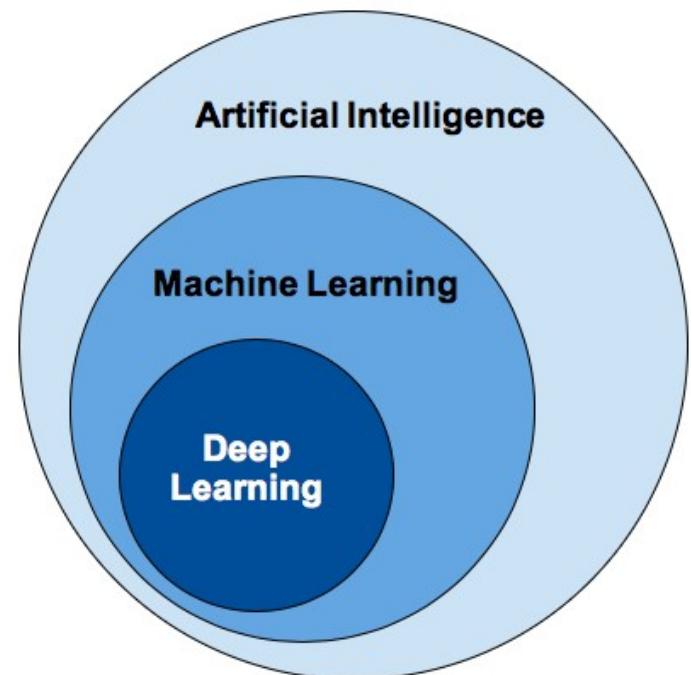


# Redes profundas y redes convolucionales

Basado en  
charlas de  
Lucas C. Uzal y  
otros



# “Revolution of Depth”

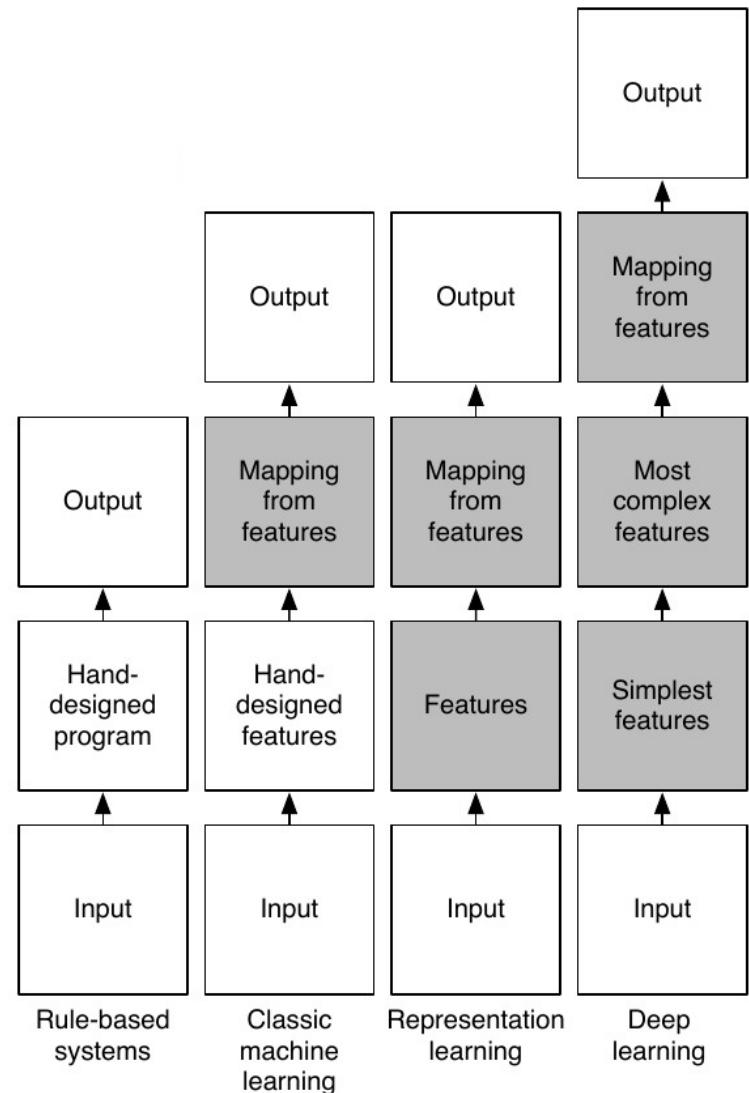
## Situación a mitad de los “2000”

- Machine Learning está establecido como campo de investigación con alto potencial tecnológico
- Las SVM dominan el campo
- Tipo de solución basado en el conocimiento de expertos:
  - Medir cosas (features) sobre los datos, sugeridas por conocimiento del problema
  - Aplicar un método de ML a esos features para decidir
- Las redes neuronales están en un “invierno”:
  - Se sabe que las redes neuronales son aproximadores universales.
  - Se sabe que las redes tienen el potencial de aprender su propia representación.
  - Problemas irresolubles: limitaciones para resolver problemas reales, dificultad para utilizarlas, falta de garantías de convergencia, lentitud de los entrenamiento

# “Revolution of Depth”

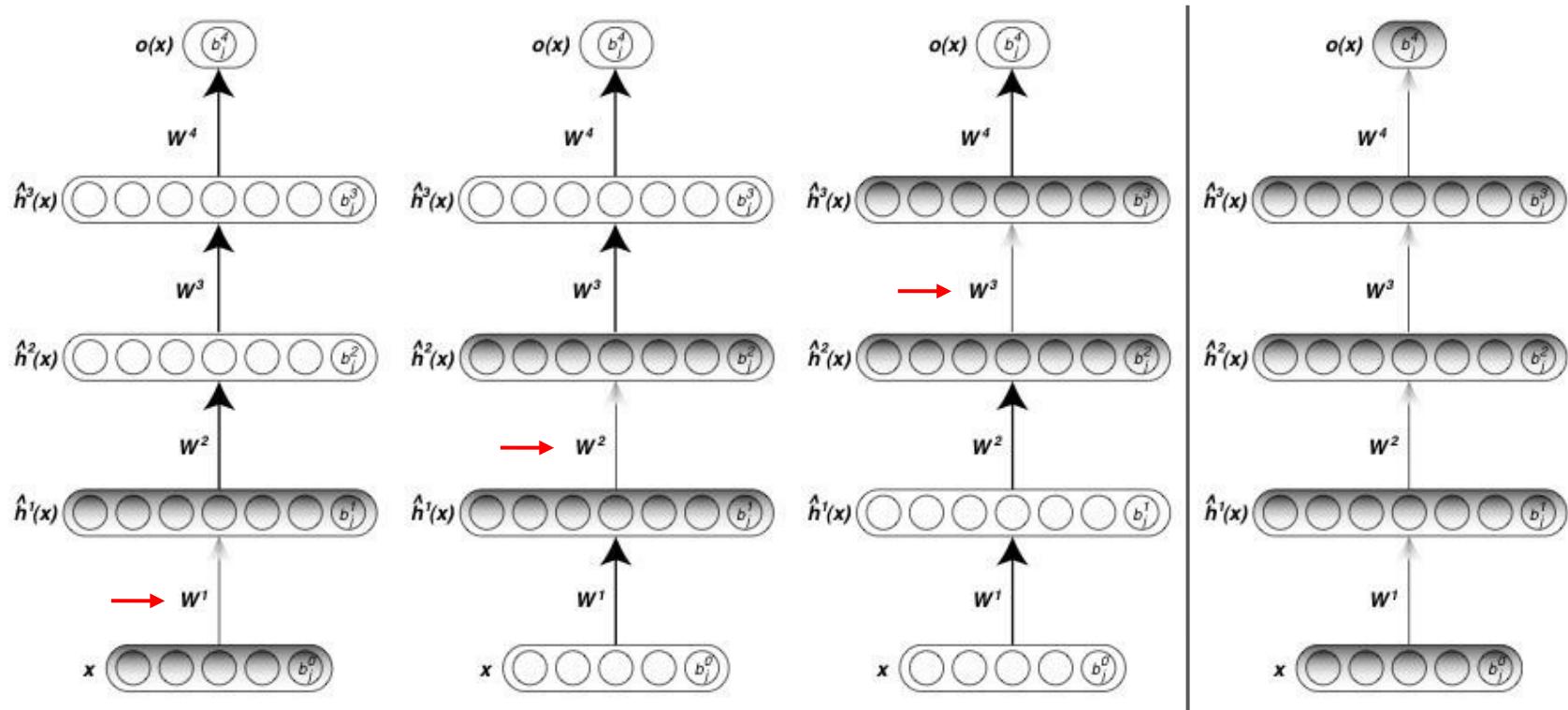
Por qué no seguir el mismo camino?

- “Aprender” las reglas produjo sistemas que solucionan las deficiencias de los sistemas expertos tradicionales.
- “Aprender” los features podría solucionar los problemas actuales.
- Si los métodos actuales no aprenden bien, por qué no usar métodos que puedan ir de lo simple a lo complejo?



No se sabía cómo hacer que las redes aprendan de esta forma!

# 2006: Greedy layer-wise training of deep networks



Hinton, Geoffrey E, and Ruslan R Salakhutdinov. "Reducing the dimensionality of data with neural networks." *Science* 313.5786 (2006): 504-507.

Hinton, Geoffrey, Simon Osindero, and Yee-Whye Teh. "A fast learning algorithm for deep belief nets." *Neural computation* 18.7 (2006): 1527-1554.

Bengio, Yoshua et al. "Greedy layer-wise training of deep networks." *Advances in neural information processing systems* 19 (2007): 153.

# “Revolution of Depth”

Primera solución al problema de entrenar una red profunda.

- Hinton muestra que se puede encontrar una primera solución entrenando de manera no supervizada, y afinando la solución (fine tuning) después, pero en redes muy particulares.
- Resurgen las redes como método, y hay un avance exponencial que lleva a la actual revolución tecnológica.
- En pocos años se encuentra como entrenar redes generales, y como extender todo lo que ya se sabía hacer a modelos profundos.

**3 claves: Datos (internet), Potencia (GPUs), Métodos (activaciones, regularización).**

# Google AI algorithm masters ancient game of Go

Deep-learning software defeats human professional for first time. x



**January 2016**

*A computer has beaten a human professional for the first time at Go — an ancient board game that has long been viewed as one of the greatest challenges for artificial intelligence (AI)*

*"We pass in the board position as a  $19 \times 19$  image and use convolutional layers to construct a representation of the position."*

Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., van den Driessche, G., ... & Dieleman, S. (2016). Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587), 484-489.



## Dermatologist-level classification of skin cancer with deep neural networks

Andre Esteva, Brett Kuprel, Roberto A. Novoa, Justin Ko, Susan M. Swetter, Helen M. Blau & Sebastian Thrun

Affiliations | Contributions | Corresponding authors

Nature 542, 115–118 (02 February 2017) | doi:10.1038/nature21056

Received 28 June 2016 | Accepted 14 December 2016 | Published online 25 January 2017

LA NACION



en/Flickr bajo licencia CC 2.0

## Un nuevo algoritmo es capaz de diagnosticar cáncer de piel con la misma eficacia que un dermatólogo

Hicieron que un sistema de inteligencia artificial analizara miles de imágenes para aprender a reconocer melanomas; por ahora es un prototipo

MIÉRCOLES 25 DE ENERO DE 2017 • 18:04

Un nuevo algoritmo desarrollado por la universidad de Stanford es capaz de diagnosticar un cáncer de piel con la misma precisión que un médico humano. No es el

## Attention Is All You Need

Ashish Vaswani\*  
Google Brain  
avaswani@google.com

Noam Shazeer\*  
Google Brain  
noam@google.com

Niki Parmar\*  
Google Research  
nikip@google.com

Jakob Uszkoreit\*  
Google Research  
usz@google.com

Llion Jones\*  
Google Research  
llion@google.com

Aidan N. Gomez\* †  
University of Toronto  
aidan@cs.toronto.edu

Lukasz Kaiser\*  
Google Brain  
lukasz.kaiser@google.com

Illia Polosukhin\* ‡  
illia.pолосухин@gmail.com

### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.





A.M.  
**TURING**  
AWARD  
2018

YOSHUA BENGIO,  
GEOFFREY E. HINTON  
AND YANN LECUN

For conceptual and engineering breakthroughs that have made deep neural networks a critical component of computing



# Grandes áreas

Deep  
Learning



Datos  
Secuenciales

Imágenes

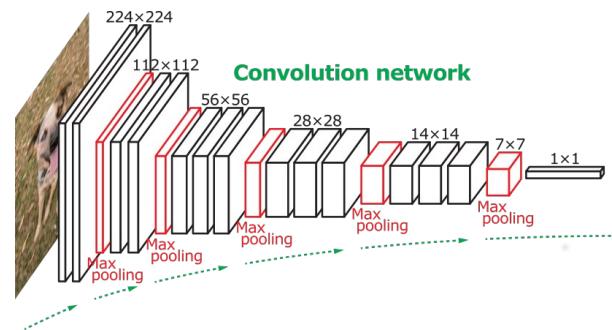
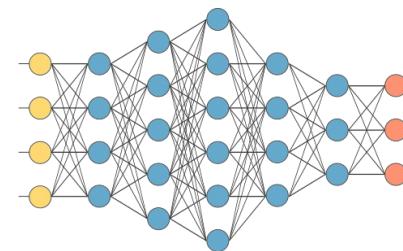
# REDES NEURONALES



