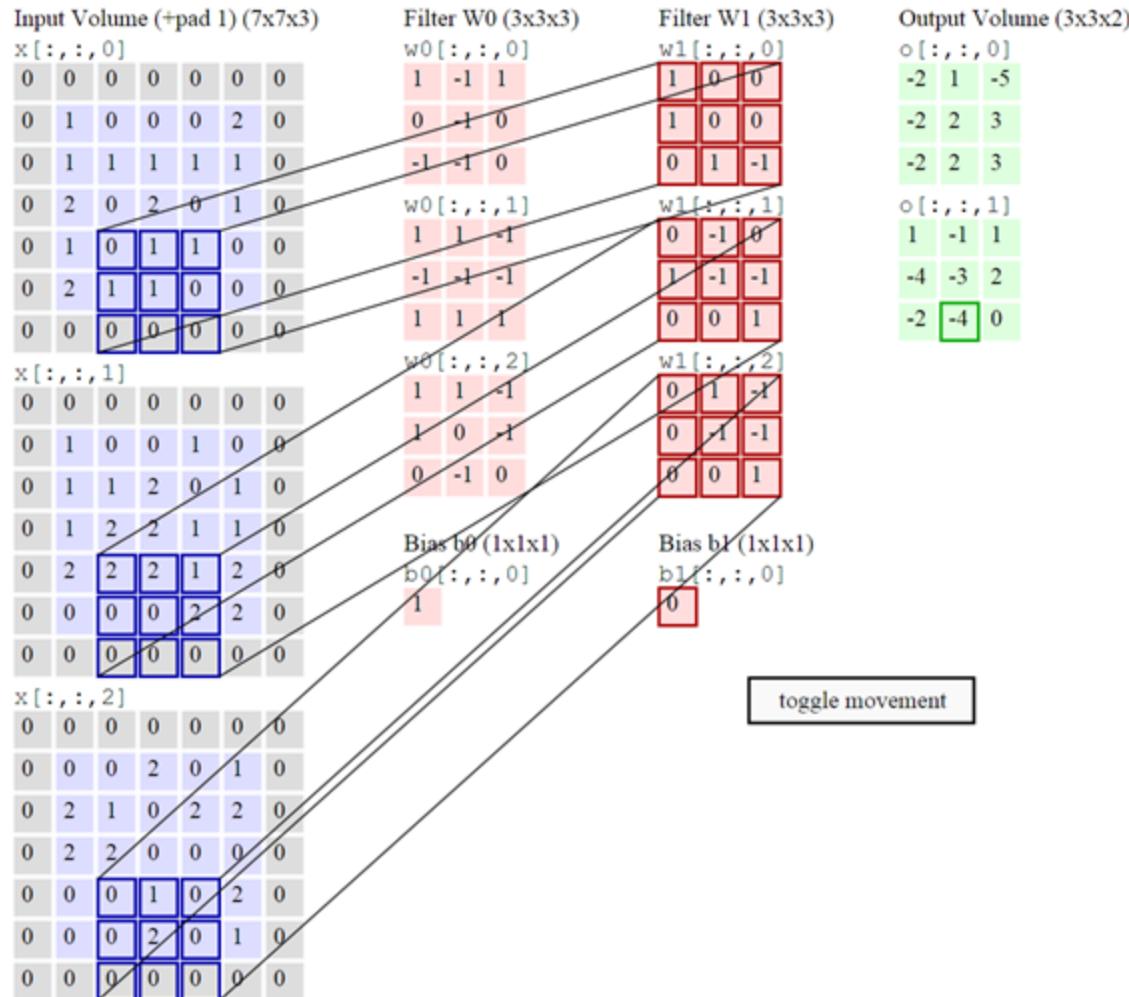
# Redes Neurais Convolucionais



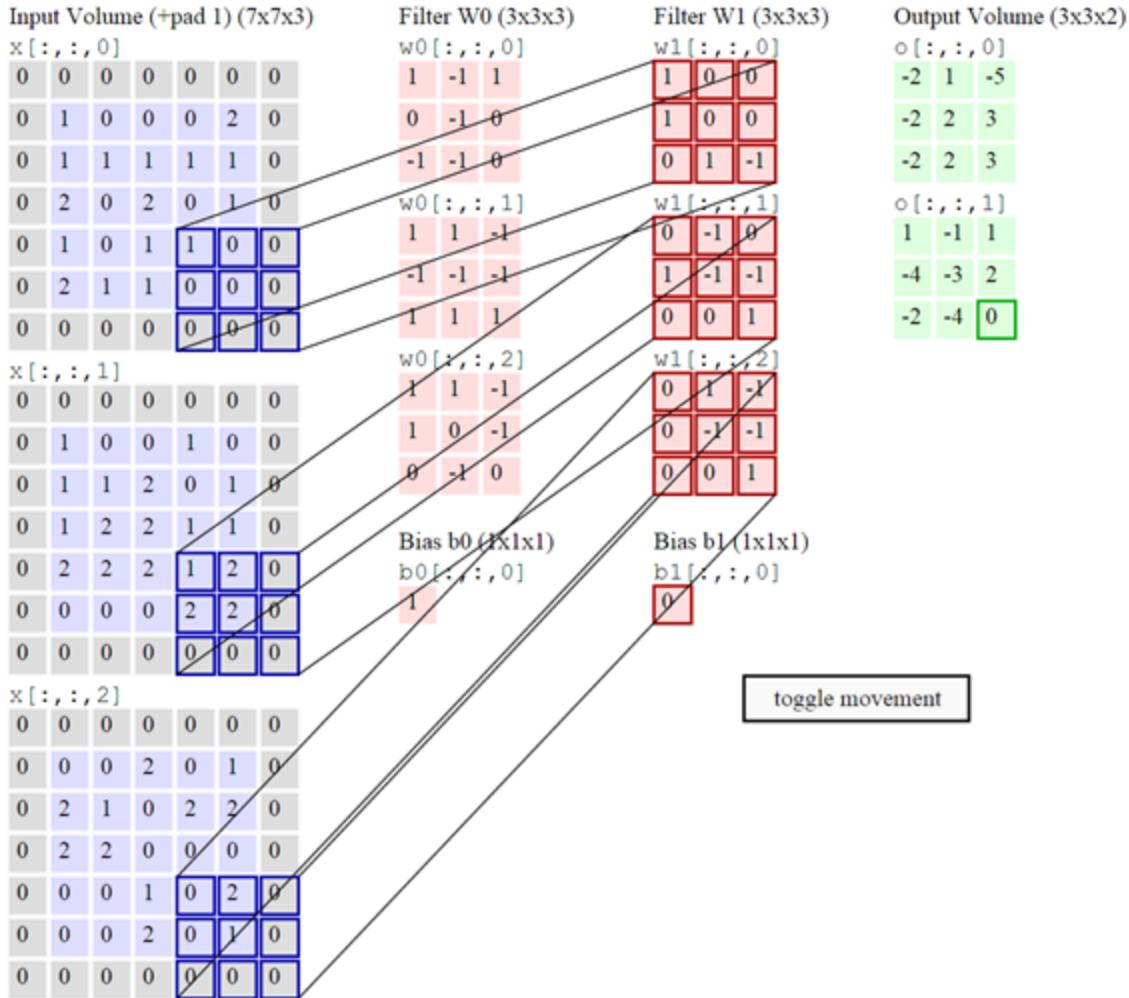
# Redes Neurais Convolucionais

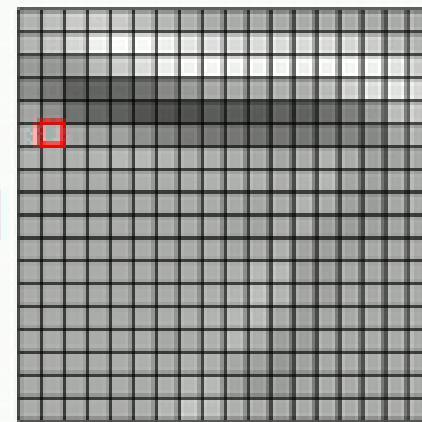
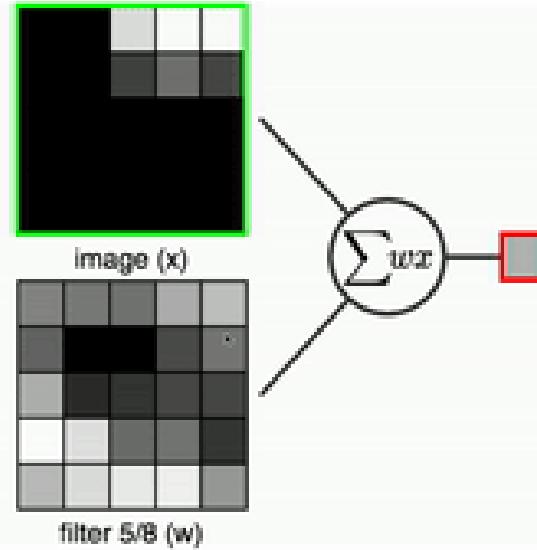
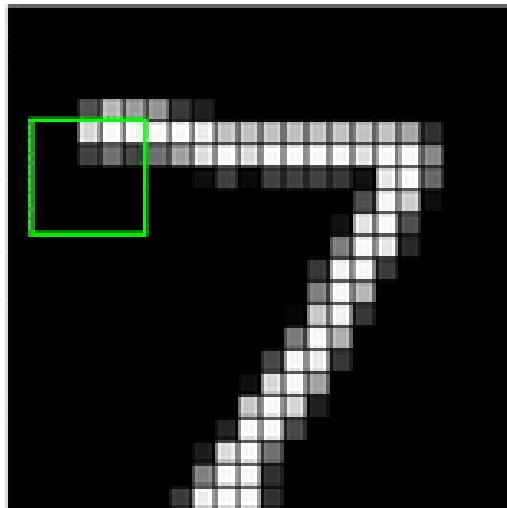


# Redes Neurais Convolucionais



# Redes Neurais Convolucionais

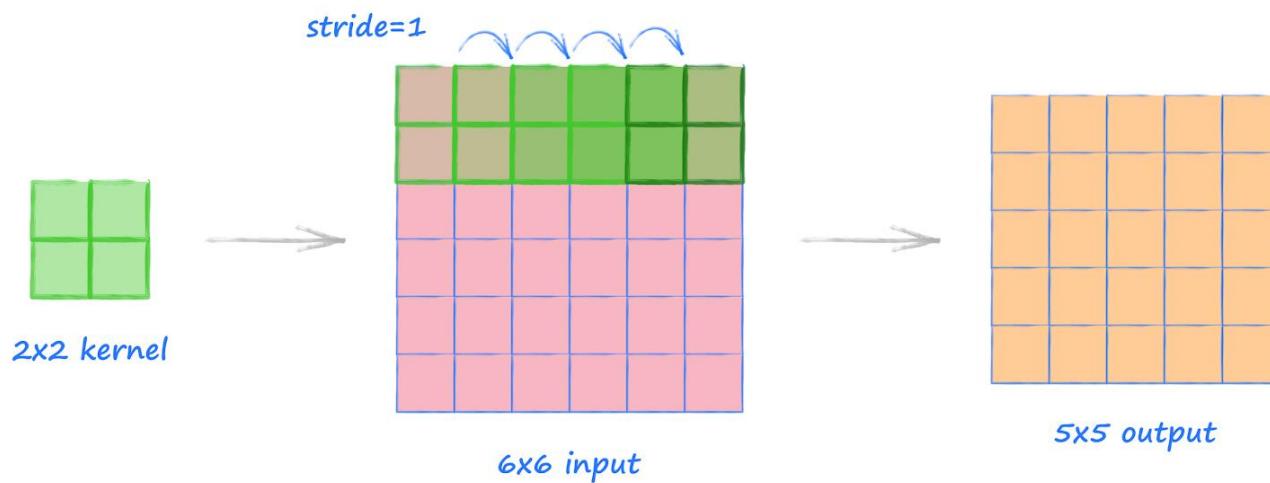




[http://ml4a.github.io/dev/demos/demo\\_convolution.html](http://ml4a.github.io/dev/demos/demo_convolution.html)

# Conceitos

- Padding = Preencher as bordas com 0
- Stride = Passo que é dado pelo filtro aplicado



# Vamos praticar?

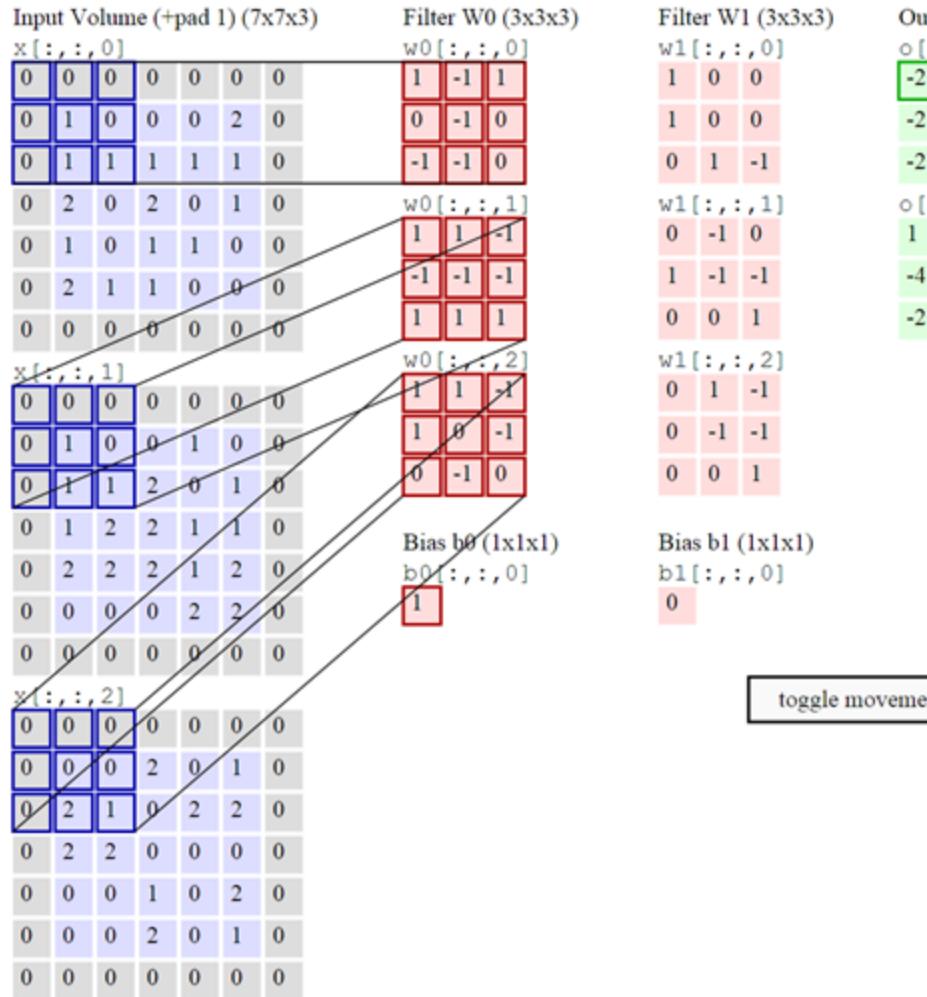
Input Volume (+pad 1) (7x7x3)	Filter W0 (3x3x3)
$x[:, :, 0]$  $\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 2 & 0 \\ 0 & 1 & 1 & 1 & 1 & 1 & 0 \\ 0 & 2 & 0 & 2 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 1 & 0 & 0 \\ 0 & 2 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$	$w0[:, :, 0]$ $\begin{bmatrix} 1 & -1 & 1 \\ 0 & -1 & 0 \\ -1 & -1 & 0 \end{bmatrix}$
$x[:, :, 1]$  $\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 2 & 0 & 1 & 0 \\ 0 & 1 & 2 & 2 & 1 & 1 & 0 \\ 0 & 2 & 2 & 2 & 1 & 2 & 0 \\ 0 & 0 & 0 & 0 & 2 & 2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$	$w0[:, :, 1]$ $\begin{bmatrix} 1 & 1 & -1 \\ 1 & 0 & -1 \\ 0 & -1 & 0 \end{bmatrix}$
$x[:, :, 2]$  $\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 & 1 & 0 \\ 0 & 2 & 1 & 0 & 2 & 2 & 0 \\ 0 & 2 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 2 & 0 \\ 0 & 0 & 0 & 2 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$	$w0[:, :, 2]$ $\begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & -1 \\ 0 & -1 & 0 \end{bmatrix}$
	Bias b0 (1x1x1) $b0[:, :, 0]$ $0$

Filter W1 (3x3x3)	Output Volume (3x3x2)
$w1[:, :, 0]$ $\begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & -1 \end{bmatrix}$	$o[:, :, 0]$ $\begin{bmatrix} -2 & 1 & -5 \\ -2 & 2 & 3 \\ -2 & 2 & 3 \end{bmatrix}$
$w1[:, :, 1]$ $\begin{bmatrix} 0 & -1 & 0 \\ 1 & -1 & -1 \\ 0 & 0 & 1 \end{bmatrix}$	$o[:, :, 1]$ $\begin{bmatrix} 1 & -1 & 1 \\ -4 & -3 & 2 \\ -2 & -4 & 0 \end{bmatrix}$
$w1[:, :, 2]$ $\begin{bmatrix} 0 & 1 & -1 \\ 0 & -1 & -1 \\ 0 & 0 & 1 \end{bmatrix}$	
Bias b1 (1x1x1) $b1[:, :, 0]$ $0$	

toggle movement

Número de parâmetros  
desta camada?

# Vamos praticar?

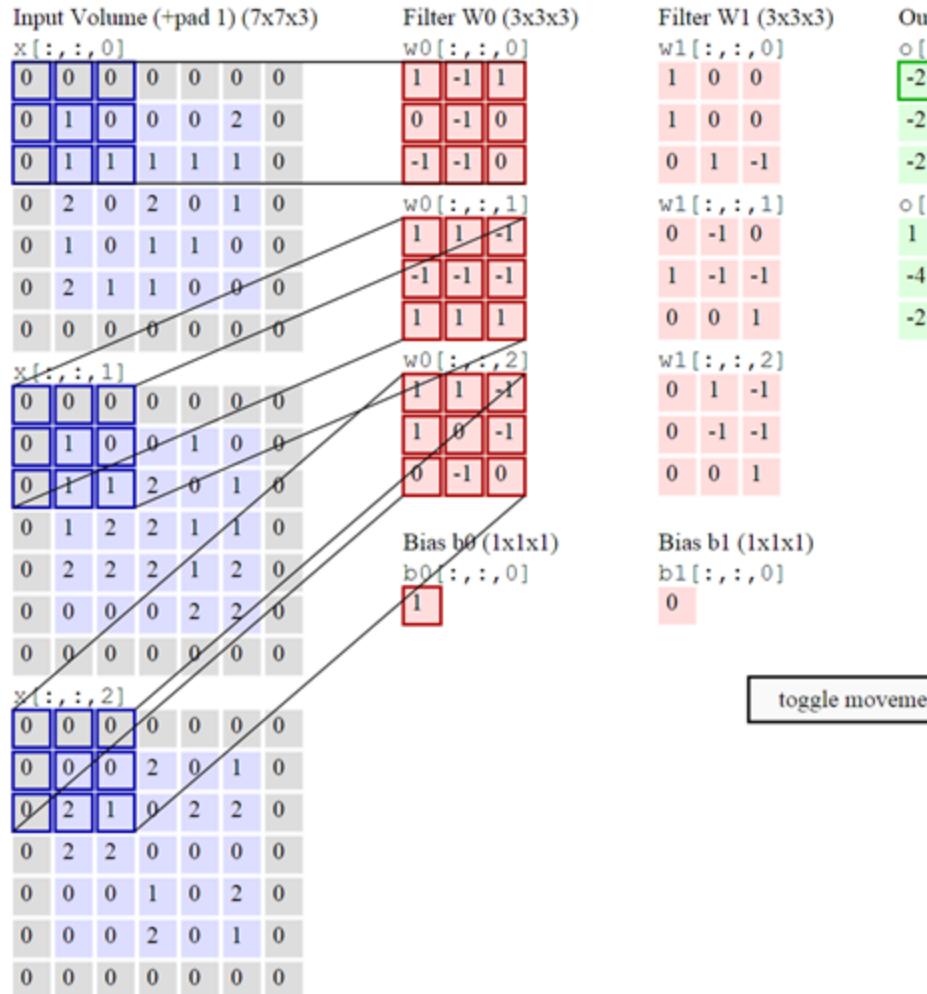


Número de parâmetros  
desta camada?