## IMÁGENES



# REDES NEURONALES CONVOLUCIONALES



# Algo de historia



# Primeras redes neuronales

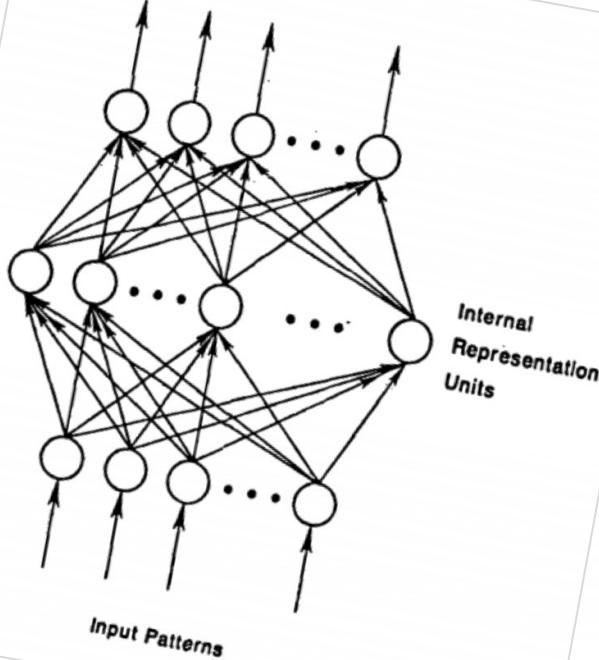
A 164 453

## LEARNING INTERNAL REPRESENTATIONS BY ERROR PROPAGATION

David E. Rumelhart, Geoffrey E. Hinton,  
and Ronald J. Williams

September 1985

Report 8506



S F

COGNITIVE  
SCIENCE



THE 80s

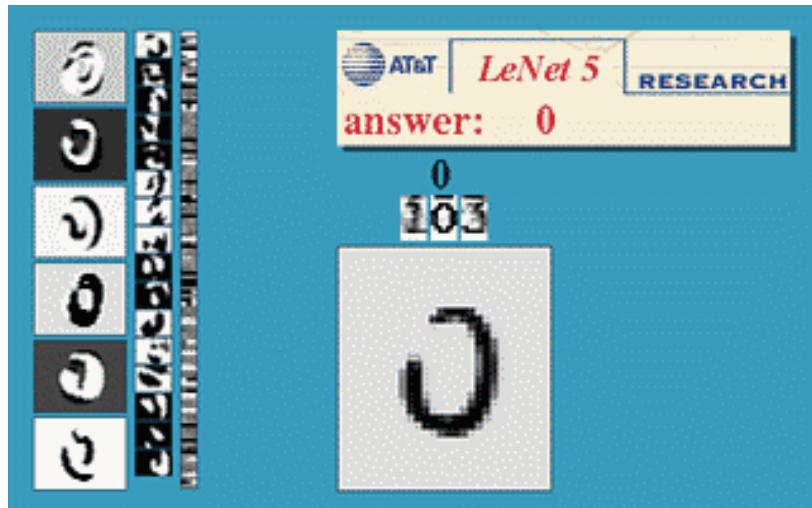


Geoffrey E.  
Hinton

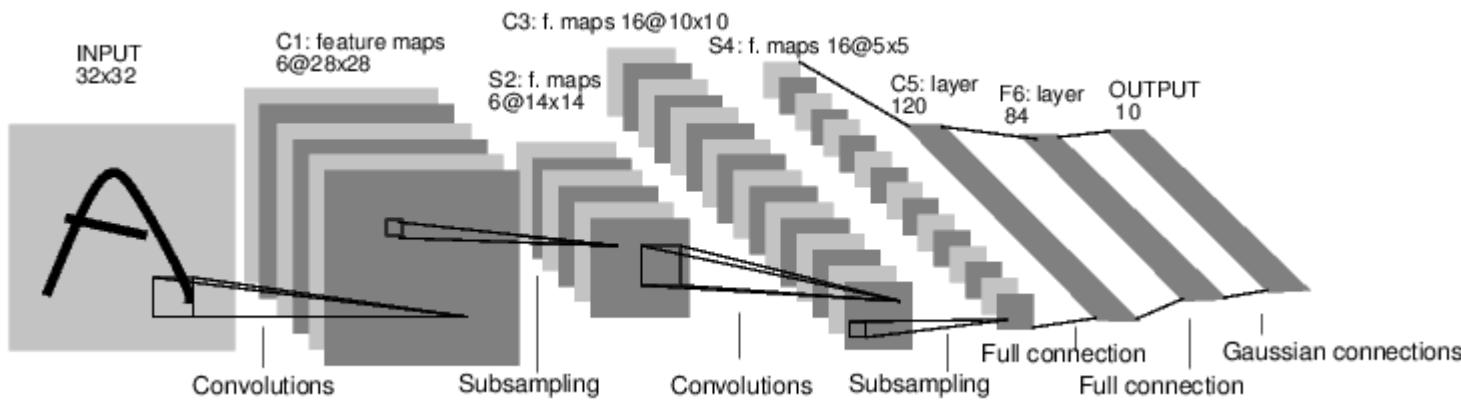


DTIC  
SELECTED  
FEB 2 0 1986  
S D

# Primera red convolucional



the 90s



Yann LeCun



# Lectura de código postal

60626

86 – (60626)

97222

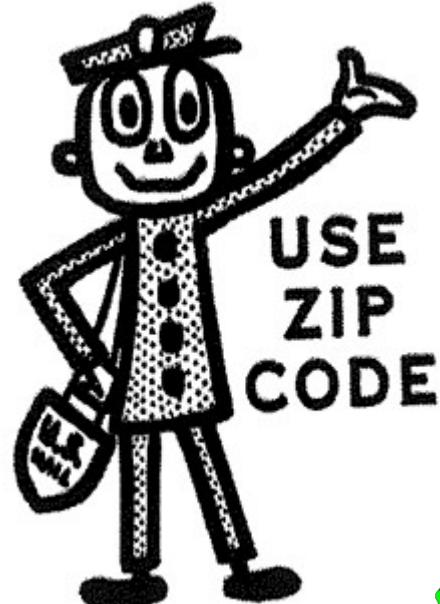
92 – (97222)

82071

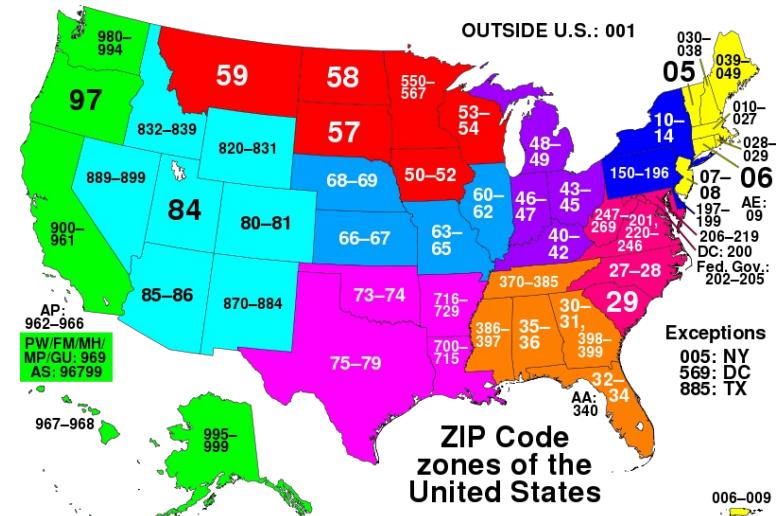
119 – (82071)

19801

469 – (19801)

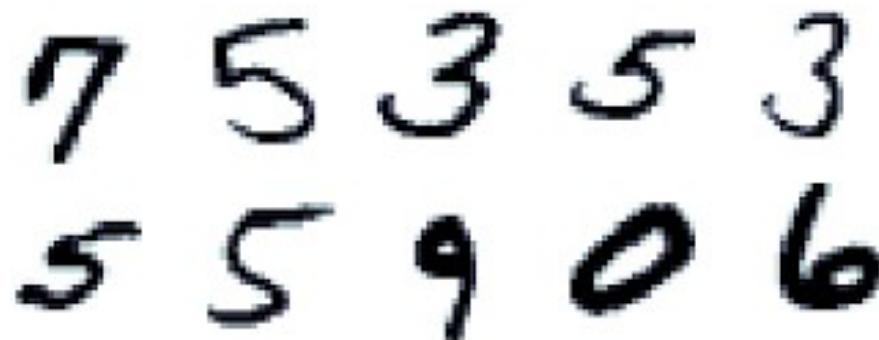


the 90s



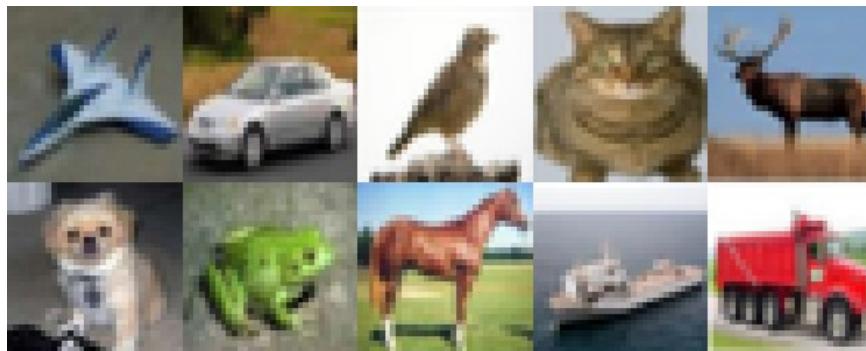
# El problema de la resolución

**28x28**



**MNIST**

**32x32**

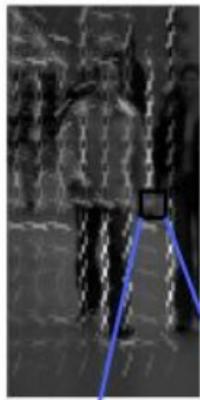


**CIFAR-10**

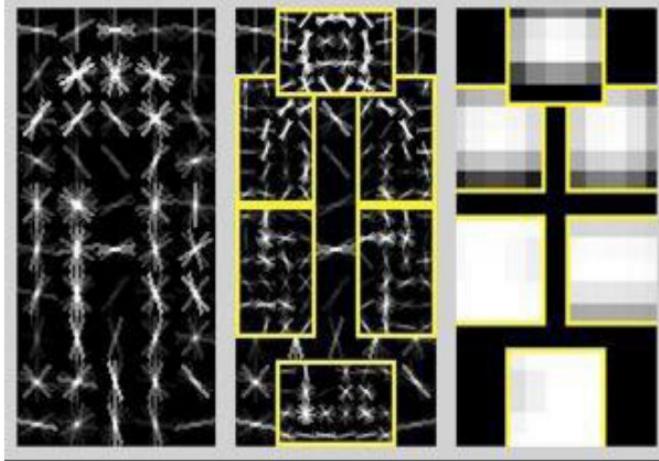
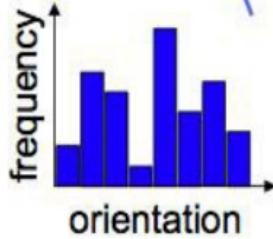
**96x96**



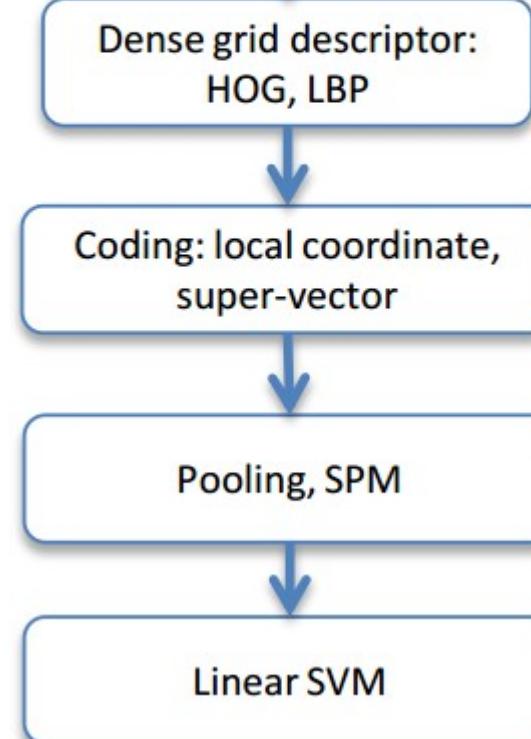
**STL-10**



Histogram of  
Gradients  
(HoG)  
Dalal & Triggs,  
2005



Deformable Part Model Felzenswalb,  
McAllester, Ramanan, 2009



2010



# Large Scale Visual Recognition Challenge

Steel drum:

The Image Classification Challenge:  
1,000 object classes  
1,431,167 images



**Output:**  
Scale  
T-shirt  
Steel drum  
Drumstick  
Mud turtle

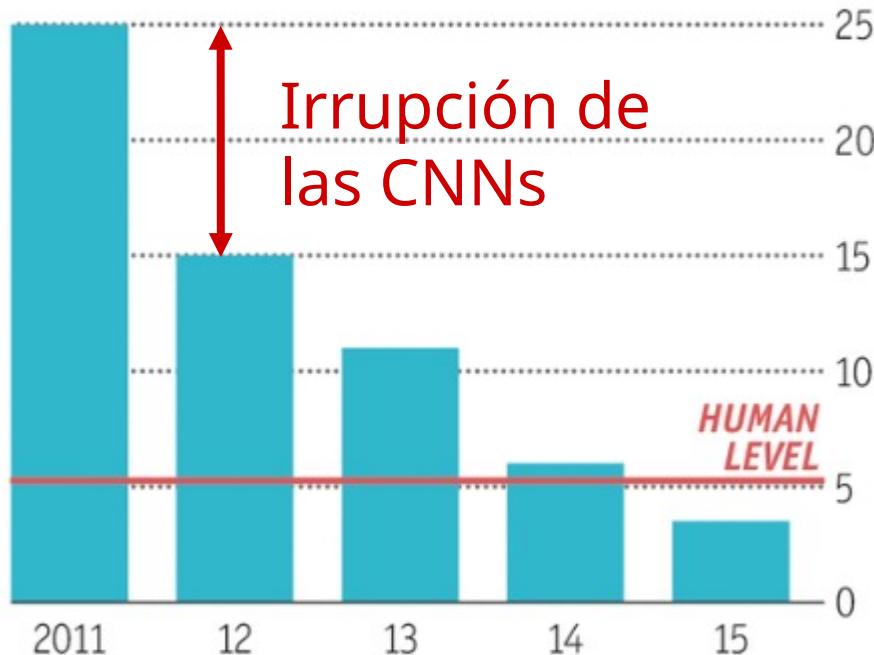


**Output:**  
Scale  
T-shirt  
Giant panda  
Drumstick  
Mud turtle



# Las redes convolucionales ganan la carrera

Error rates on ImageNet Visual Recognition Challenge, %



Sources: ImageNet; Stanford Vision Lab

## ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky  
University of Toronto  
kriz@cs.utoronto.ca

Ilya Sutskever  
University of Toronto  
ilya@cs.utoronto.ca

Geoffrey E. Hinton  
University of Toronto  
hinton@cs.utoronto.ca

Yann LeCun



Geoffrey E.  
Hinton



# Las redes convolucionales ganan la carrera

Year 2010

NEC-UIUC



Dense grid descriptor:  
HOG, LBP

Coding: local coordinate,  
super-vector

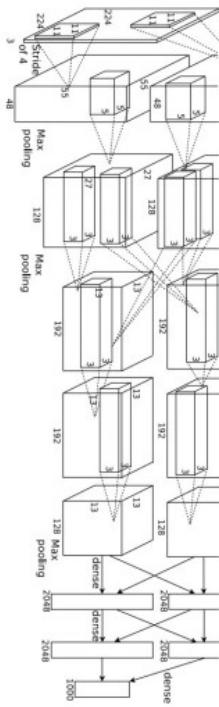
Pooling, SPM

Linear SVM

[Lin CVPR 2011]