$$(3 \times 3 \times 3 + 1) \times 2 \\ = 56$$

toggle movement

# Vamos praticar?

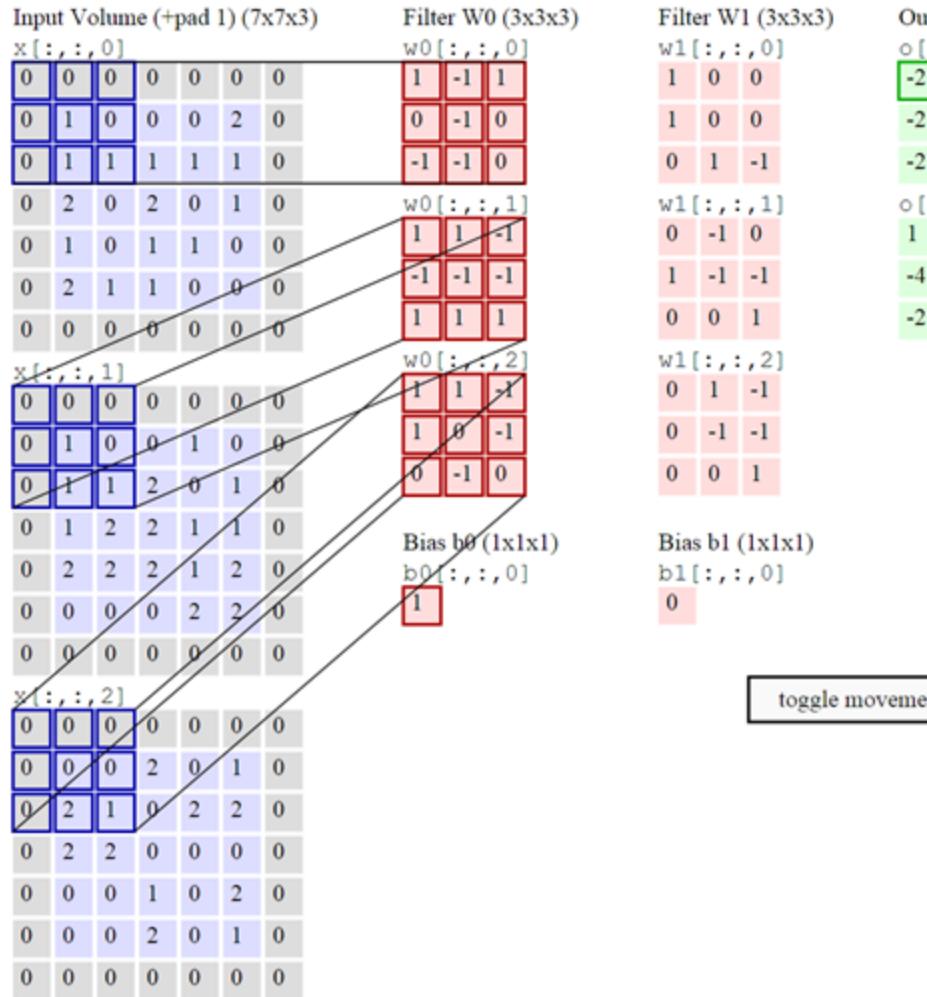


Número de parâmetros  
desta camada?

$$(3 \times 3 \times 3 + 1) \times 2 \\ = 56$$

E se utilizássemos uma  
camada totalmente  
conectada com 18 neurônios?

# Vamos praticar?

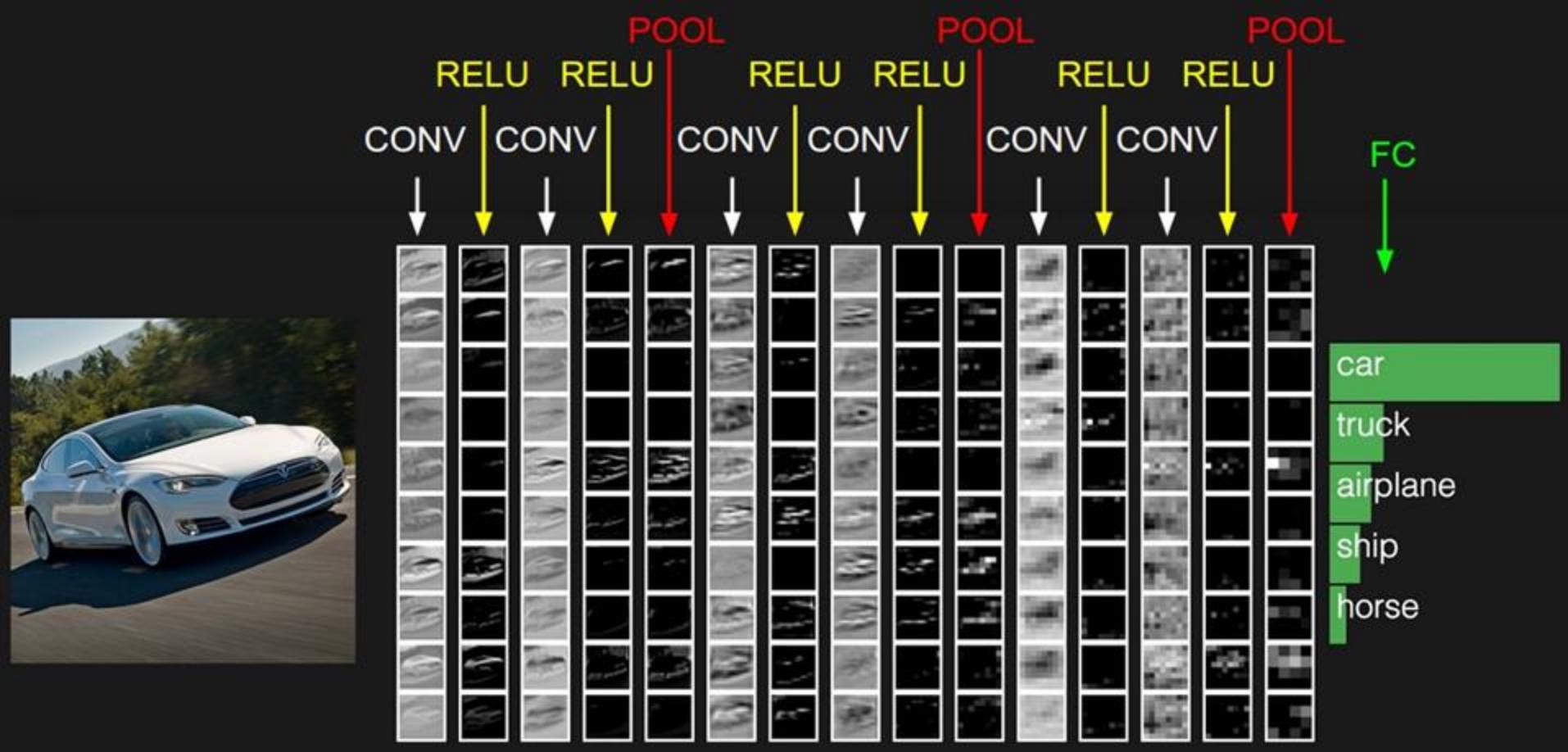


Número de parâmetros  
desta camada?

$$(3 \times 3 \times 3 + 1) \times 2 \\ = 56$$

E se utilizássemos uma  
camada totalmente  
conectada com 18 neurônios?

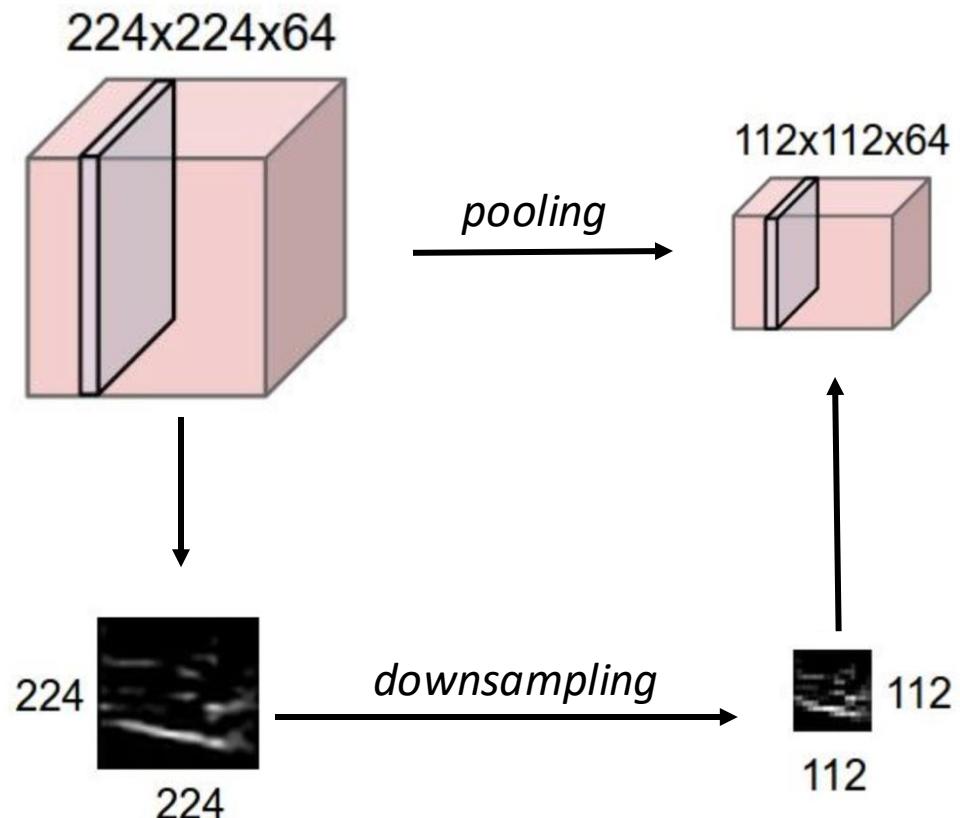
$$18 \times (5 \times 5 \times 3) + 18 \\ = 1368$$



Próxima camada: Pooling!

# Camada de *Pooling*

- **Reduz os volumes** para que se tornem gerenciáveis
- Atua sobre cada mapa de ativação de **maneira independente**



# MAX Pooling

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pooling com  
filtro  $2 \times 2$  e stride 2