Year 2012

SuperVision



[Krizhevsky NIPS 2012]

Year 2014

GoogLeNet

VGG

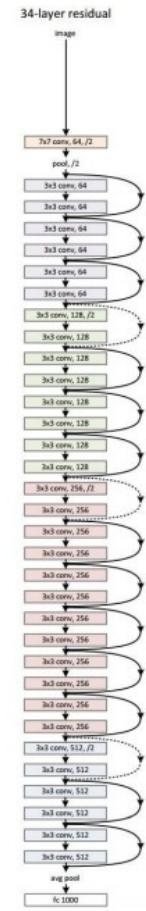


image  
conv-64  
conv-64  
maxpool  
conv-128  
conv-128  
maxpool  
conv-256  
conv-256  
maxpool  
conv-512  
conv-512  
maxpool  
conv-512  
conv-512  
maxpool  
FC-4096  
FC-4096  
FC-1000  
softmax

[Szegedy arxiv 2014] [Simonyan arxiv 2014]

Year 2015

MSRA



# “Revolution of Depth”

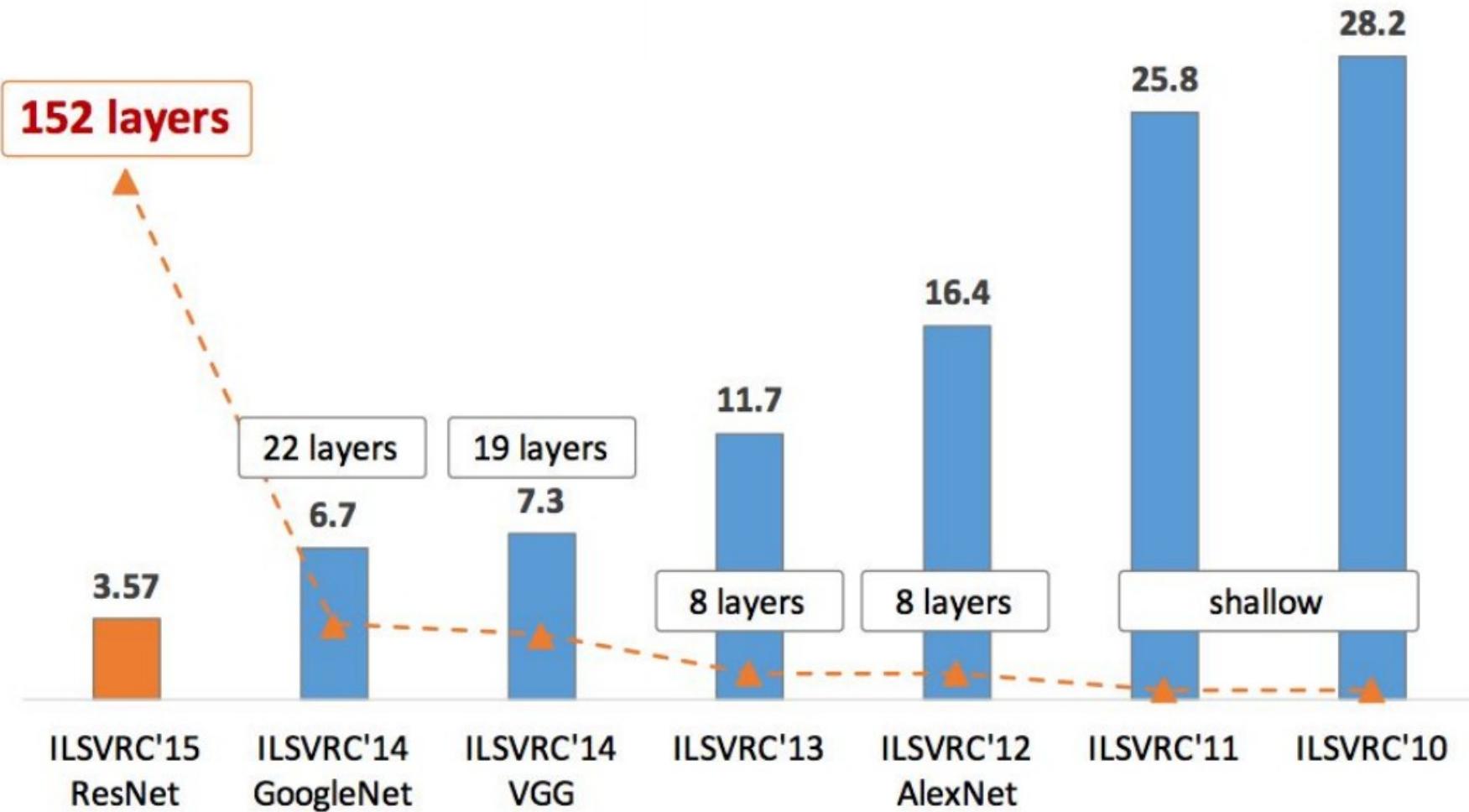


Figure copyright Kaiming He, 2016.

# Redes convolucionales

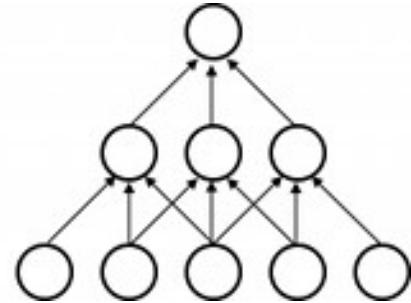
# Convolutional Neural Networks

Sparse connectivity

layer  $m+1$

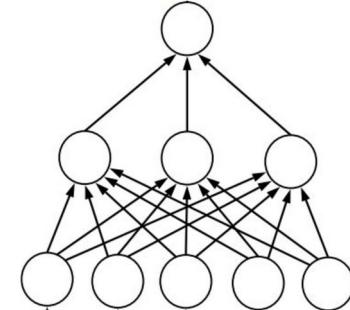
layer  $m$

layer  $m-1$

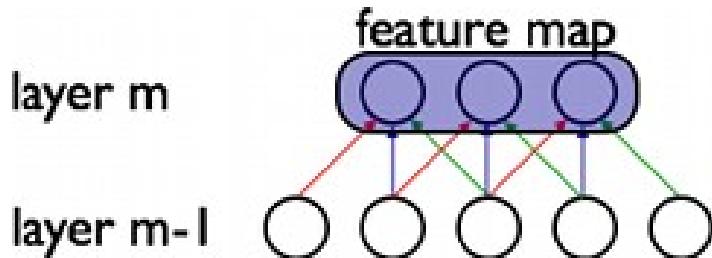


Dense connectivity

versus



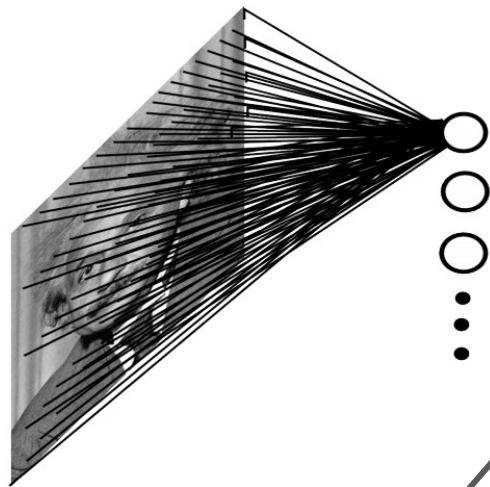
Shared weights



# Convolutional Neural Networks

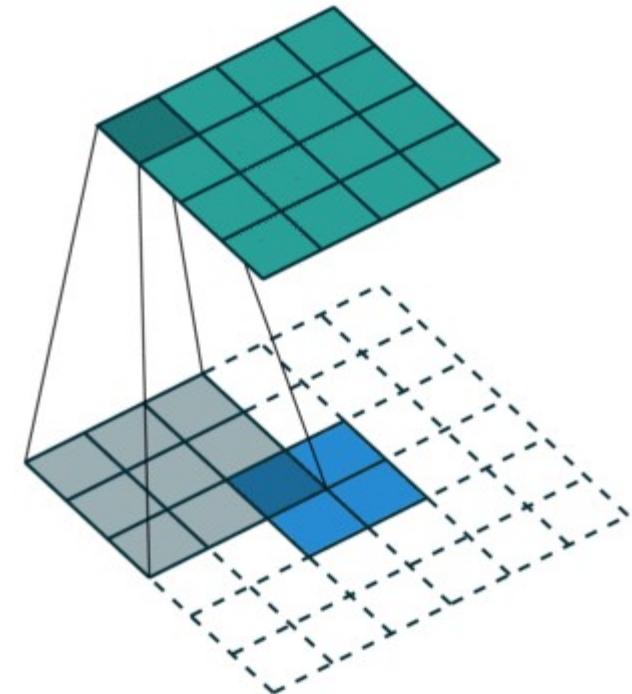
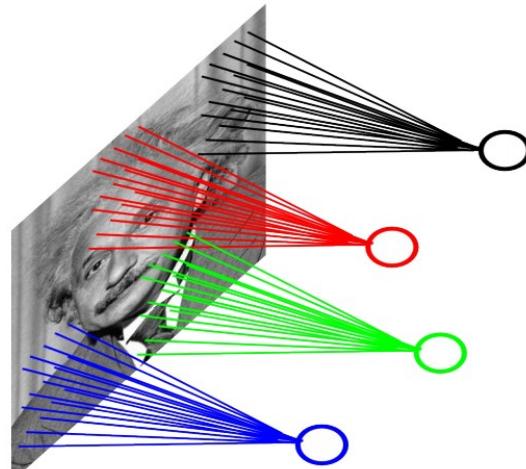
dot product + bias

$$h_k = \text{ReLU}(W_k \cdot x + b_k)$$



convolution + bias

$$h_{ij}^k = \text{ReLU}((W_k * x)_{ij} + b_k)$$



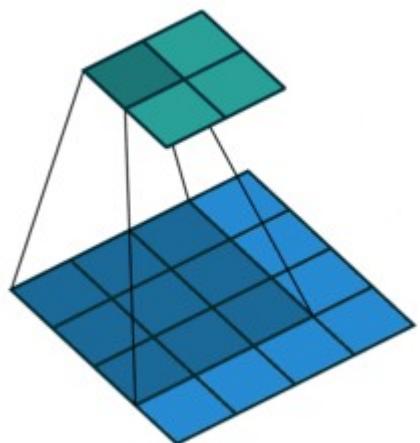
# Convolutional Neural Networks

## Por qué?

- **Invariancia translacional.** La visión responde igual en cualquier lugar, una cara es una cara en cualquier lugar de la imagen
- **Localidad.** Entender qué hay en una imagen depende del lugar en que se concentra la atención.
- **Reducir sobreajuste.** Al tener menos pesos ajustables, se reduce la posibilidad de aprender particularidades.

# Convolutional Neural Networks

## Ejemplo

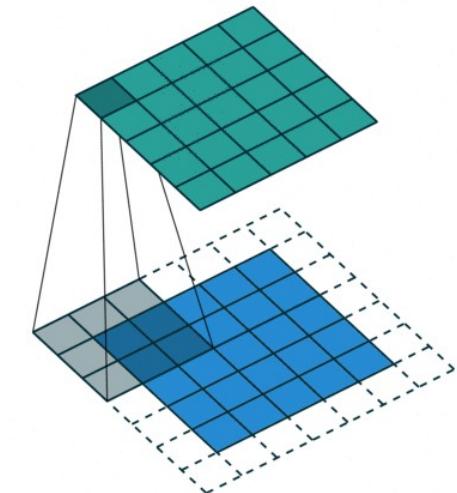


Input	Kernel	Output
$\begin{matrix} 0 & 1 & 2 \\ 3 & 4 & 5 \\ 6 & 7 & 8 \end{matrix}$	$*$	$=$
	$\begin{matrix} 0 & 1 \\ 2 & 3 \end{matrix}$	$\begin{matrix} 19 & 25 \\ 37 & 43 \end{matrix}$

# Padding y Stride

Input                    Kernel                    Output

$$\begin{array}{|c|c|c|c|c|}\hline 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 1 & 2 & 0 \\ \hline 0 & 3 & 4 & 5 & 0 \\ \hline 0 & 6 & 7 & 8 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 \\ \hline\end{array} * \begin{array}{|c|c|}\hline 0 & 1 \\ \hline 2 & 3 \\ \hline\end{array} = \begin{array}{|c|c|c|c|}\hline 0 & 3 & 8 & 4 \\ \hline 9 & 19 & 25 & 10 \\ \hline 21 & 37 & 43 & 16 \\ \hline 6 & 7 & 8 & 0 \\ \hline\end{array}$$



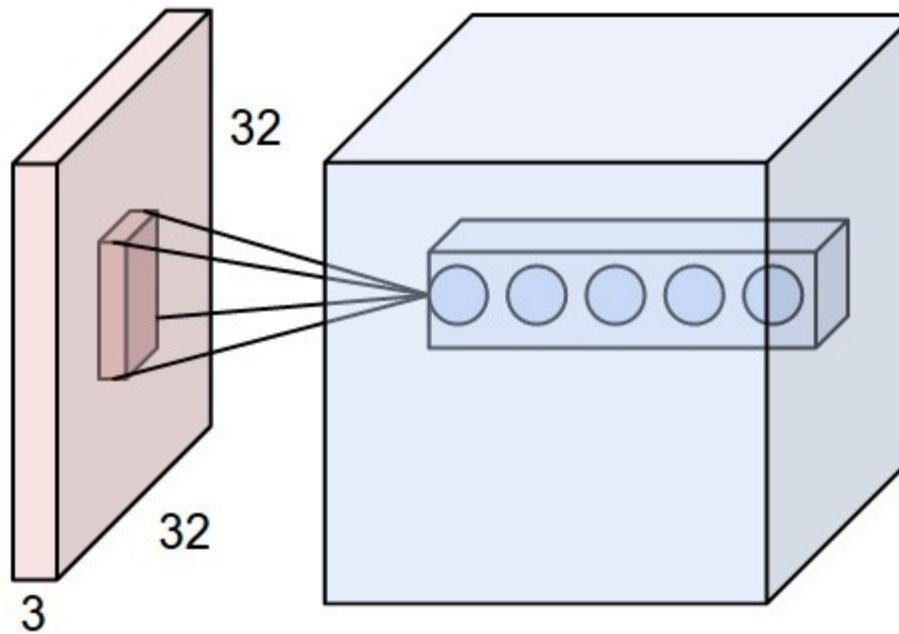
**Stride** is the number of “unit” the kernel shifted per slide over rows/columns.