6	8
3	4

1 único mapa de ativação

# AVERAGE (AVG) Pooling

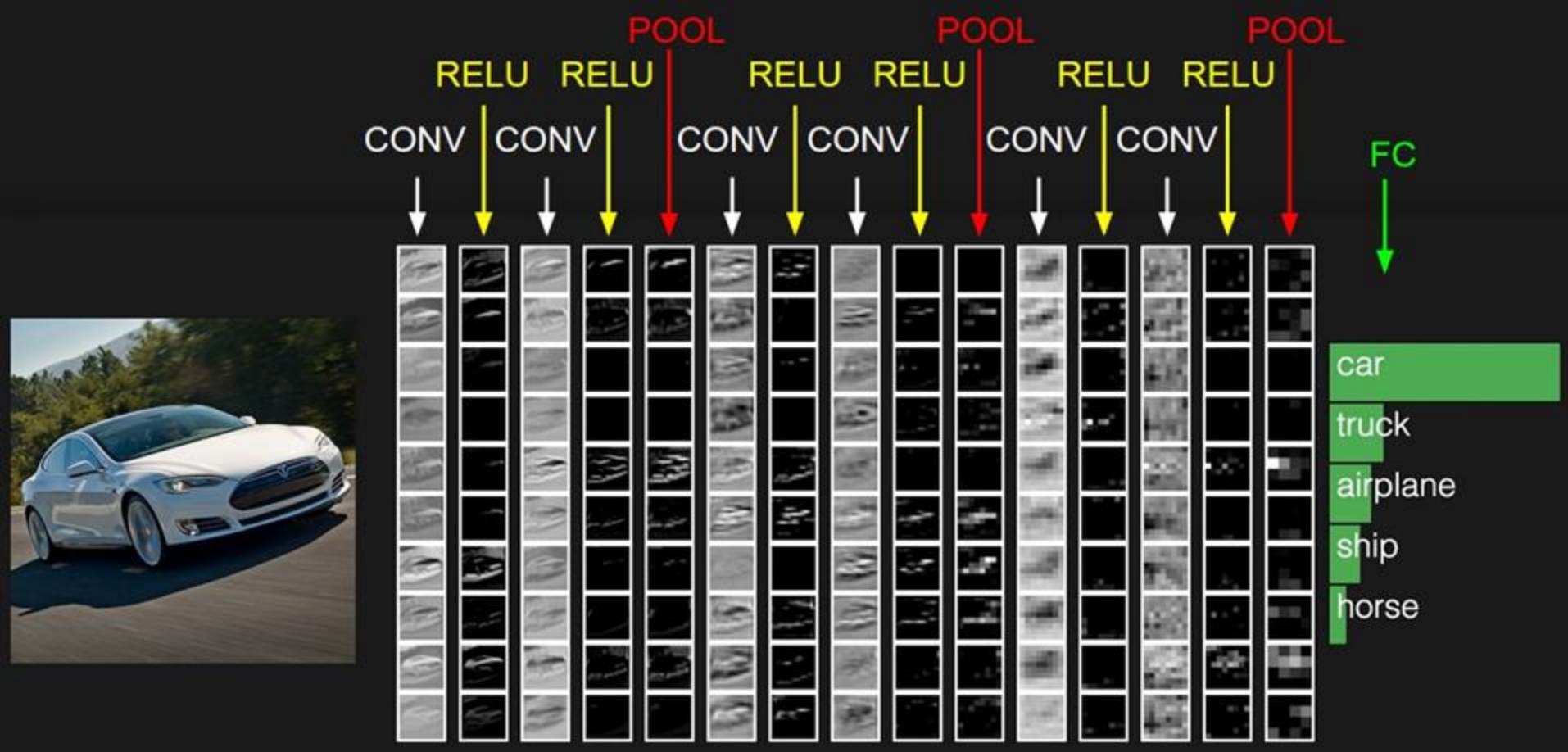
2	3	2	5
5	6	9	8
3	2	1	2
1	2	0	1

*avg pooling com  
filtros  $2 \times 2$  e stride 2*



4	6
2	1

**1 único mapa de ativação**

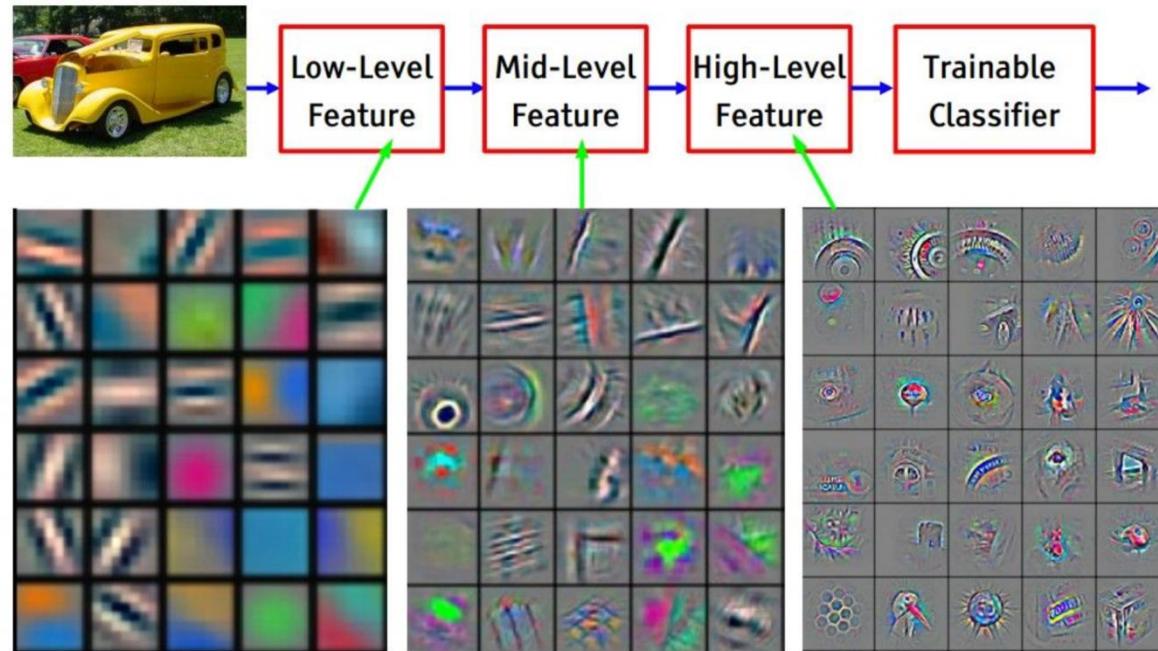


## Finalmente: FC (Fully-Connected Layer)

Uma rede convolucional pode conter uma (ou várias) camadas totalmente conectadas, onde os neurônios estão conectados com toda a entrada (redes neurais tradicionais)

# Hierarquia de Features

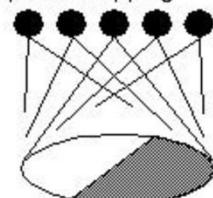
Slides Yann LeCun



Visualização de Features (Rede treinada no ImageNet) [Zeiler & Fergus, 2013]

## Hubel & Weisel

topographical mapping

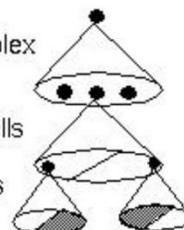


## featural hierarchy

hyper-complex  
cells

complex cells

simple cells



high level

mid level

low level

# Exemplos de Arquiteturas: LeNet-5

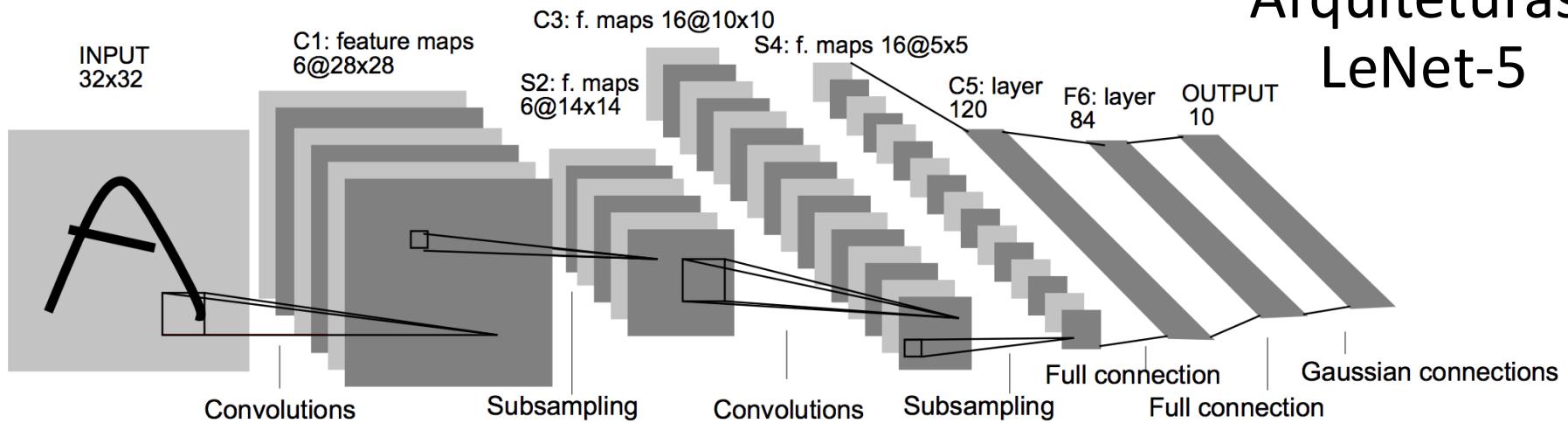


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Filtros Convolucionais:  $5 \times 5$   
Stride: 1

Filtros de Pooling:  $2 \times 2$   
Stride: 2

Arquitetura: [CONV-POOL-CONV-POOL-FC1-FC2-FC3]

Y. LeCun, L. D. Jackel, L. Bottou, A. Brunot, C. Cortes, J. S. Denker, H. Drucker, I. Guyon, U. A. Muller, E. Sackinger, P. Simard, and V. Vapnik. Comparison of learning algorithms for handwritten digit recognition. In F. Fogelman and P. Gallinari, editors, *International Conference on Artificial Neural Networks*, pages 53-60, Paris, 1995.

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, november 1998.

# Exemplos de Arquiteturas: VGGNet

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 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	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 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Table 2: Number of parameters (in millions).

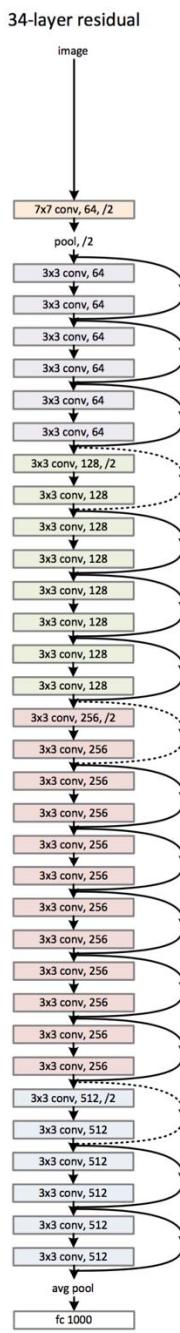
Network	A,A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144

Apenas convoluções  $3 \times 3$  (stride 1, pad 1) e *max pooling*  $2 \times 2$  (stride 2)

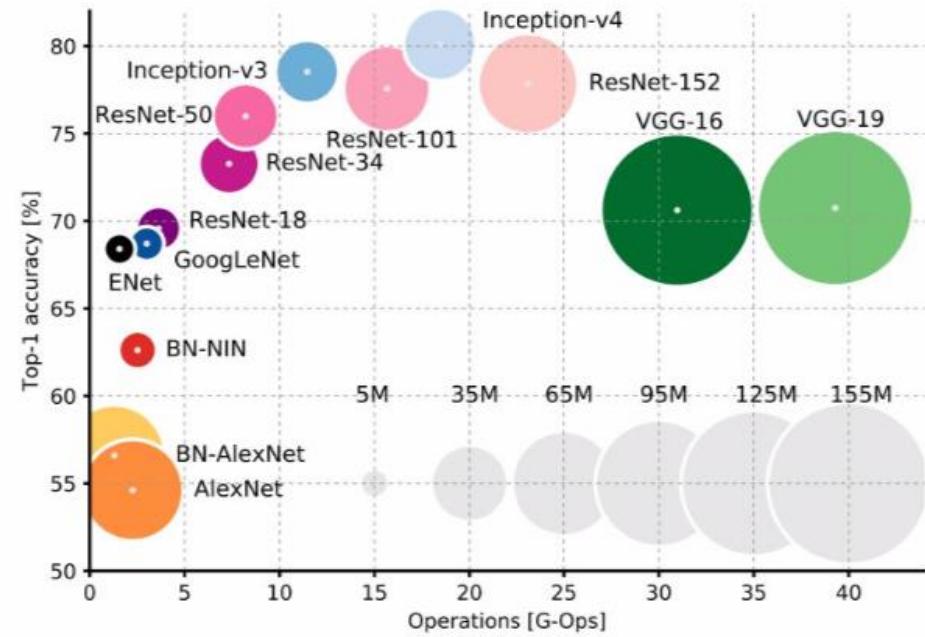
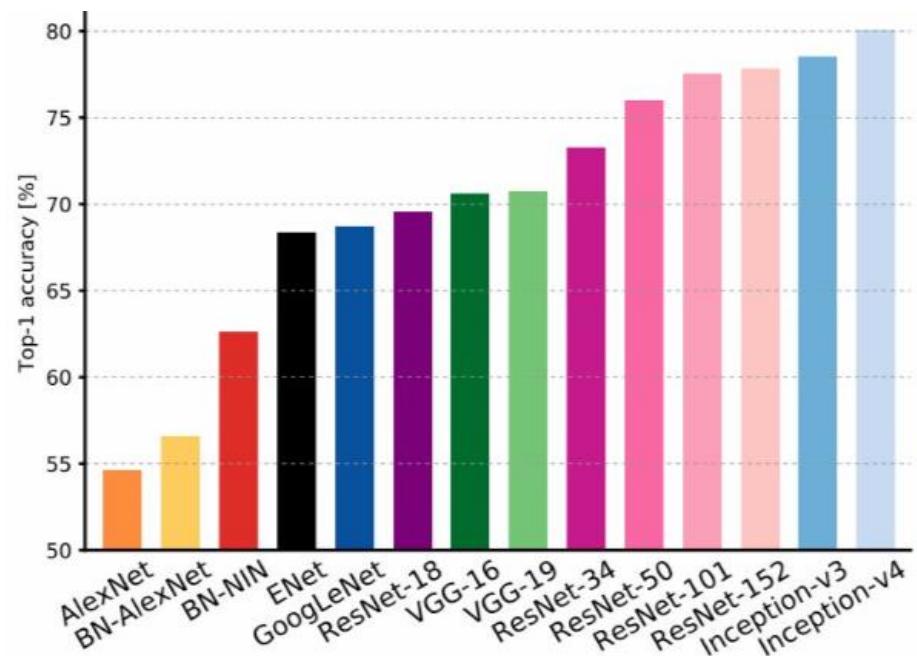
melhor modelo

Erro top-5 ImageNet:  
7.0% (uma única rede)  
6.8% (ensemble com 2 redes)

# Exemplos de Arquiteturas: ResNet



# Comparação de Complexidade



# *Transfer Learning*

“Precisamos utilizar uma base de dados gigante”

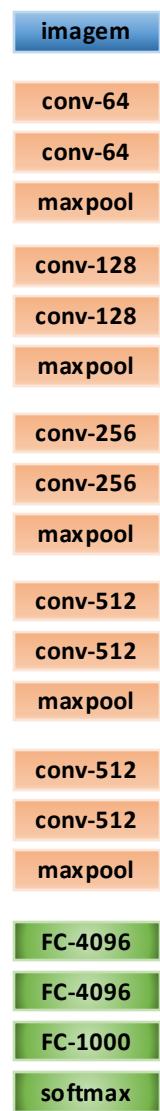


“Precisamos treinar uma ConvNet do zero”



*“Transfer learning (inductive transfer) é uma área do aprendizado de máquina que foca em armazenar conhecimento adquirido ao se resolver determinada tarefa para posterior aplicação em problema diferente mas relacionado.”*

# Transfer Learning em ConvNets



(Pré-) treinar no ImageNet (ILSVRC-2012)

<https://github.com/tensorflow/models>

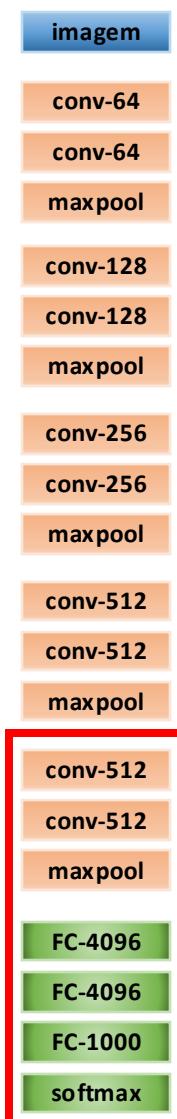
<https://github.com/pytorch/vision>



Se dataset for pequeno, utilizar rede como **extrator de features!**

"congelar" estas camadas

Treinar aqui!  
(classificador linear ou jogar *features* num SVM)



Dataset médio: **fine-tuning!**

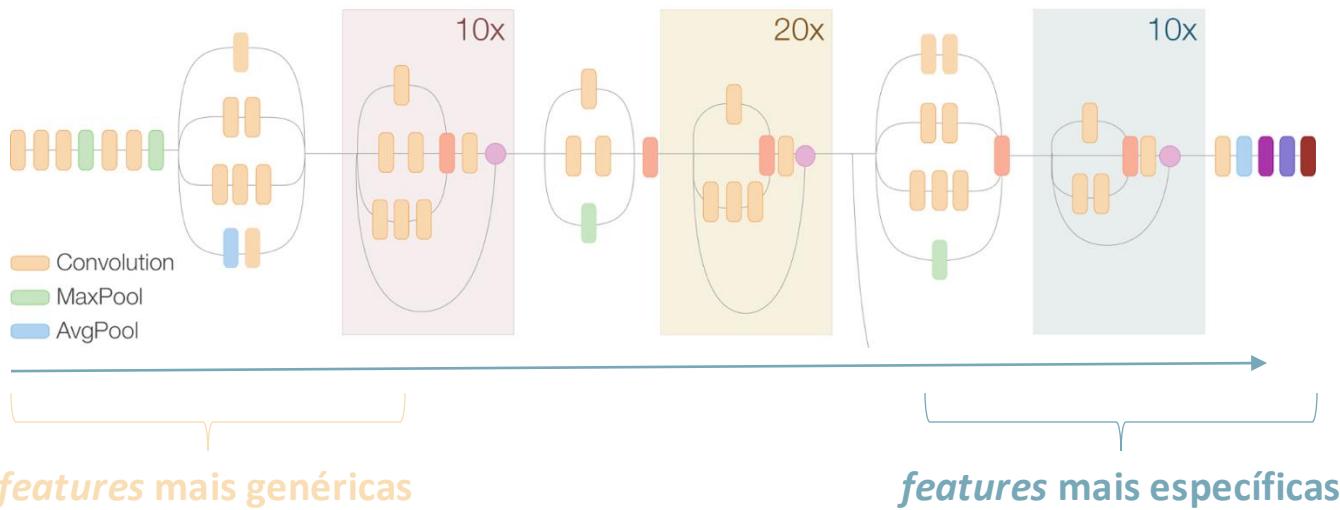
Quanto mais dados, mais rede pode ser re-treinada...

"congelar" estas camadas

Treinar aqui!