E.g., Strides of 3 for height and 2 for width

Input                    Kernel                    Output

$$\begin{array}{|c|c|c|c|c|}\hline 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 1 & 2 & 0 \\ \hline 0 & 3 & 4 & 5 & 0 \\ \hline 0 & 6 & 7 & 8 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 \\ \hline\end{array} * \begin{array}{|c|c|}\hline 0 & 1 \\ \hline 2 & 3 \\ \hline\end{array} = \begin{array}{|c|c|}\hline 0 & 8 \\ \hline 6 & 8 \\ \hline\end{array}$$

# Una capa convolucional transforma el volumen



Entrada: imagen, 32x32, 3 canales RGB: 32x32x3

Filtro: tamaño 3x3, x 3 canales.

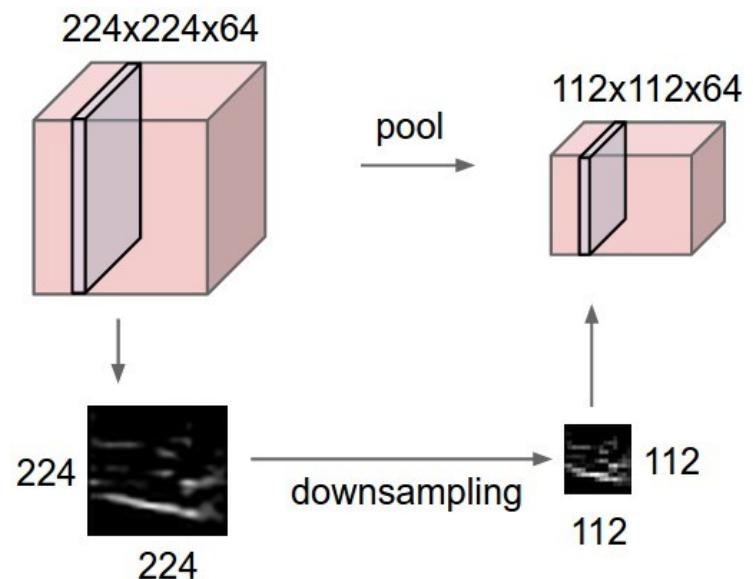
Salida: aplico 5 filtros convolucionales distintos,  
obtengo por cada uno una “imagen” de 32x32: 32x32x5

# Reducción de volumen: Pooling

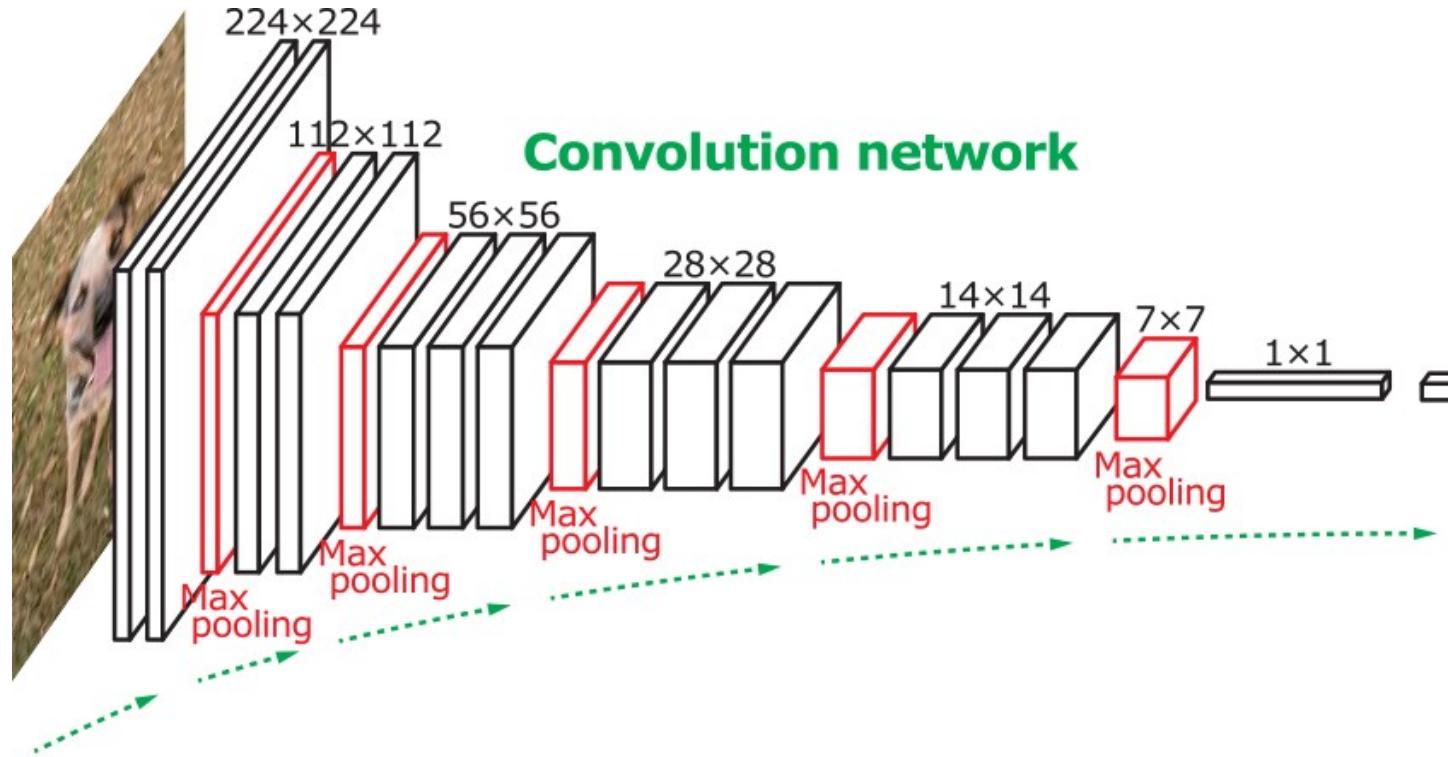
1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4



6	8
3	4



# Ejemplo: VGG net ("versión D")



dataset: ImageNet (ILSVRC 2013)

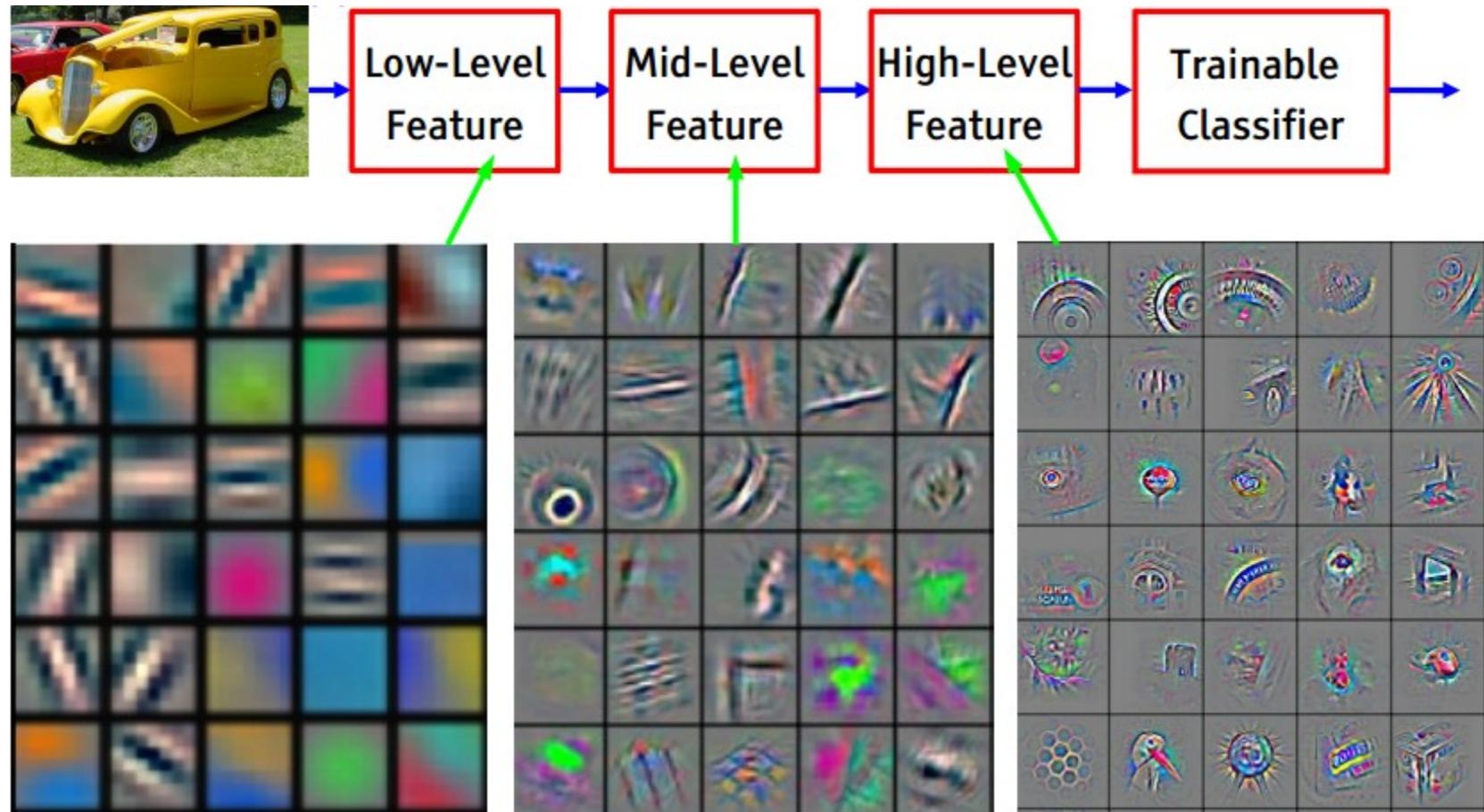
Input  $224 \times 224$   
conv3-64  
conv3-64  
**maxpool**  
conv3-128  
conv3-128  
**maxpool**  
conv3-256  
conv3-256  
conv3-256  
**maxpool**  
conv3-512  
conv3-512  
conv3-512  
**maxpool**  
conv3-512  
conv3-512  
conv3-512  
**maxpool**  
FC-4096  
FC-4096  
FC-1000  
softmax

# ¿Cómo se entrena una CNN?

## Backpropagation.

- Derivación automática.
- Learning rate adaptativos.
- Minibatch.
- Unidades Relu para mantener el gradiente.

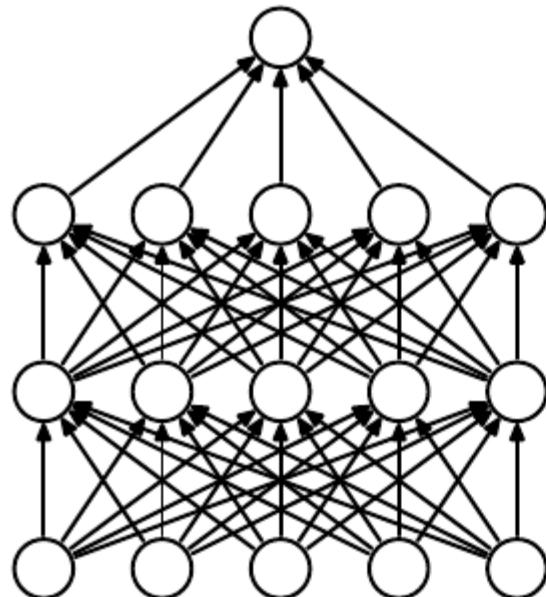
# ¿Qué ve una CNN ya entrenada?



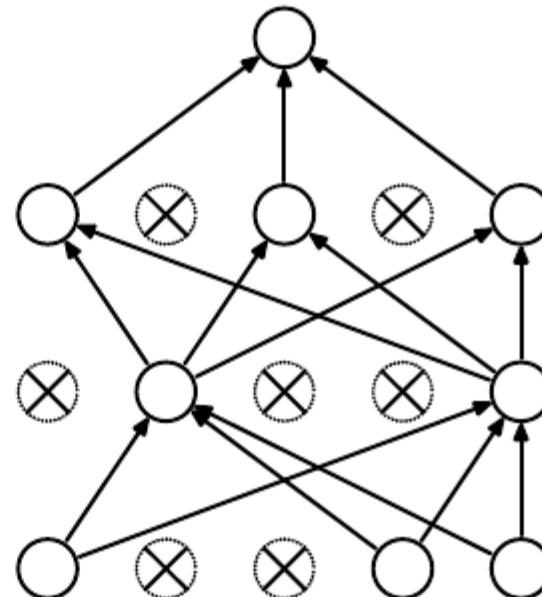
# Regularización

Dropout y Data augmentation

# Dropout Neural Net Model



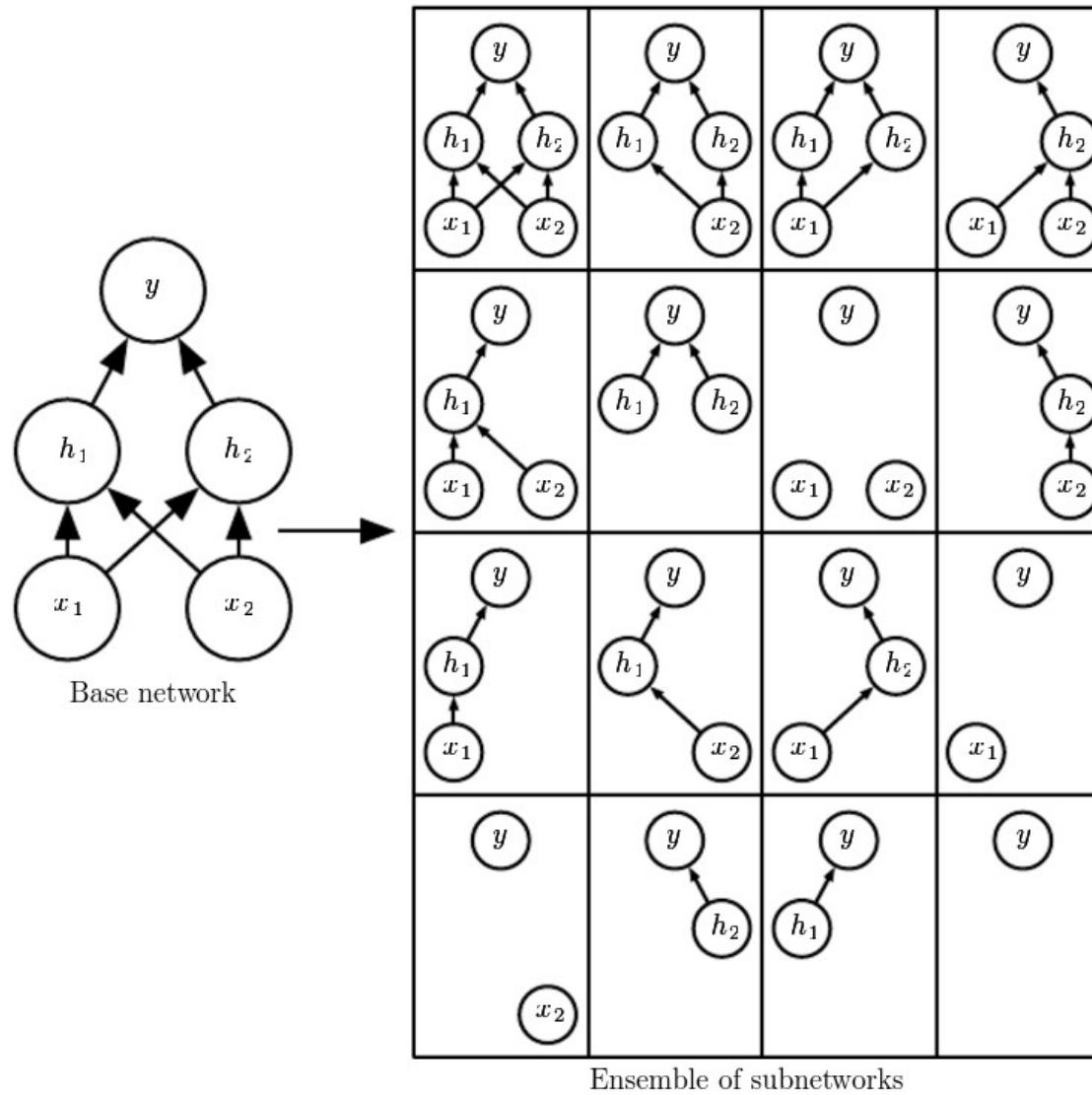
(a) Standard Neural Net



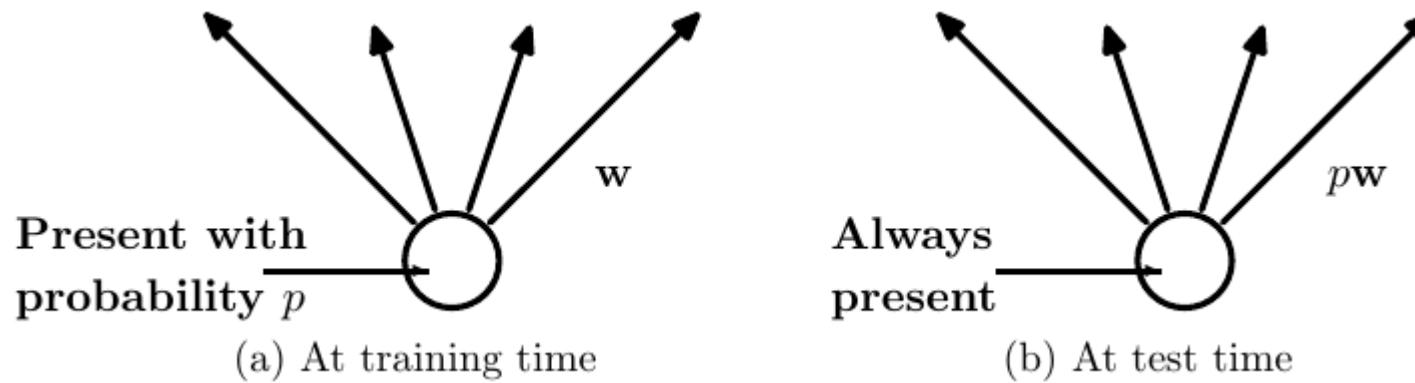
(b) After applying dropout.

Dropout Neural Net Model. **Left:** A standard neural net with 2 hidden layers. **Right:** An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.

# Dropout Neural Net Model

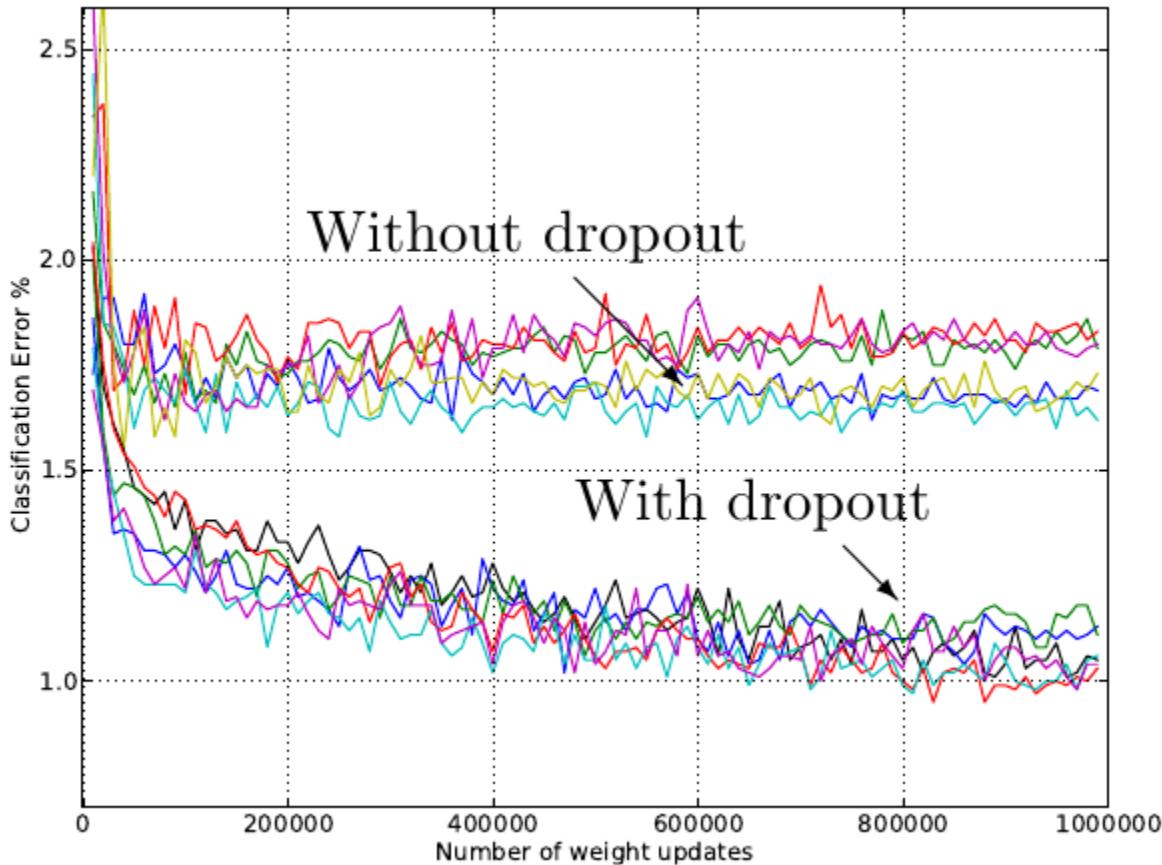


# Dropout units



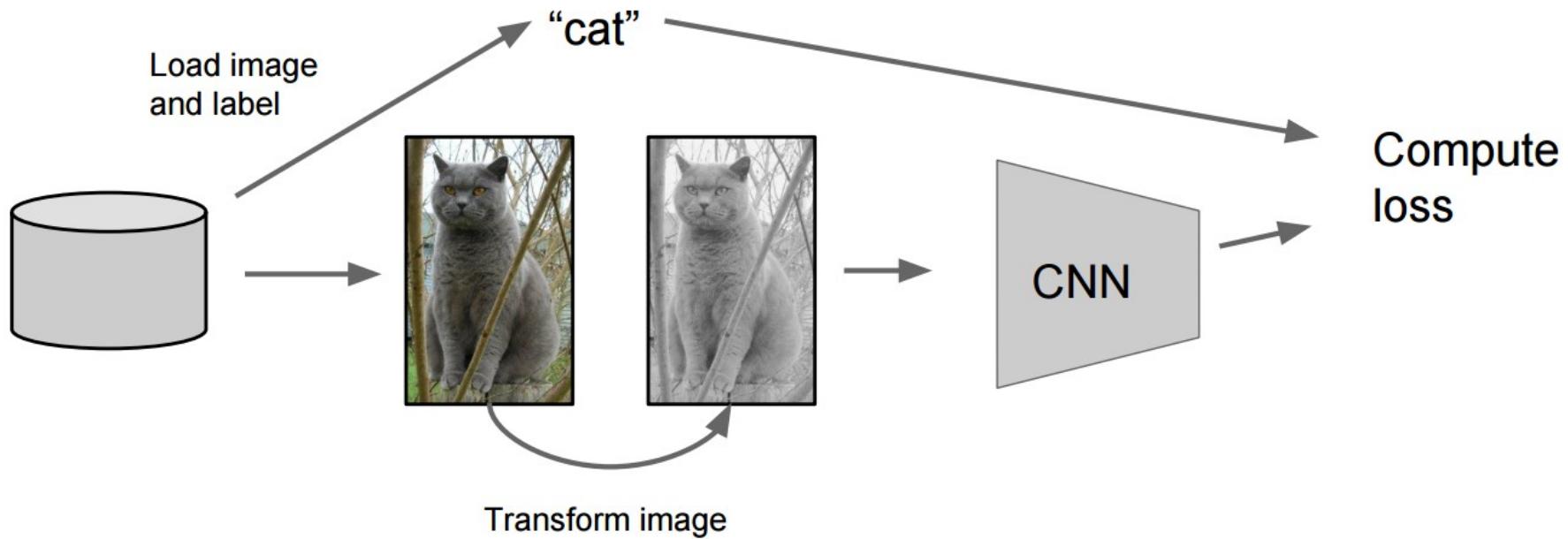
**Left:** A unit at training time that is present with probability  $p$  and is connected to units in the next layer with weights  $w$ . **Right:** At test time, the unit is always present and the weights are multiplied by  $p$ . The output at test time is same as the expected output at training time.

# Dropout: Robustness

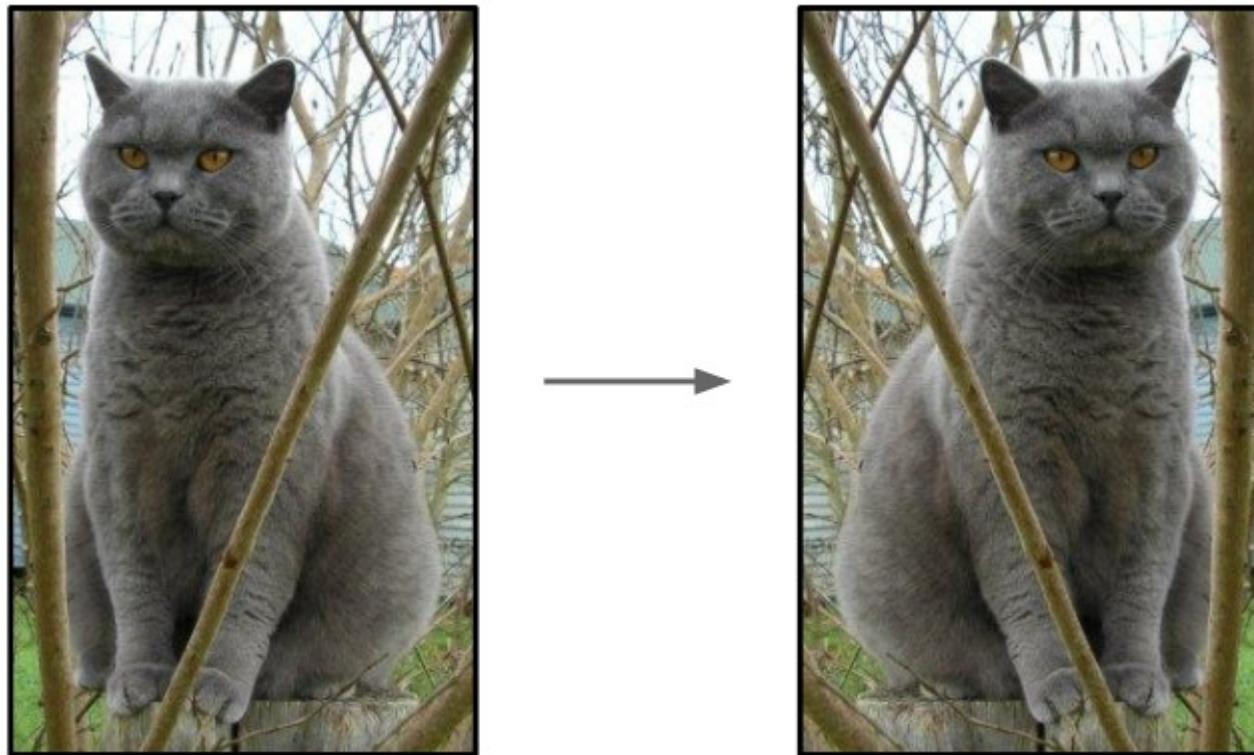


Test error for different architectures with and without dropout. The networks have 2 to 4 hidden layers each with 1024 to 2048 units.

# Data Augmentation



# Data Augmentation: Horizontal flips



## Data Augmentation: Random crops/scales

**Training:** sample random crops / scales

Por ejemplo [ResNet]:

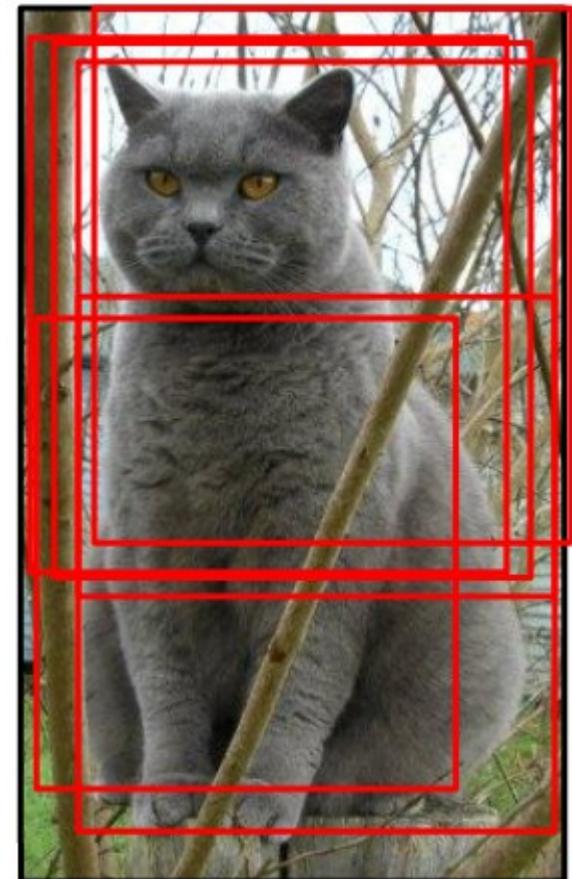
1. Pick random  $L$  in range  $[256, 480]$
2. Resize training image, short side =  $L$
3. Sample random  $224 \times 224$  patch

**Testing:** average a fixed set of crops

Por ejemplo [ResNet]

4. Resize image at 5 scales:  $\{224, 256, 384, 480, 640\}$
5. For each size, use 10  $224 \times 224$  crops: 4 corners + center, + flips

[ResNet] He, Kaiming, et al. "Deep residual learning for image recognition." arXiv preprint arXiv:1512.03385 (2015).



## Versión simple:

Alterar el contraste aleatoriamente



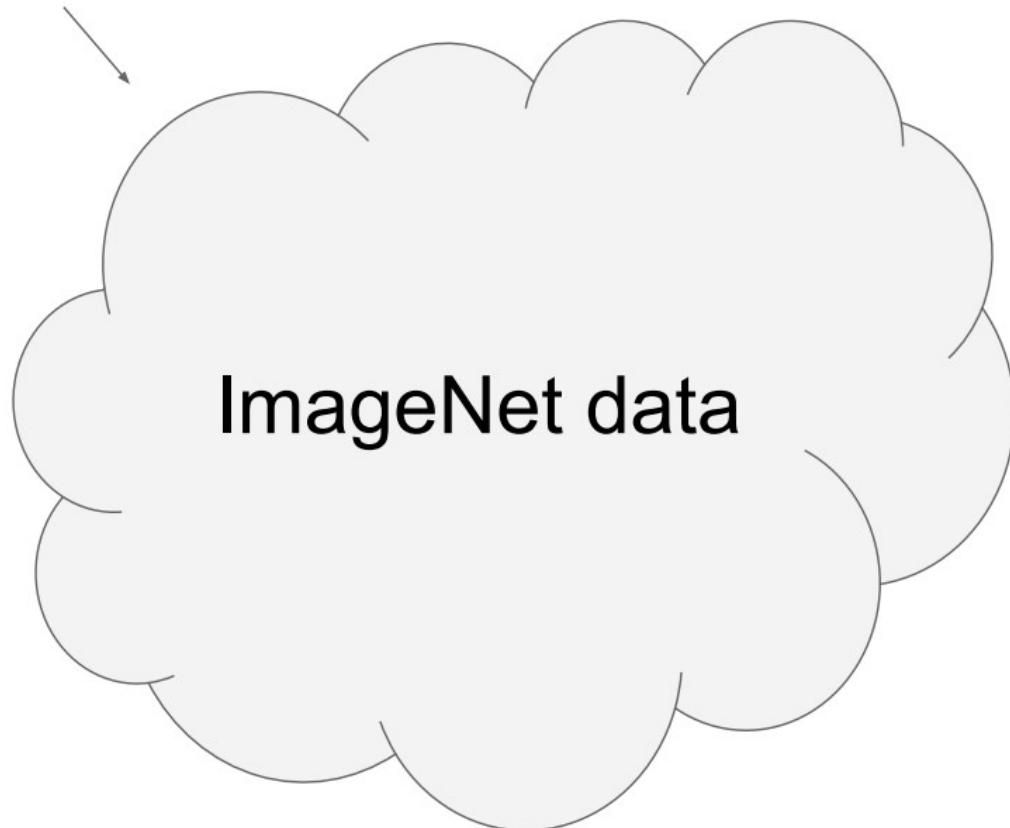
## Versión compleja:

1. Aplicar PCA a todos los píxeles [R,G,B] del training set
2. Samplear un "desplazamiento de color" a lo largo de las direcciones principales
3. Aplicarle este desplazamiento a todos los píxeles de una imagen de entrenamiento.

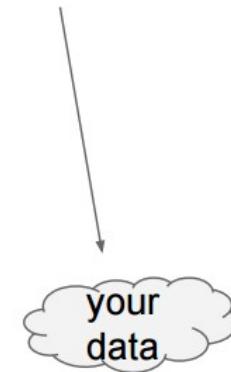
# Transfer Learning

# Pre entrenamiento sobre ImageNet

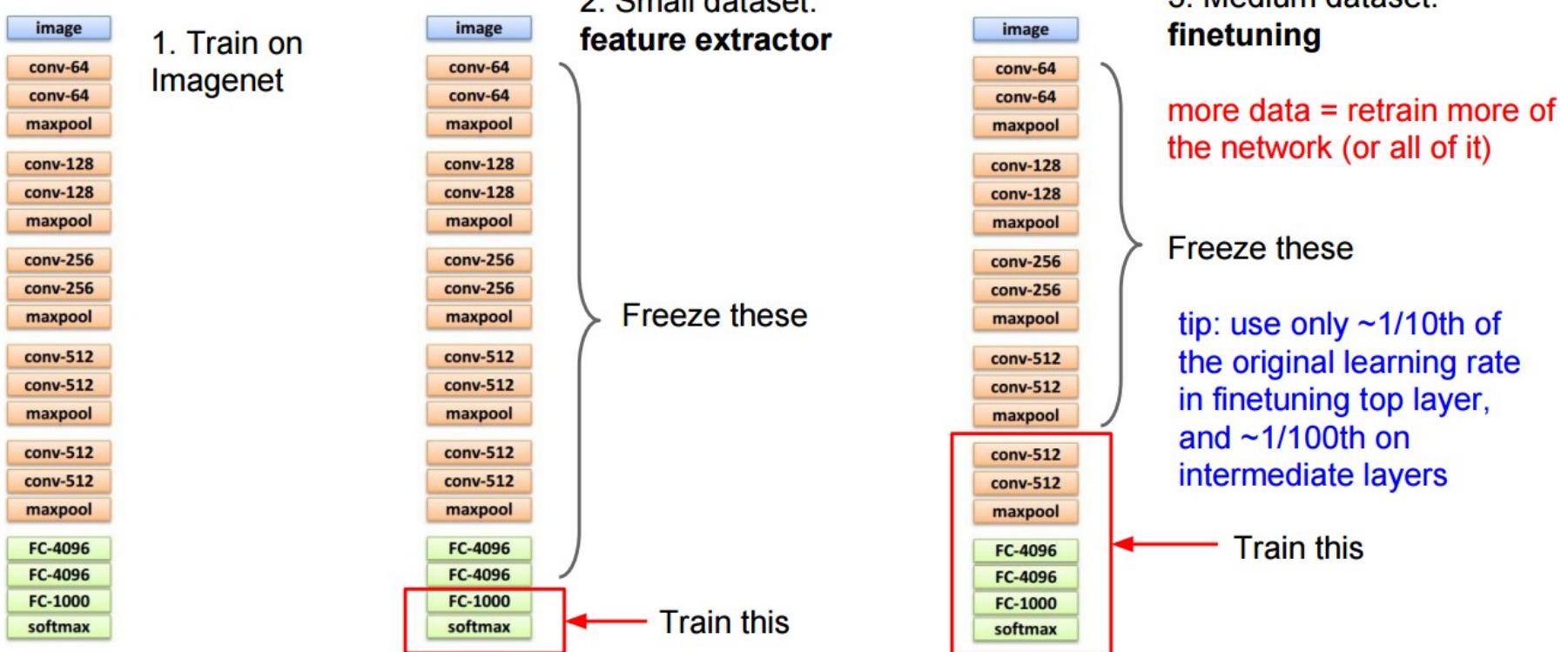
1. Train on ImageNet



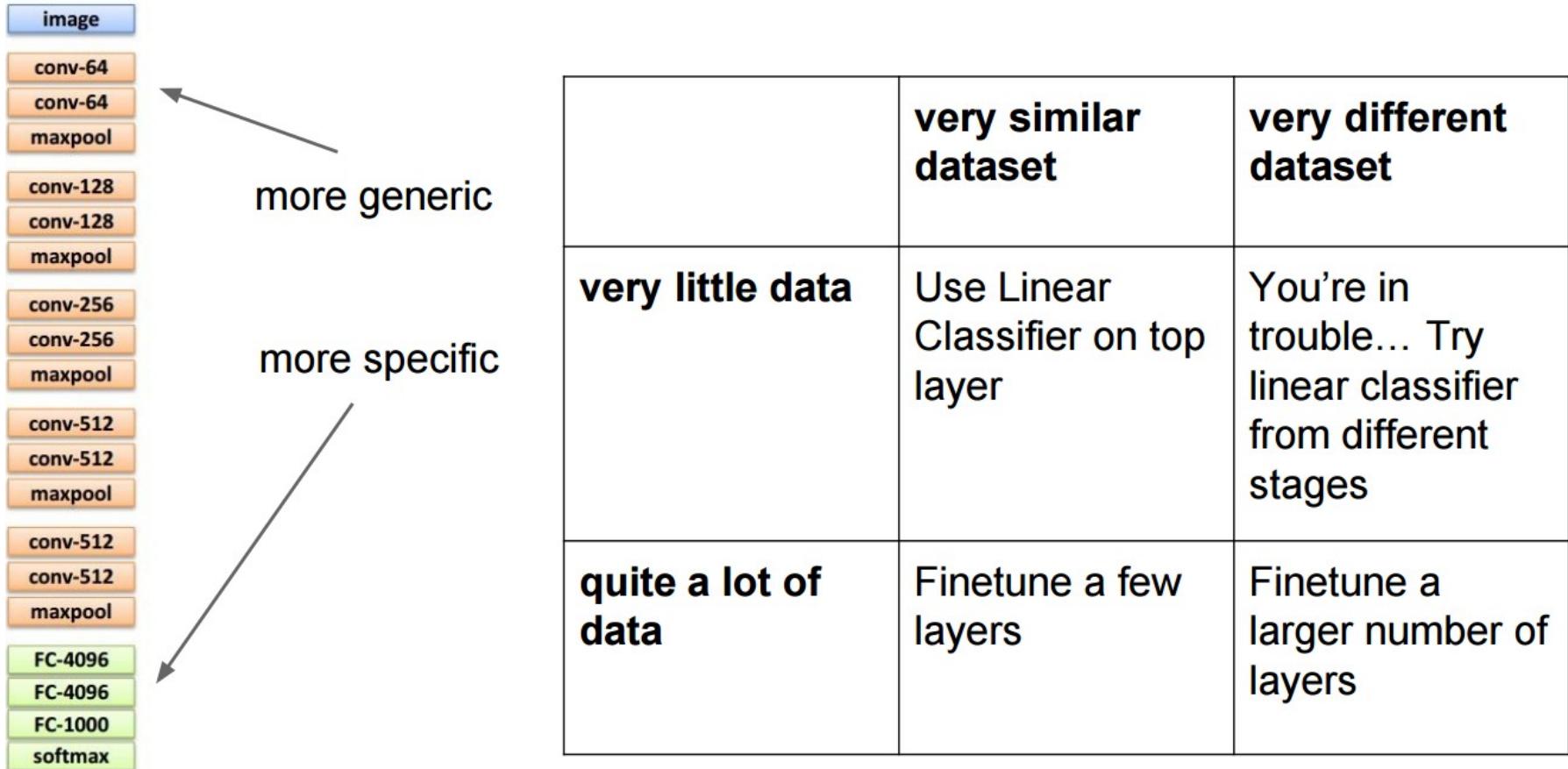
2. Finetune network on  
your own data



# Ejemplo con redes convolucionales



## Ejemplo con redes convolucionales



# Datos secuenciales

# Datos Secuenciales

time series

ia mirada panteísta donde un solo hombre inmortal es todos los hombres y a su vez ninguno. Y a partir de esta idea también puede irse, como luego veremos, que un solo texto también dos los textos. Según Borges este relato vendría a ser "el sueño de una ética para inmortales" y su tema "el efecto de la inmortalidad causaría en los hombres". Este efecto lo describe a través del autor implícito del relato, el anticuario José Artaphilus, quien narra la vida del tribuno romano Marcius Rufo. Así podremos presenciar en este relato la voz de un hombre que fue todos y a la vez fue nadie, ya que fueron "las palabras de otros [...] la pobre limosna que le dejaron las horas y siglos". El texto presente nos servirá para hacer una reflexión.

text



handwriting

Reçu De Monsieur D. Monceaux  
la somme De Nos cent francs, pour un  
timbre De la pension alimentaire pour  
le veuve André Monceaux.  
Le Dit Timbre sera remis le 1<sup>er</sup> Mai  
entrant et finissant le 31 juillet prochain

A Dijon le 1<sup>er</sup> Mai 1890

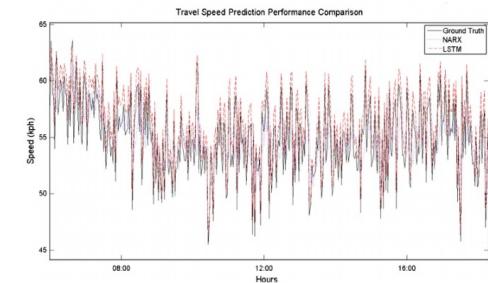
Reçu De Monsieur D. Monceaux  
la somme De Nos cent francs, pour un  
timbre De la pension alimentaire pour  
le veuve André Monceaux.  
Le Dit Timbre sera remis le 1<sup>er</sup> Mai  
entrant et finissant le 31 juillet prochain

speech

code

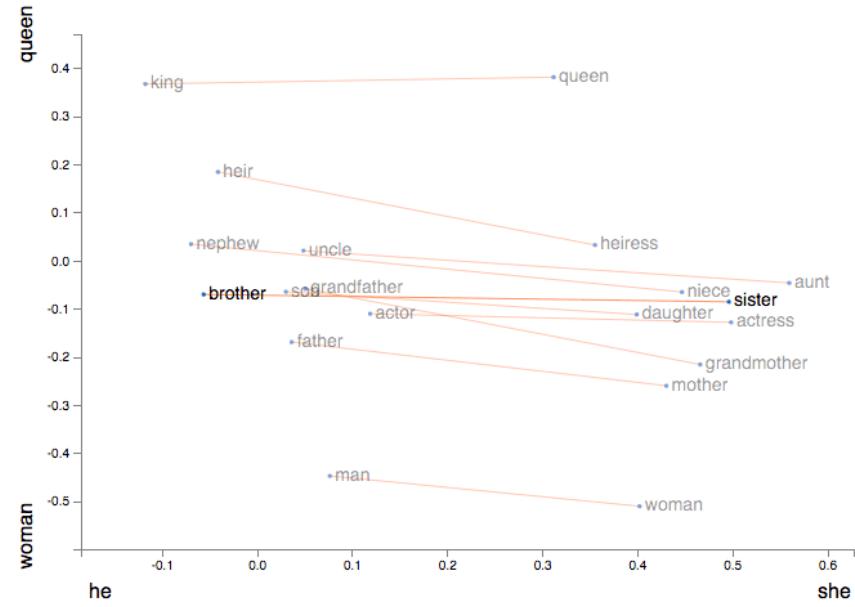
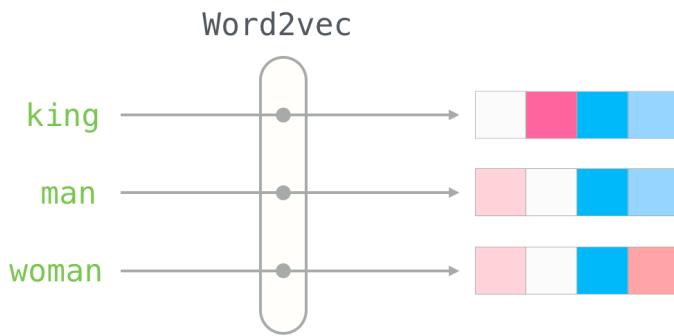
```
require_once('chorus/Shard.php');
Database::set_defaults();
array('user' => 'tumblr3', 'password' => 'm3MgHlCOKoh39AQD83TFhsBPLOMR',
      'database' => 'tumblr3', 'write_lock_tables' => '*',
      'extended_log' => (idate('G') == 17 && intval(idate('i')) == 56
if (__FILE__ == '/var/www/apps/tumblr/config/config.php' || __FILE__ == '/define('ENVIRONMENT', 'production');
if (!defined('DEFAULT_DATABASE')) define('DEFAULT_DATABASE', 'primary');
define('S3_BUCKET', 'data.tumblr.com');
define('ENABLE_PANTHER', true);
define('ENABLE_MEDIA_CDN', true);
define('ASSETS_URL', (ENABLE_MEDIA_CDN && !isset($_SERVER['HTTPS'])) &&
define('MEMCACHE_HOST', '10.252.0.68');
define('MEMCACHE_VERSION_HOST', '10.252.0.67');
define('VALIDATION_FAILURE_LOG', BASE_PATH . '/validate.log');
define('REDIRECT_403_LOG', BASE_PATH . '/403.log');
define('GOOGLE_API_KEY', isset($_SERVER['HTTP_HOST']) && $_SERVER['HTTP_HOST'] == 'ABOIAAAAJladOHJn-kbPSqUsrS6CyhTpoeXstiCMps15pU3slU-WDRPjxQts41ksQogksy
Database::add('primary', array('host' => '192.168.200.142'));
Database::add('db-tumblelogs', array('host' => '192.168.200.103'));
```

stock  
market



# Texto: Word embedding

El texto no se puede usar directamente como las imágenes. Se necesita un espacio vectorial donde palabras con sentido similar tengan una representación similar.

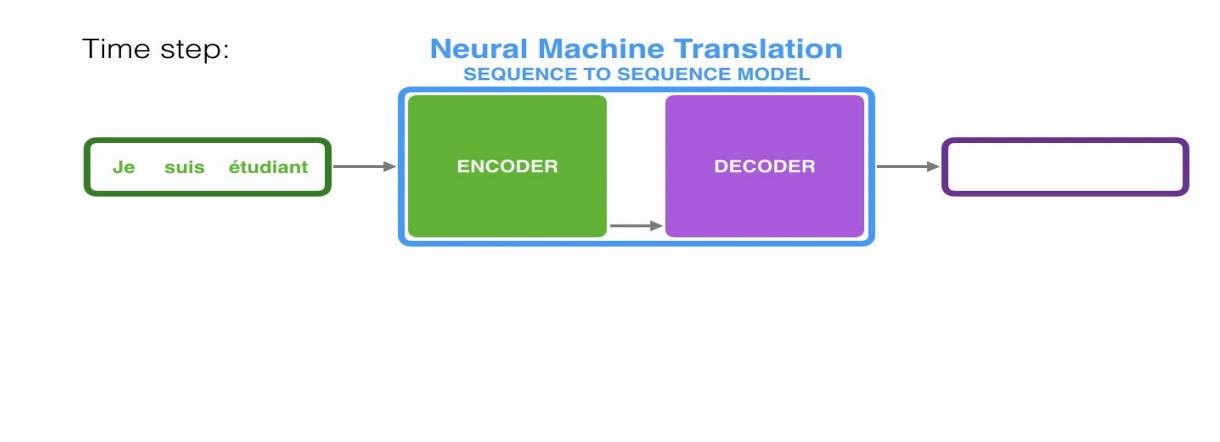


Learned: word2vec

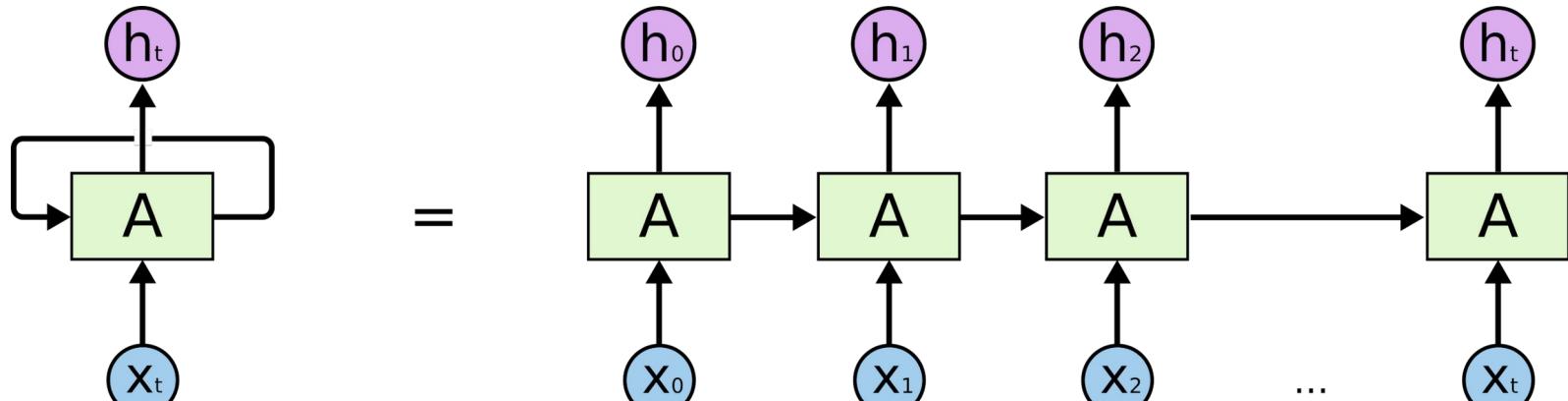
# Aprendiendo una secuencia

Para aprender una relación de secuencia-a-secuencia se usa un sistema en dos partes:

- Un encoder que traduce la secuencia a un estado interno que se aprende
- Un decoder que toma el estado interno y lo traduce a otra secuencia



# Recurrent Neural Networks

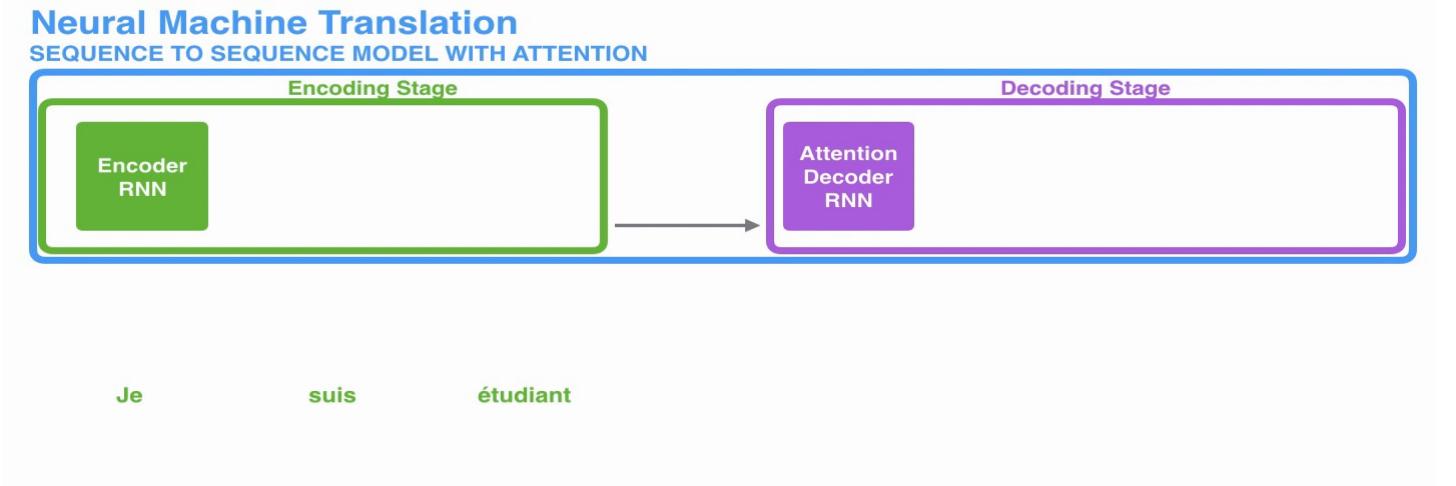


Neural Network  
with a loop

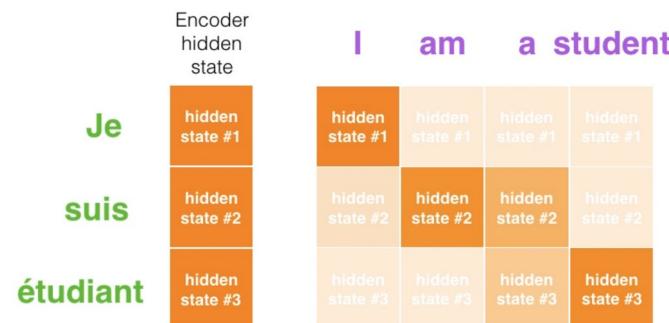
Unfolded Computational Graph

Problems: long-term dependencies + sequential learning

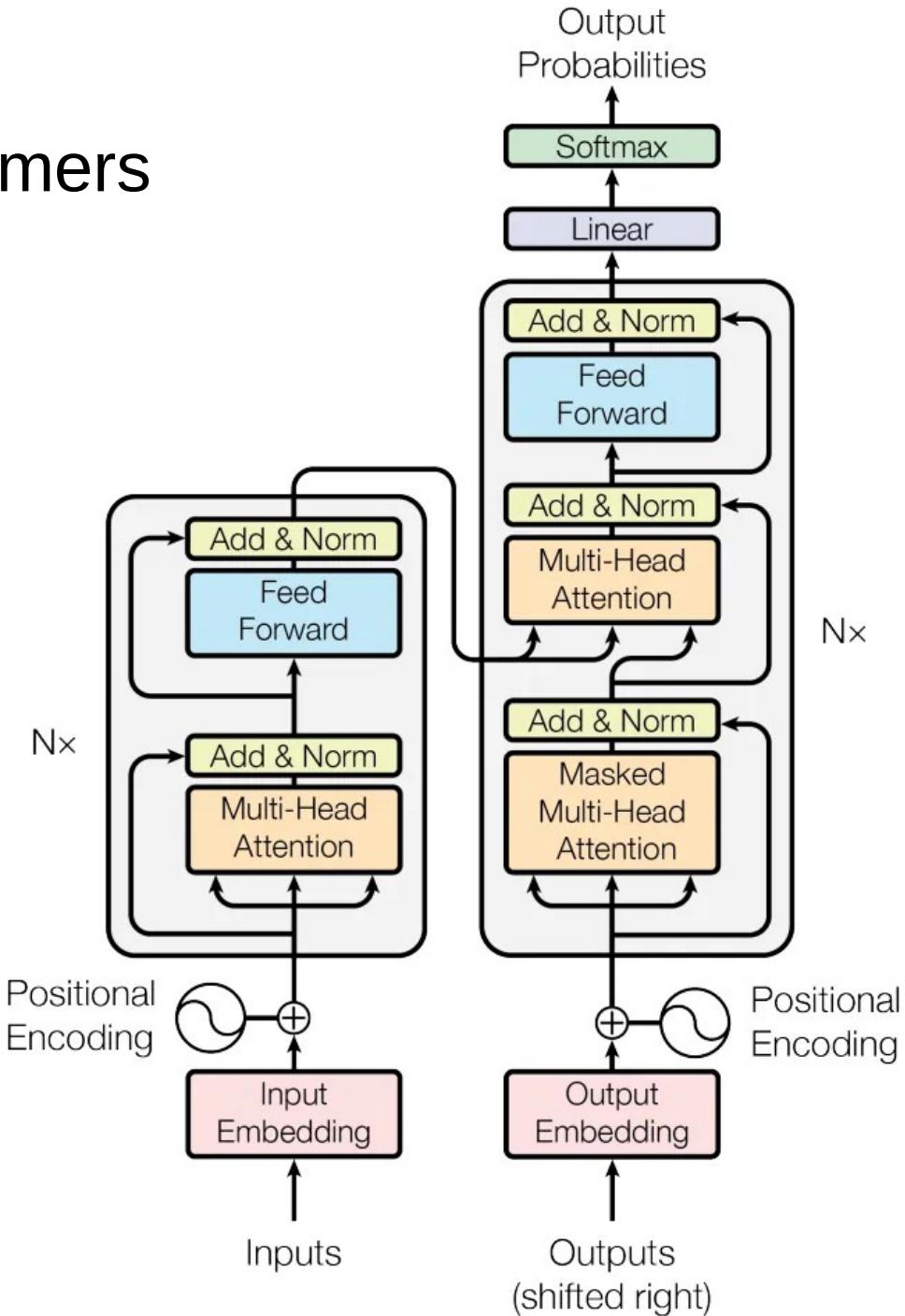
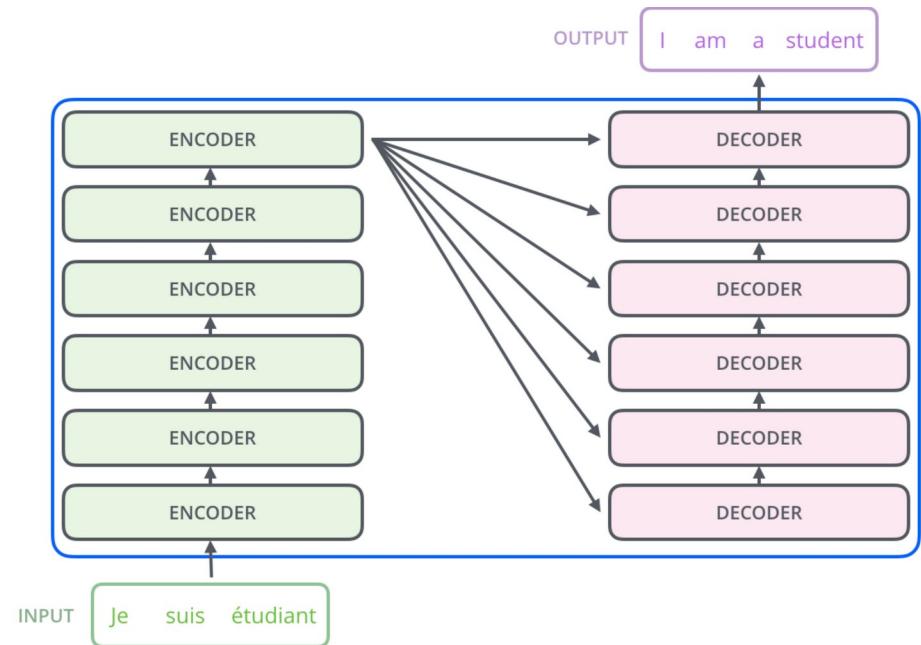
# Long term dependencies: Attention!



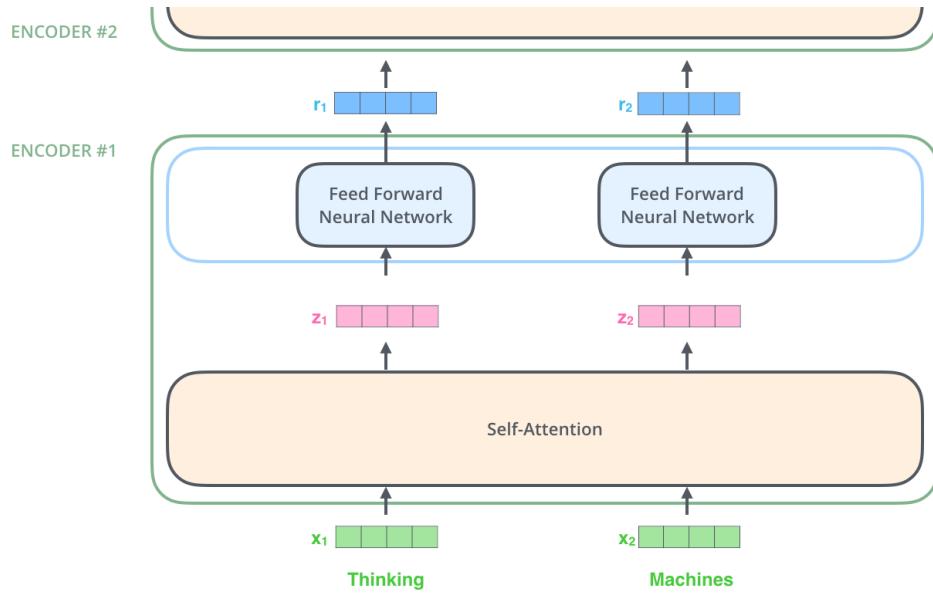
Permitimos al decoder ver todos los estados y elegir la información que usa en cada paso



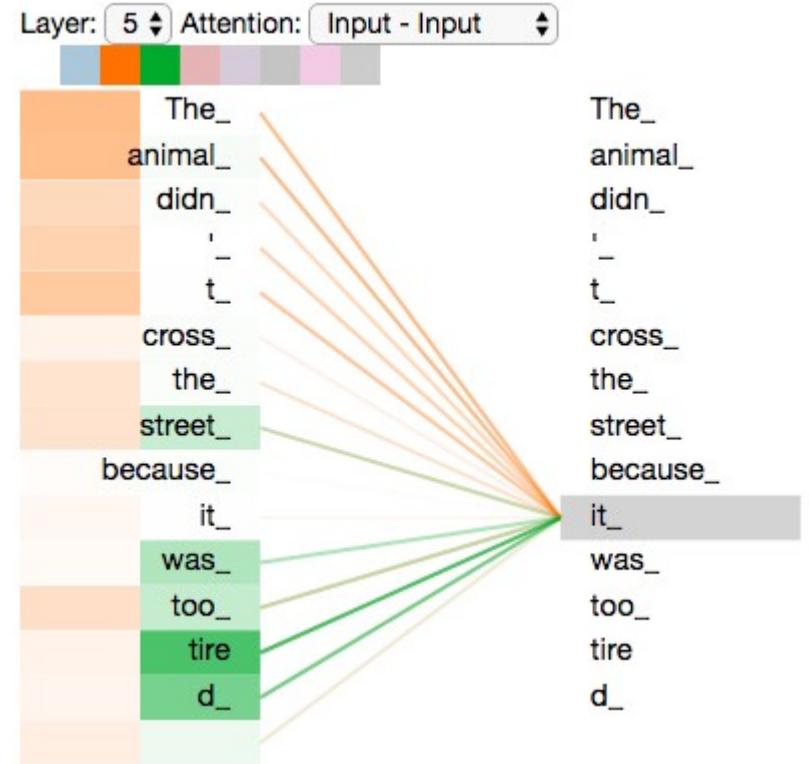
# Parallel learning: Transformers



# Transformers: key points



Multiple layers  
of encoders

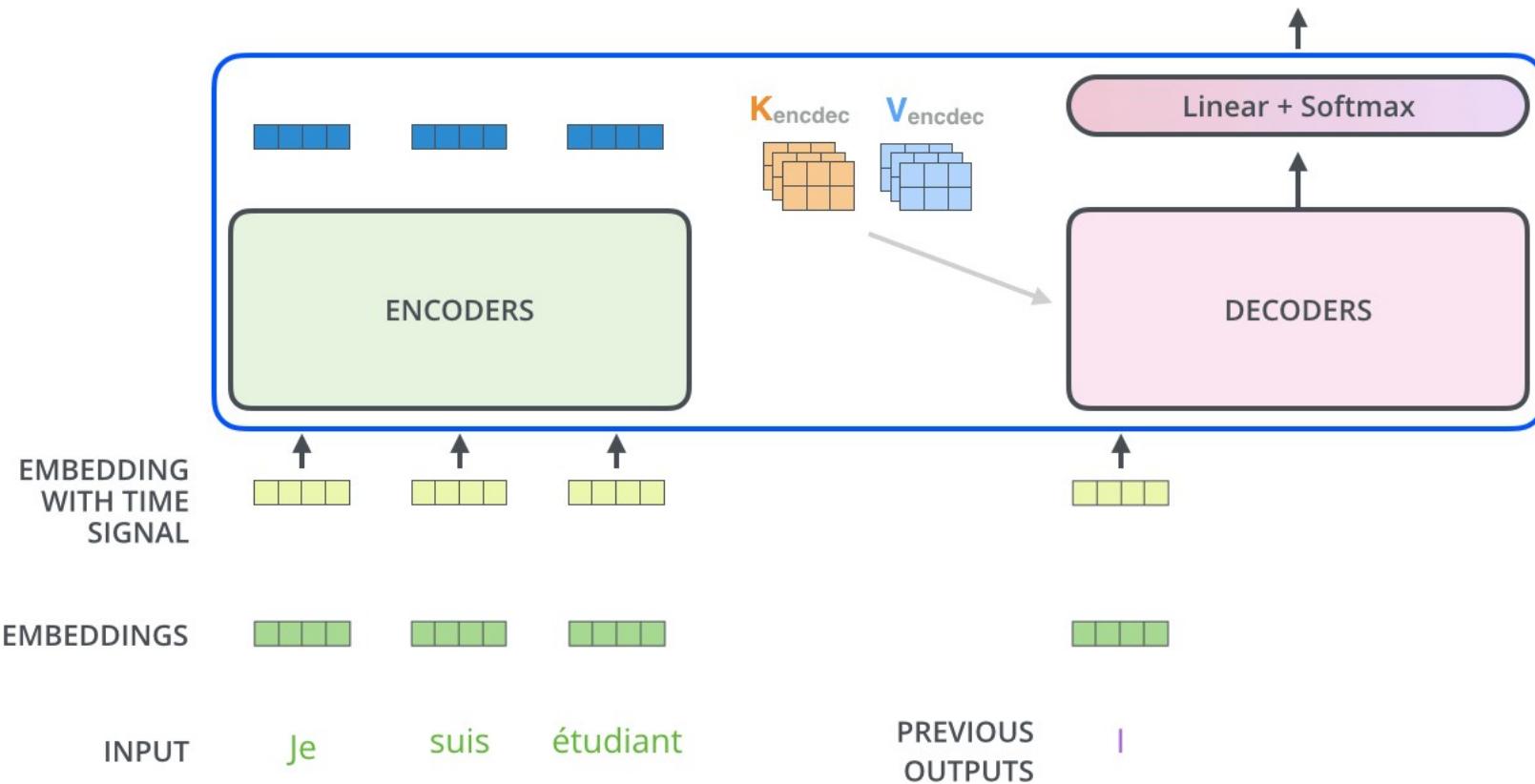


Multiple  
attention heads

# Transformers: decoding

Decoding time step: 1 2 3 4 5 6

OUTPUT |



# Transformers: images

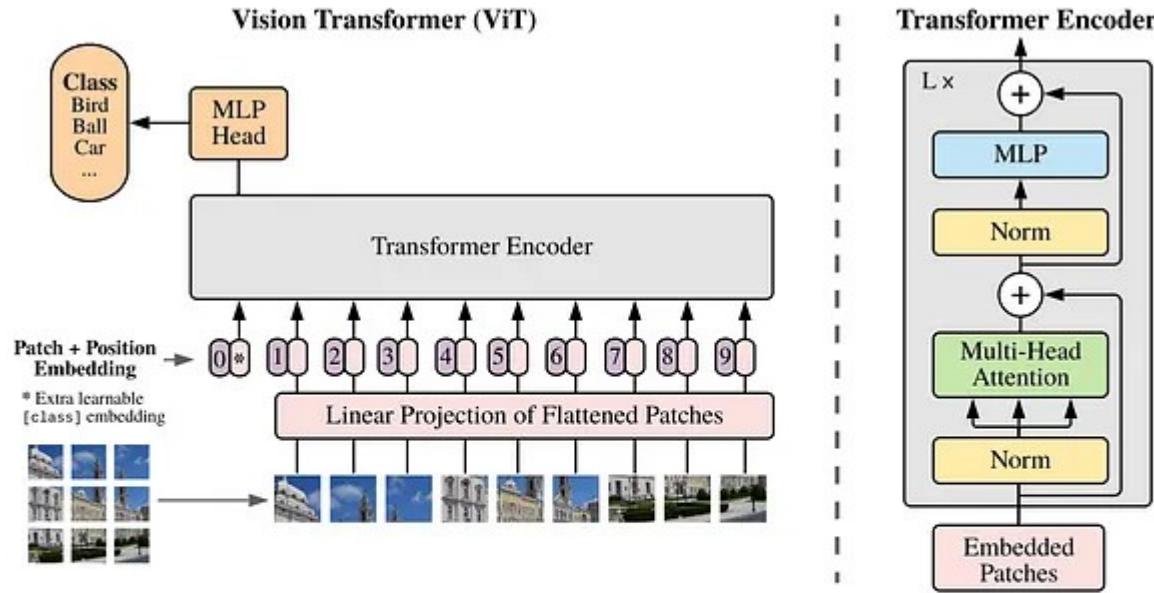


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable “classification token” to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).