# Transfer Learning em ConvNets

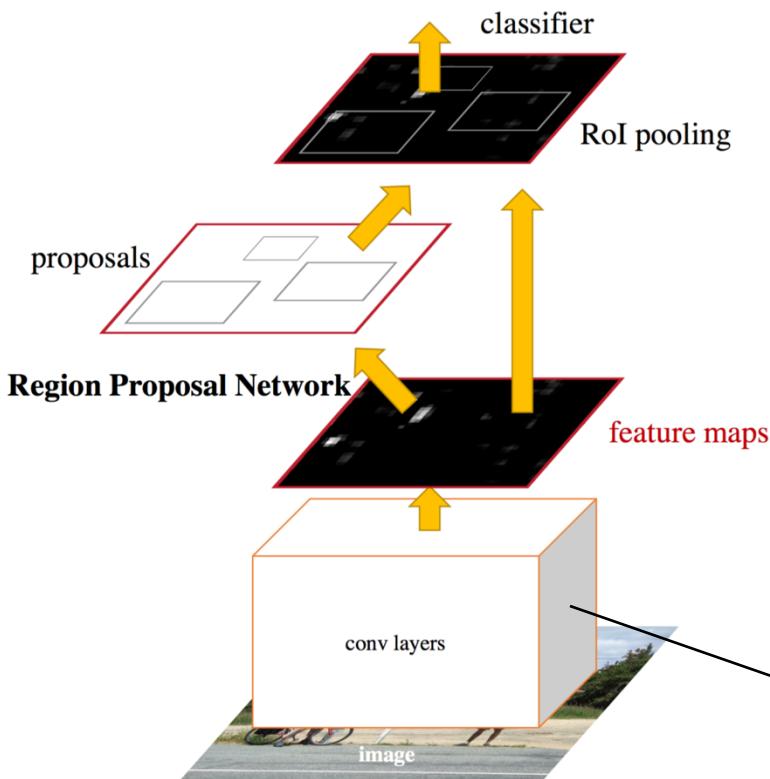
Compressed View



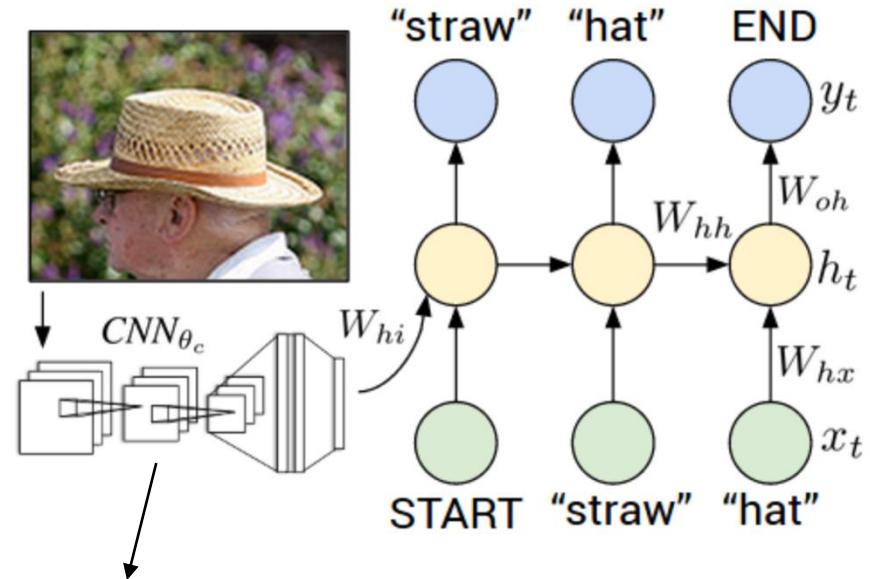
	Dados similares ao contexto em questão	Dados muito diferentes
Poucos dados	Usar um classificador linear na saída	
Muitos dados	Fine-tune algumas camadas...	Fine-tune algumas camadas a mais...

# Transfer Learning em ConvNets

Transfer Learning em ConvNets é a regra (e não exceção!!)



Andrej Karpathy, Fei-Fei Li. Deep visual-semantic alignments for generating image descriptions. CVPR 2015, 3128-3137



redes pré-treinadas no ImageNet!!

Shaoqing Ren, Kaiming He, Ross B. Girshick, Jian Sun.  
**Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks.** NIPS 2015, 91-99

# *Transfer Learning* em ConvNets

## Conclusão:

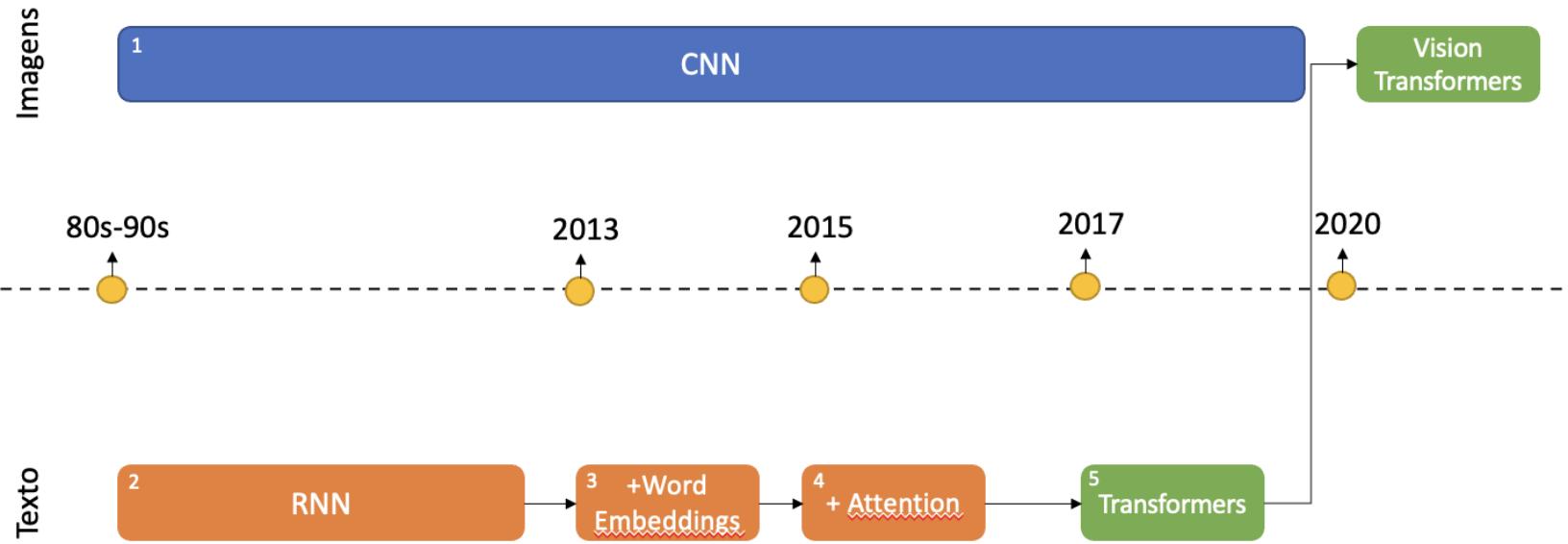
1. Encontre dataset muito grande (>1M) com dados similares aos do seu interesse e (pré-) treine uma ConvNet lá! (ou pegue o modelo treinado de algum repositório de modelos, e.g. model zoo)
2. Faça transfer learning para o seu dataset (atenção para escolha da estratégia mais adequada de acordo com o que foi apresentado)

Caffe → <https://github.com/BVLC/caffe/wiki/Model-Zoo>

TensorFlow → <https://github.com/tensorflow/models>

PyTorch → <https://github.com/pytorch/vision>

# *Deep Learning*



# Leituras Sugeridas

- Capítulo 9 do livro:

