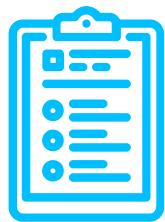


Sesión 1.1

Redes convolucionales

VGG, ResNet, MobileNet

1.

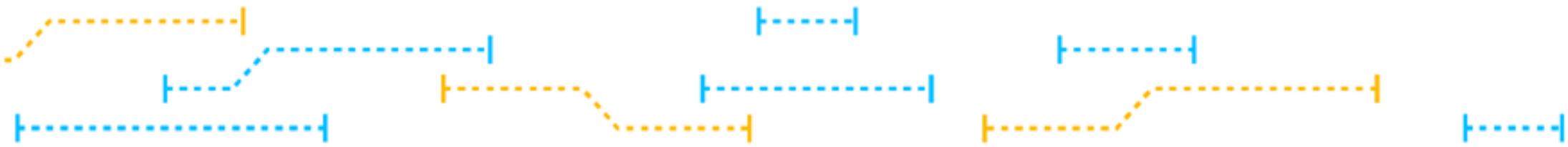
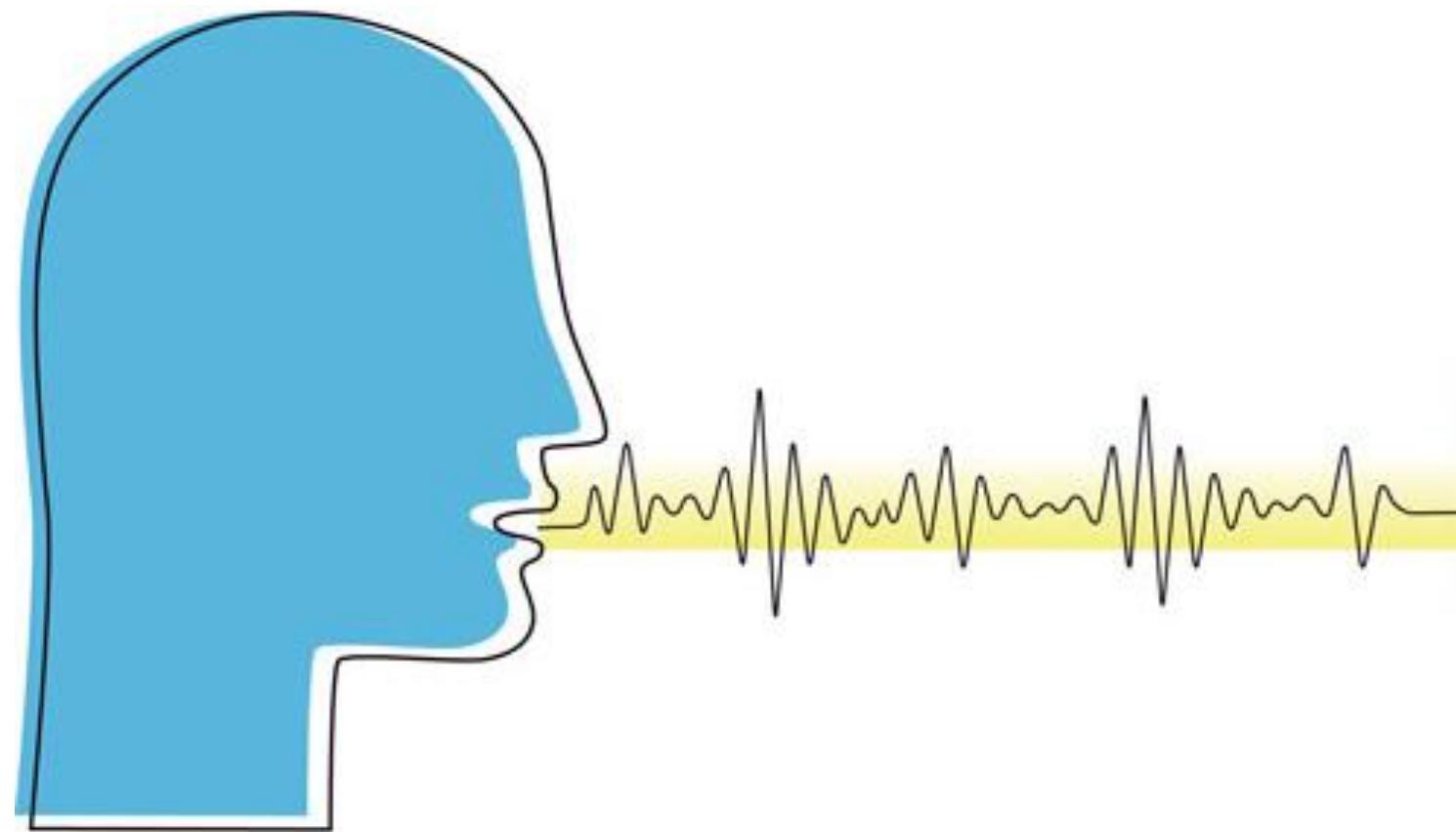


Filtros *digitales*

TRANSFORMATEC

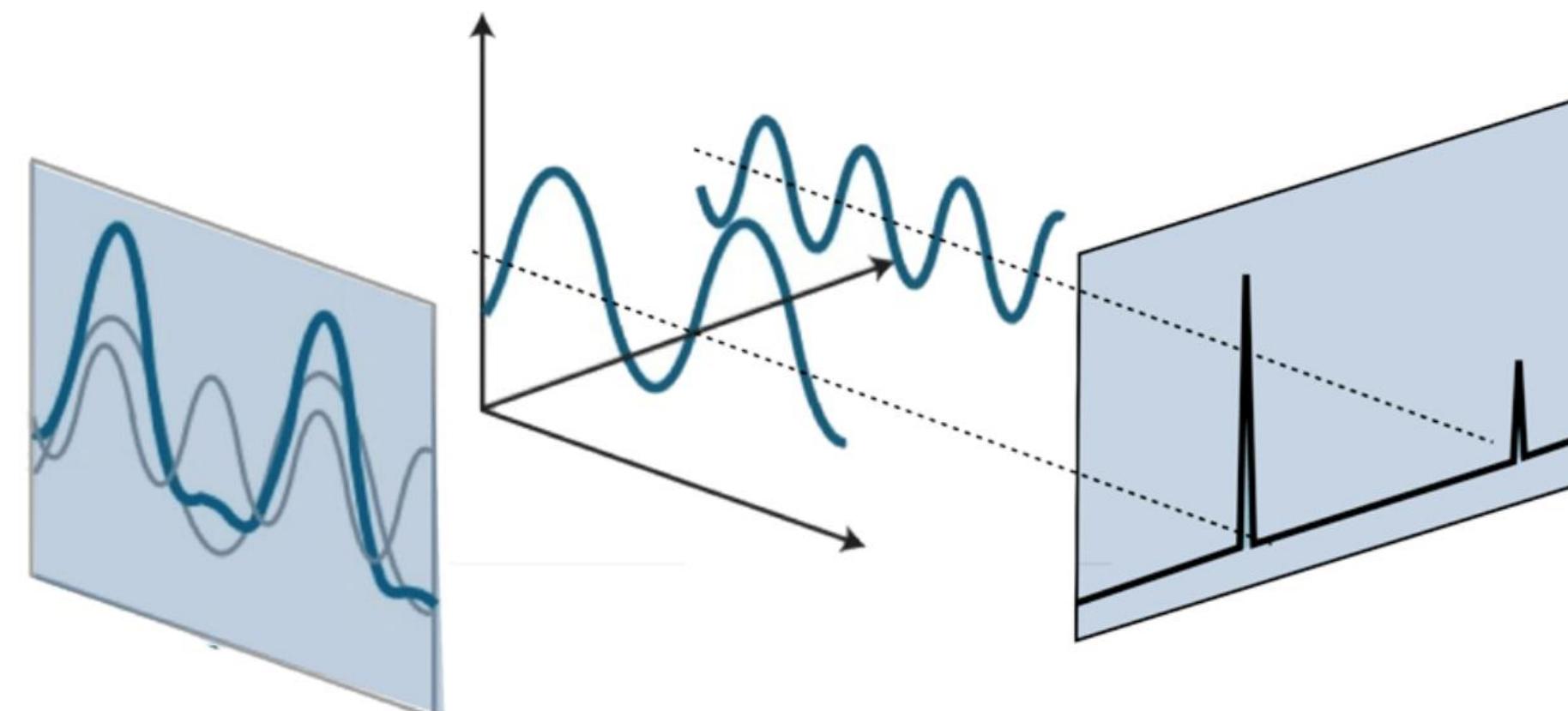


Señal



TRANSFORMATEC

Dominio de *Frecuencia*



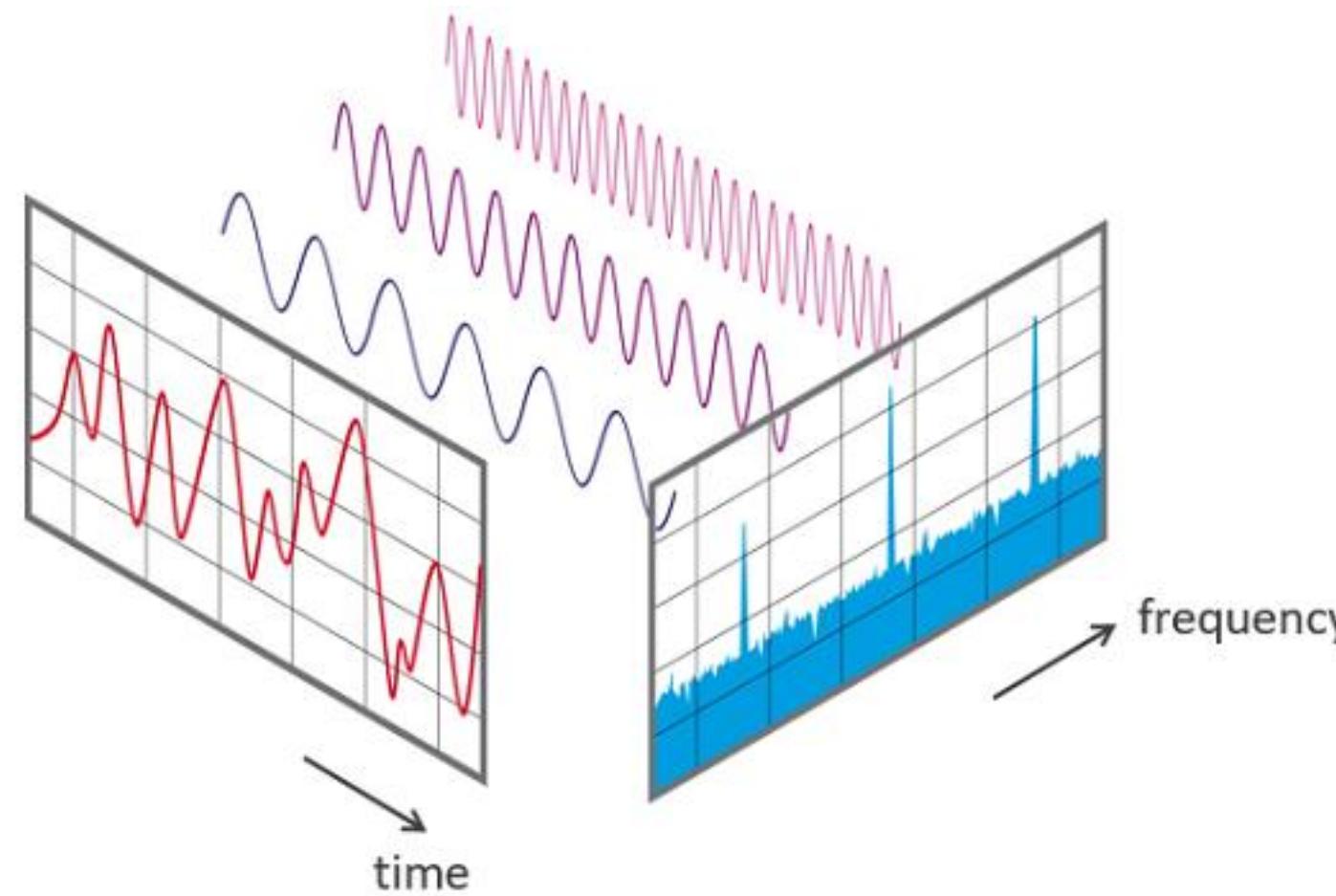
Time Domain
 $s(t)$

FT

Frequency Domain
 $S(\omega)$



Dominio de *Frecuencia*



Serie de Fourier

$$f(t) = a_0 + \sum_{n=1}^{\infty} \left(a_n \cos \frac{n\pi t}{L} + b_n \sin \frac{n\pi t}{L} \right)$$

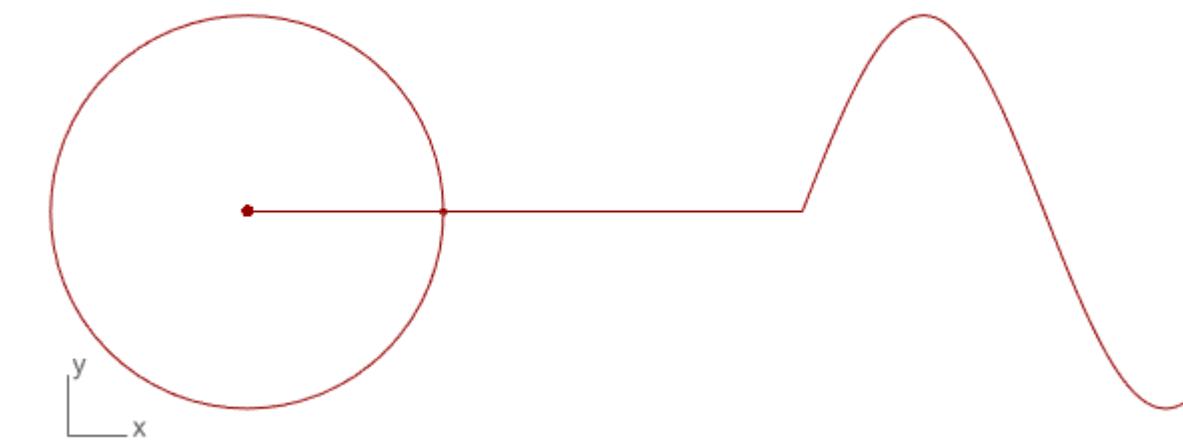
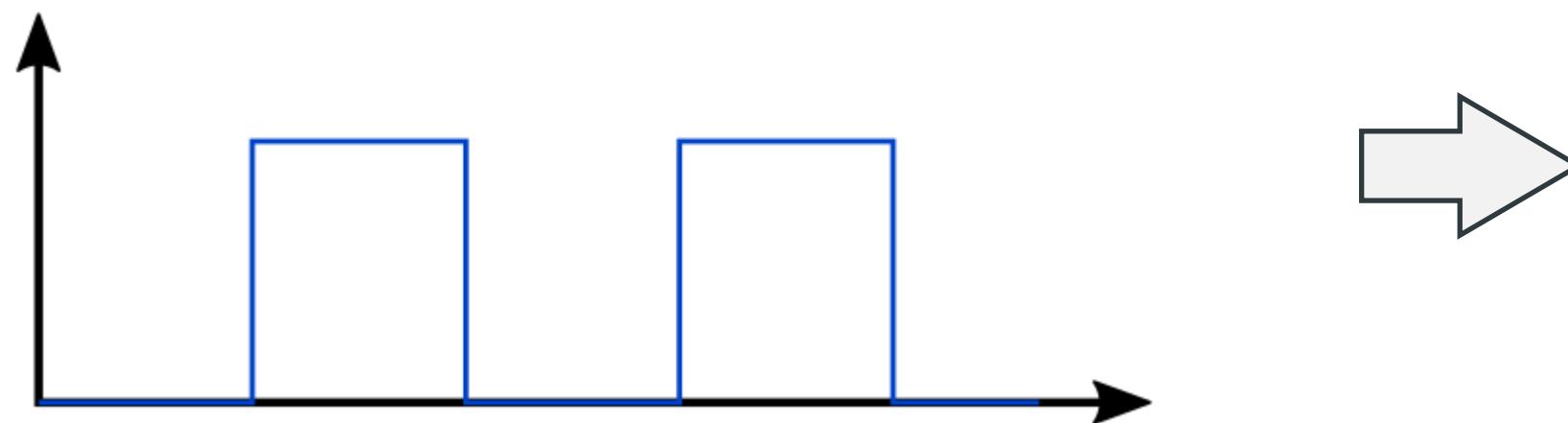
Transformada de Fourier

$$F(\omega) = \int_{-\infty}^{\infty} f(t) e^{-2\pi i \omega t} dt$$



TRANSFORMATEC

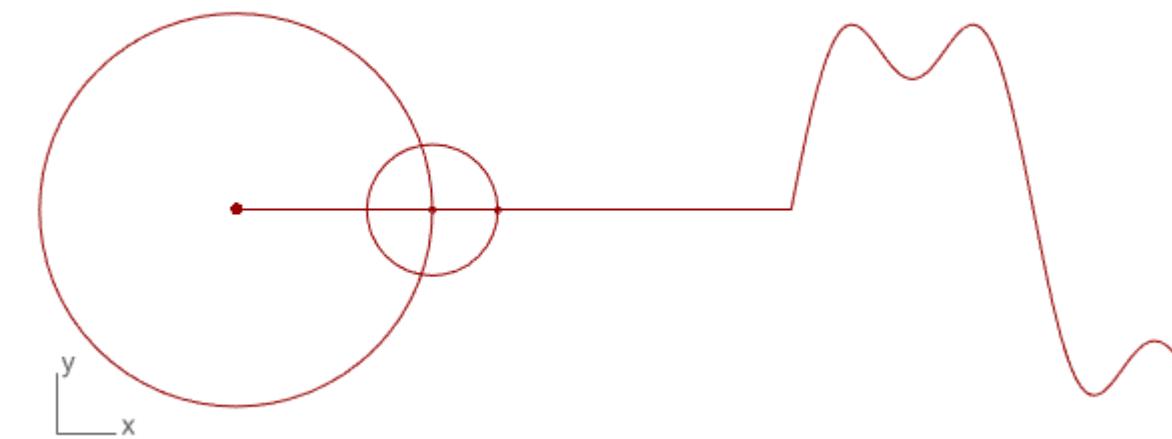
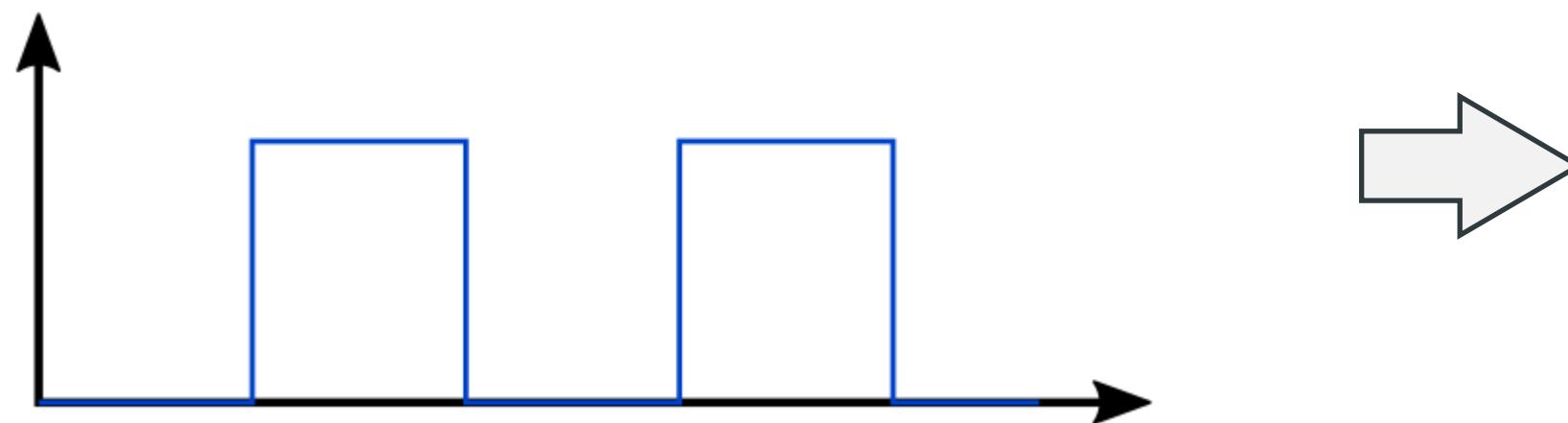
Dominio de *Frecuencia*



Primer componente



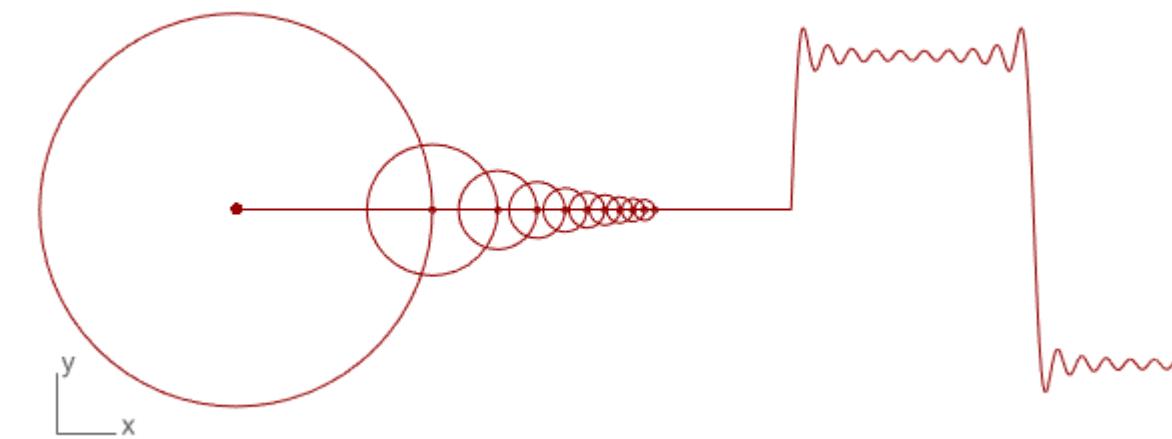
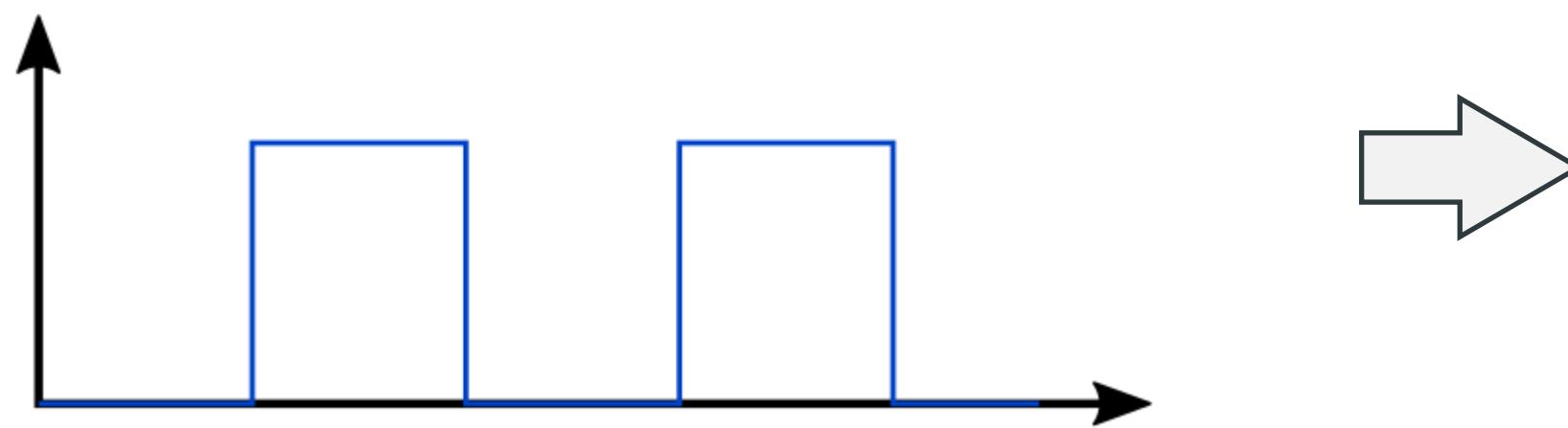
Dominio de *Frecuencia*



Primer y segundo componentes



Dominio de *Frecuencia*

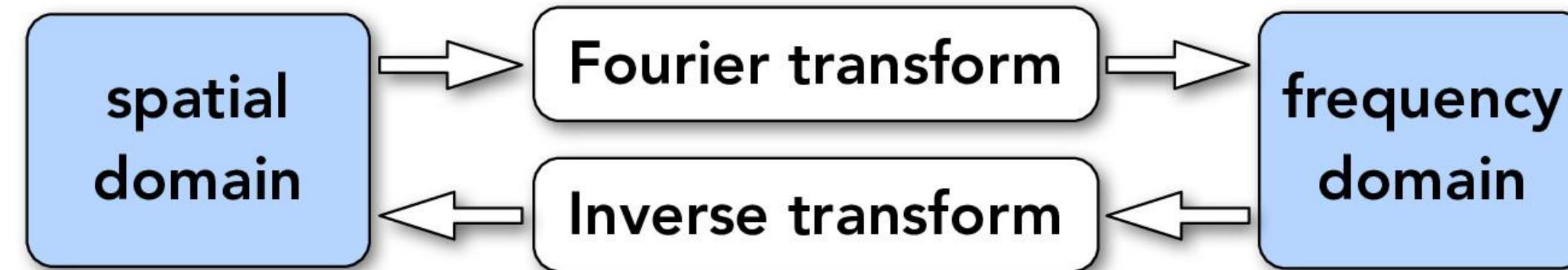


Primeros 10 componentes



Dominio de *Frecuencia*

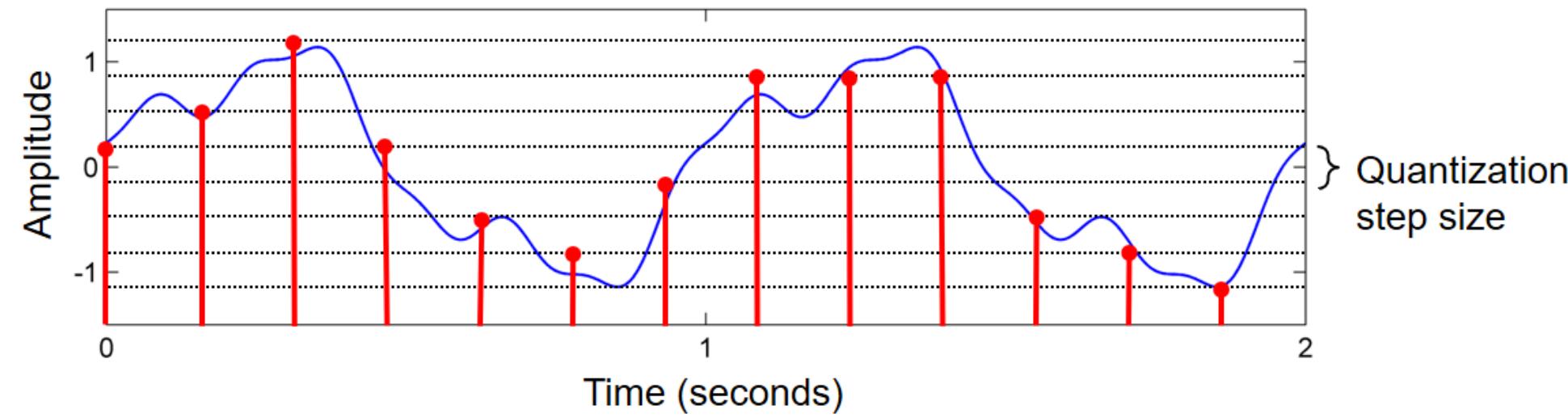
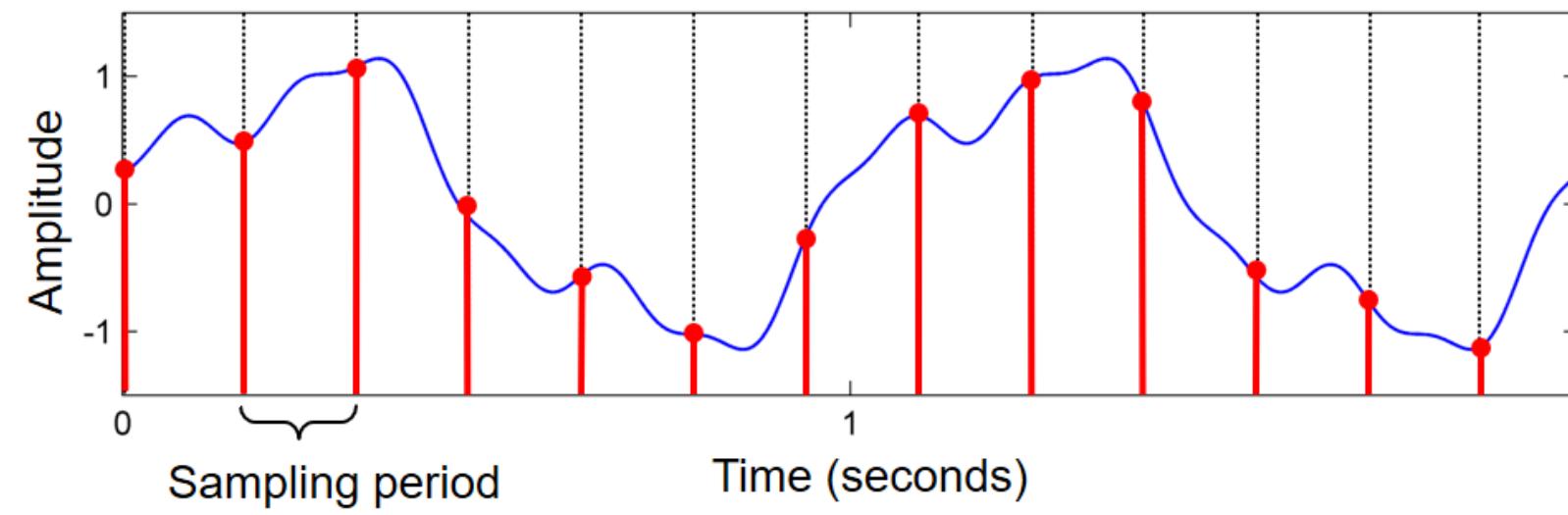
$$f(x) \quad F(\omega) = \int_{-\infty}^{\infty} f(x)e^{-2\pi i \omega x} dx \quad F(\omega)$$



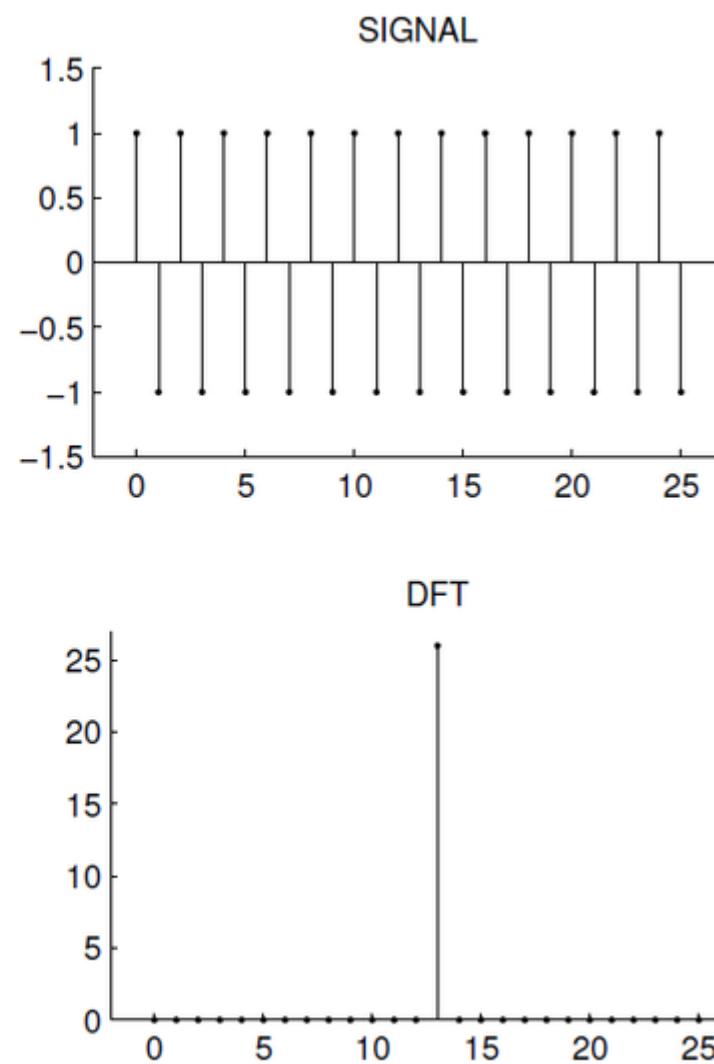
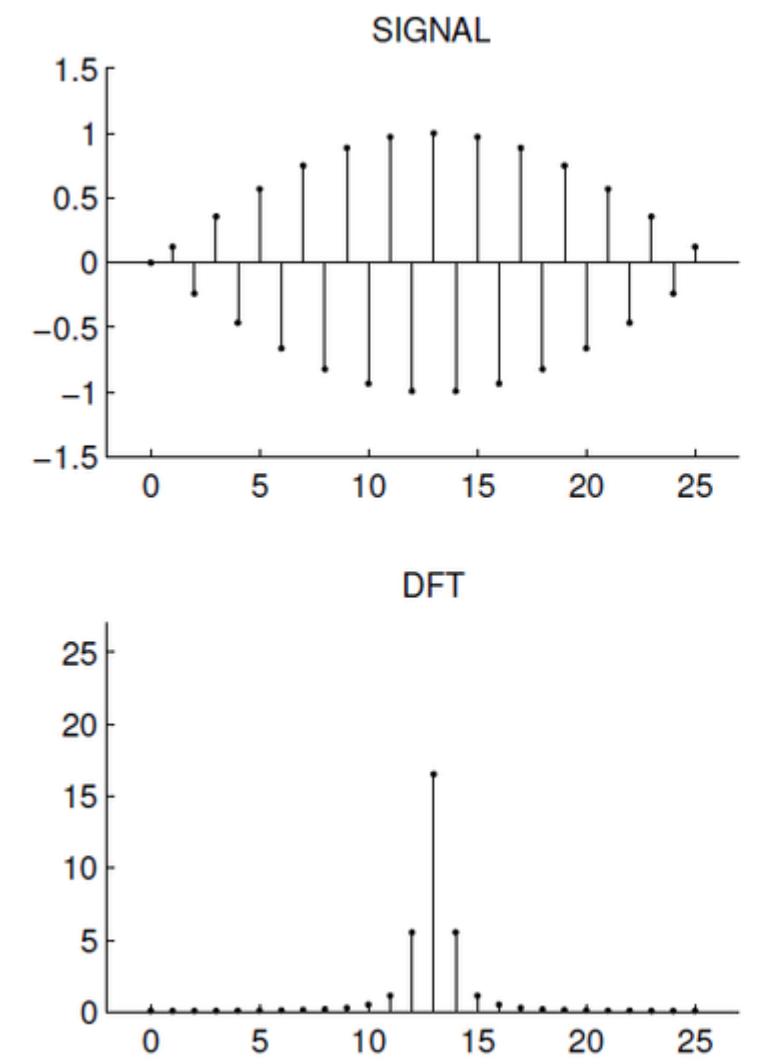
$$f(x) = \int_{-\infty}^{\infty} F(\omega)e^{2\pi i \omega x} d\omega$$



Discretización



Discrete Fourier transform



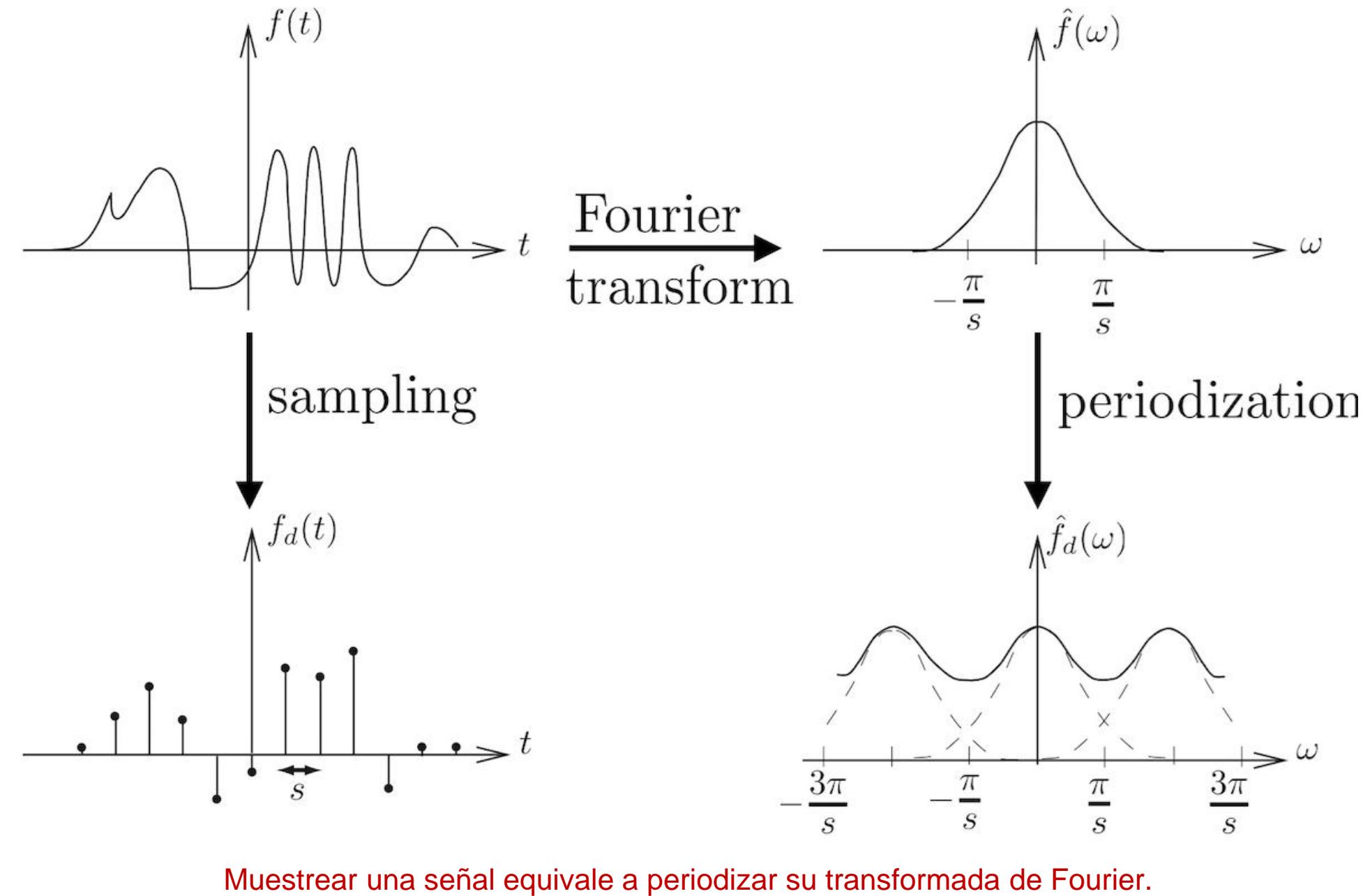
Transformada

$$X_k = \sum_{n=0}^{N-1} x_n \cdot e^{-i \frac{2\pi}{N} kn}$$

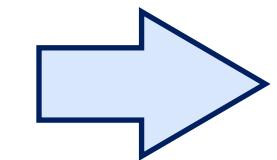
$$X_k = \sum_{n=0}^{N-1} x_n \cdot \left[\cos\left(\frac{2\pi}{N} kn\right) - i \cdot \sin\left(\frac{2\pi}{N} kn\right) \right]$$



Discrete Fourier transform



Aliasing



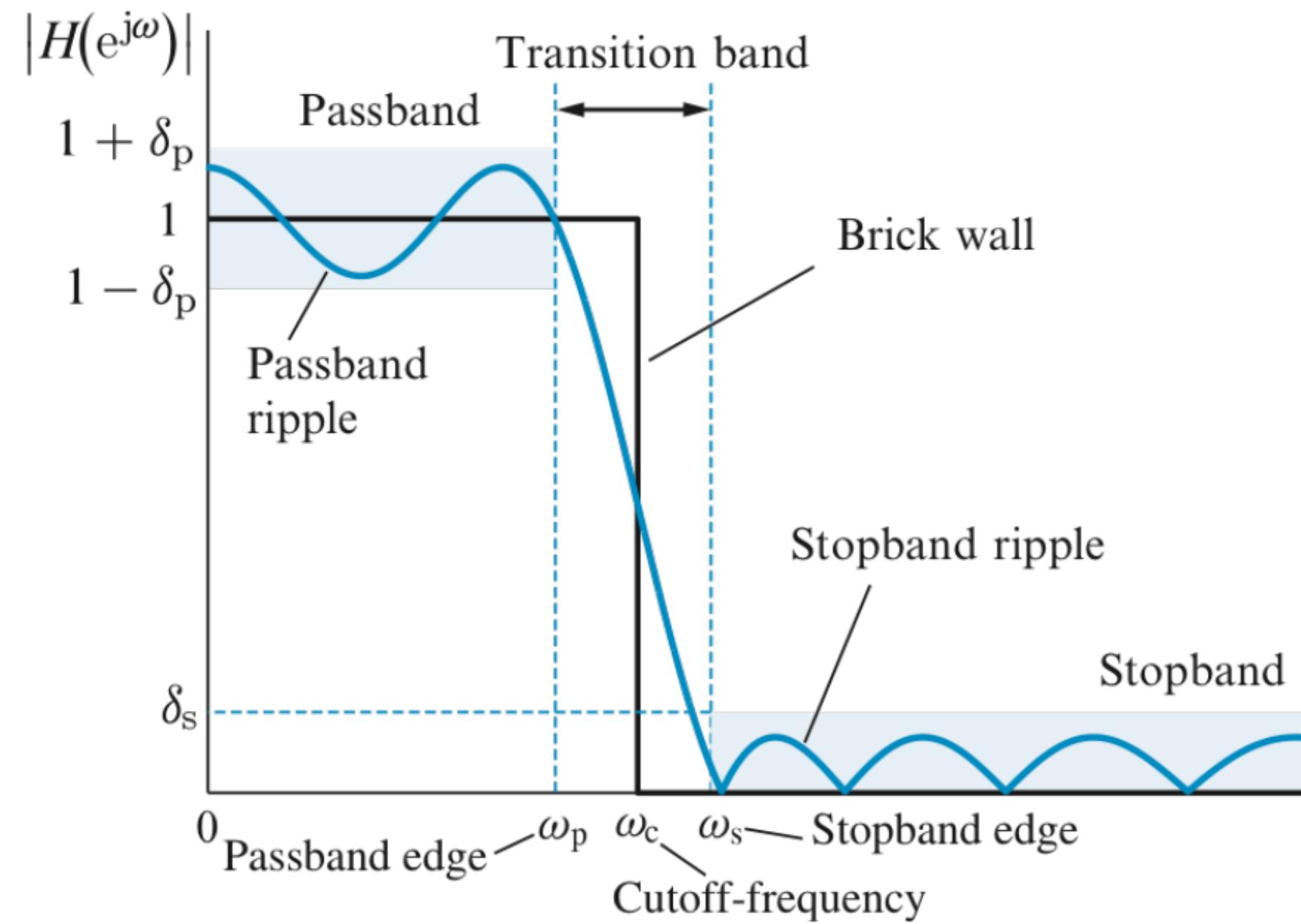
Submuestreo



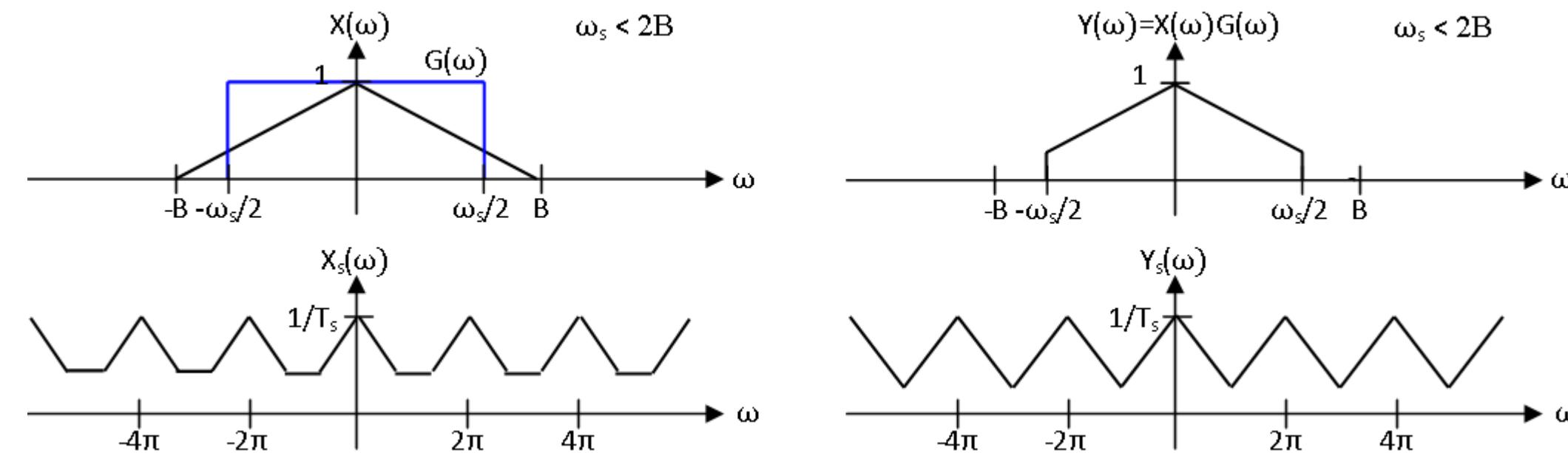
Patrón de muaré



Filtro



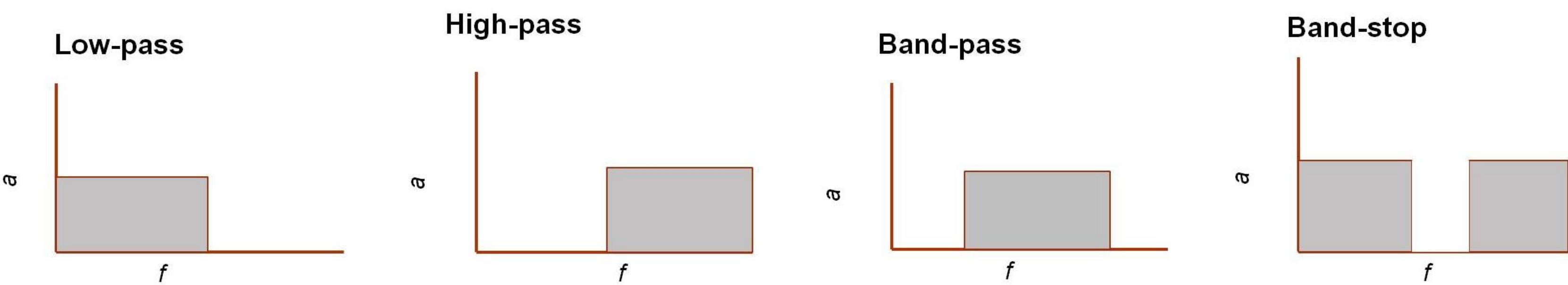
Anti-alias*sing filter*



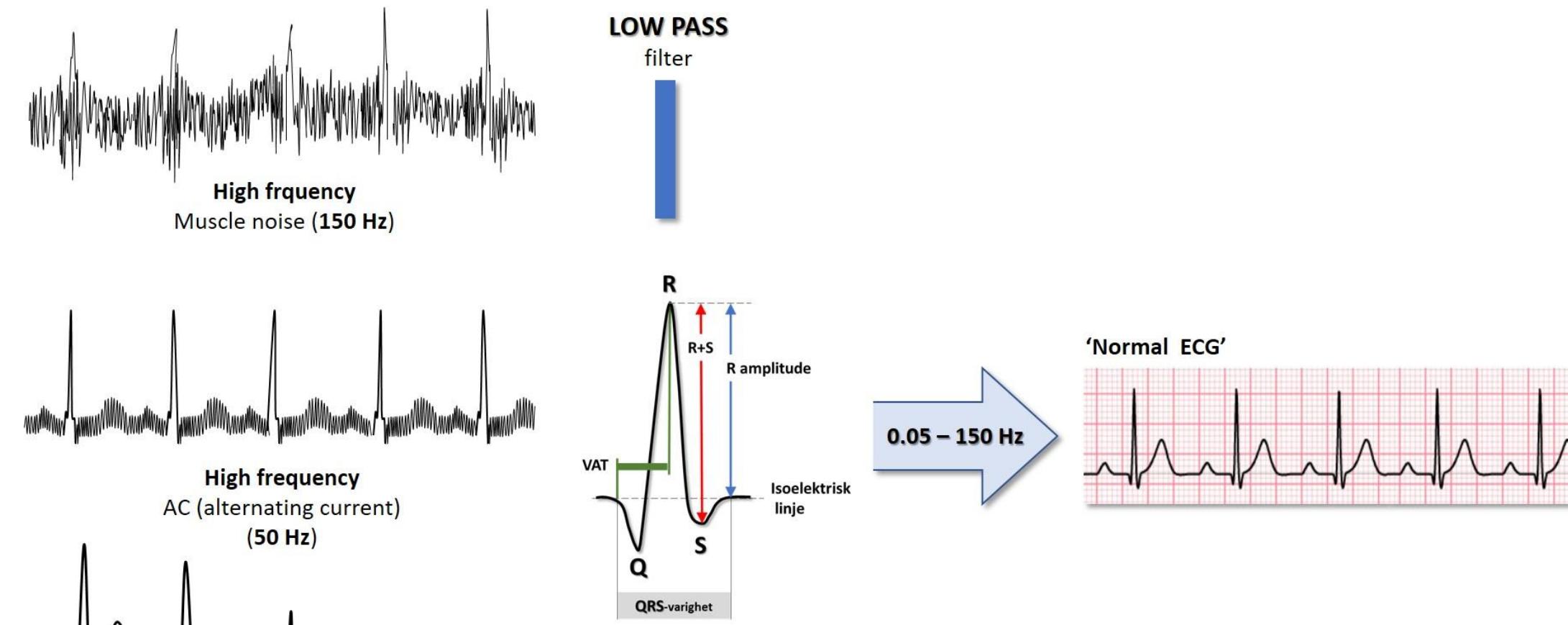
Frecuencia de Nyquist: $f_N = \frac{f_s}{2}$



Filter types



Filtro

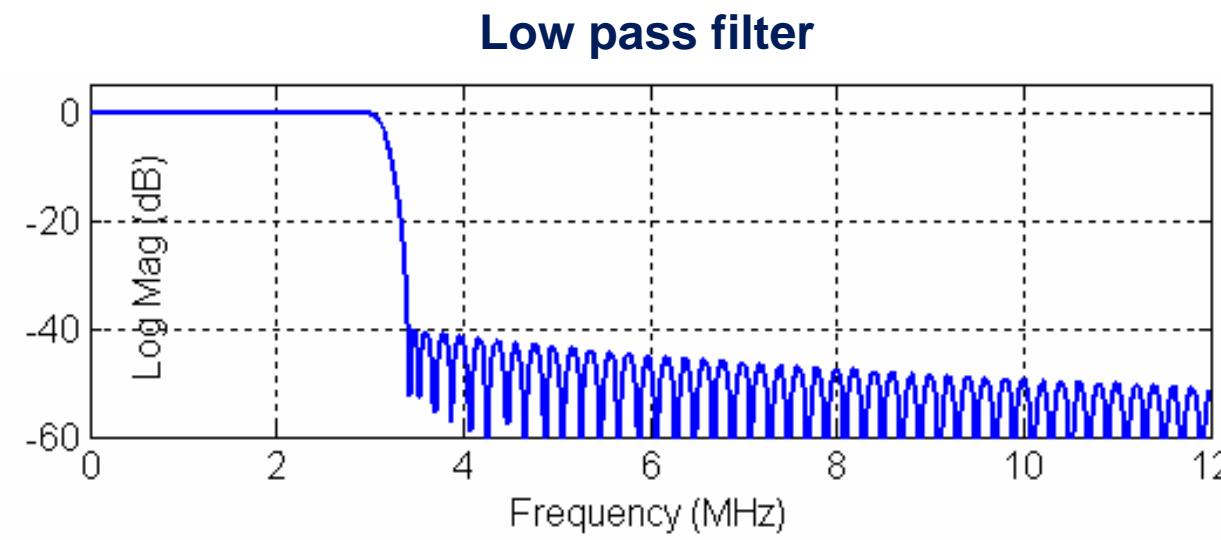


HIGH PASS filter

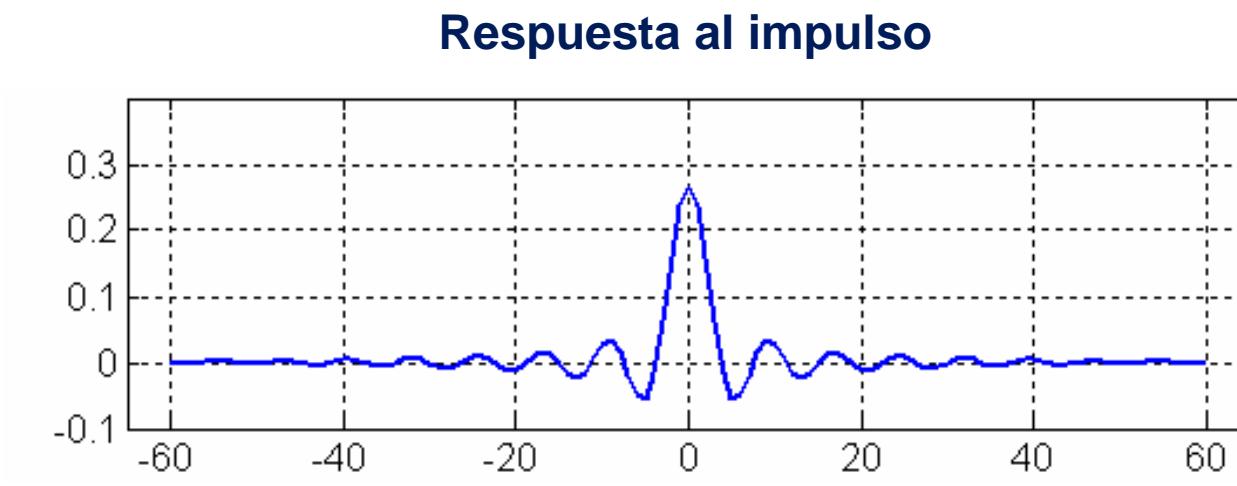
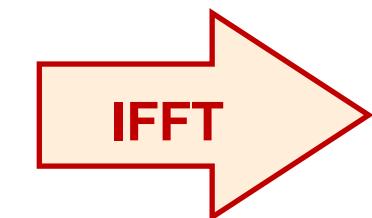
© MfBjørner



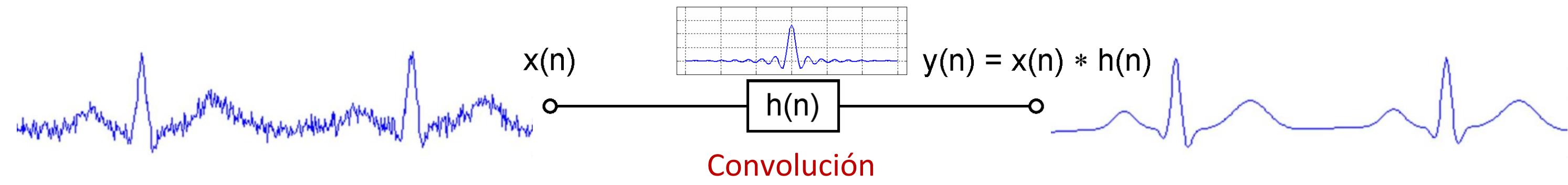
Filtro



$$H[\omega]$$



Filtro

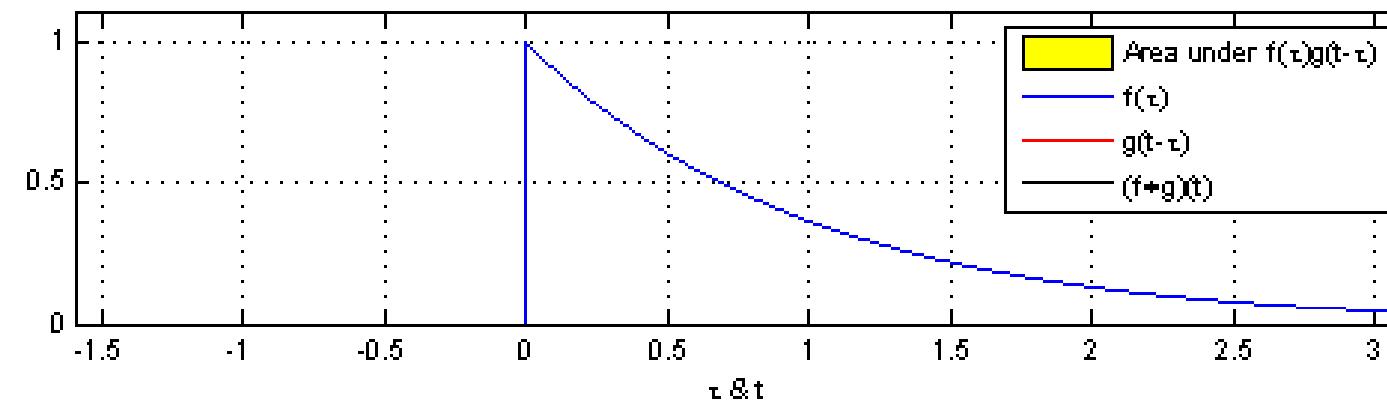


Convolución



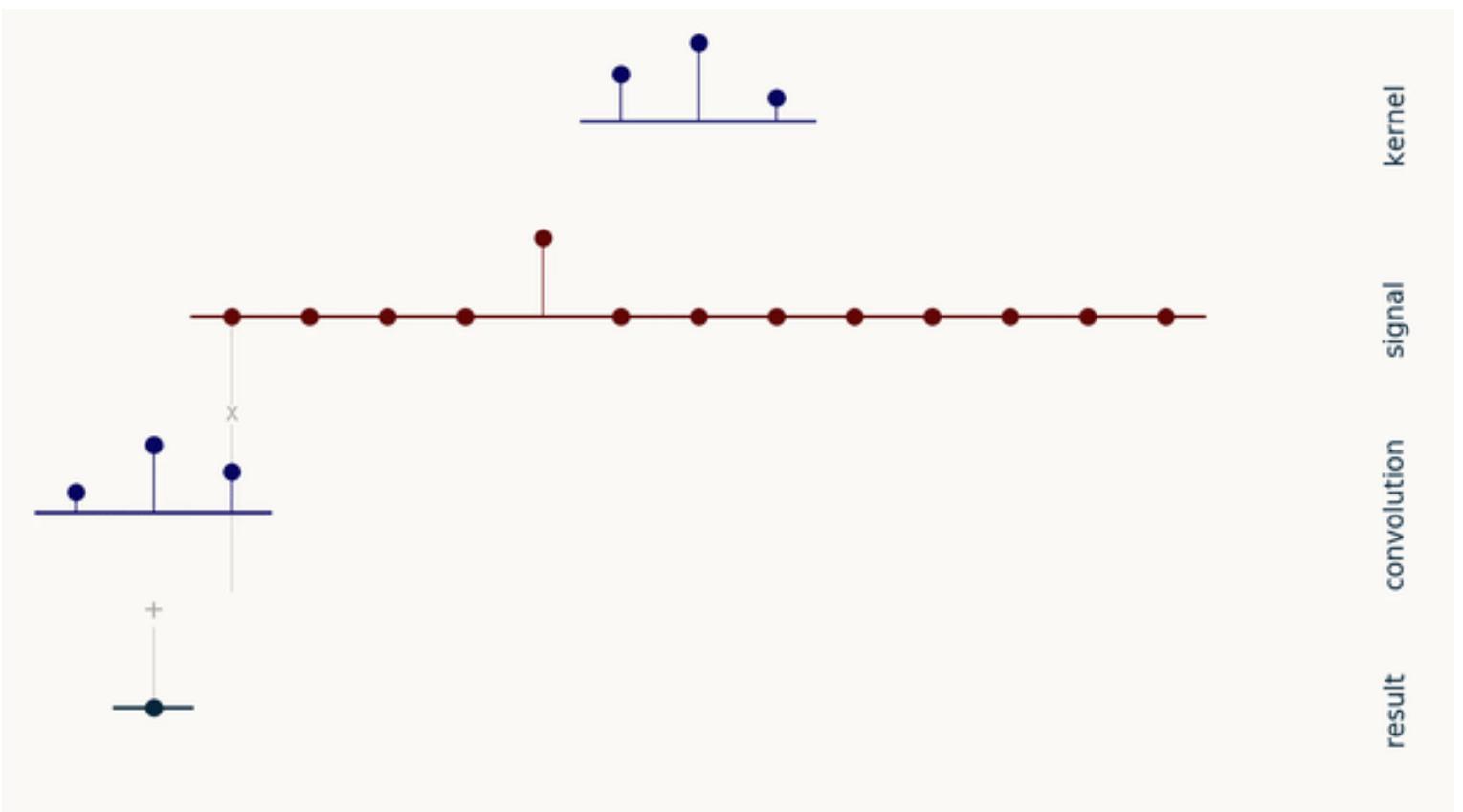
Convolution

$$(f * g)(t) = \int_{-\infty}^{\infty} f(\tau)g(t - \tau) d\tau$$



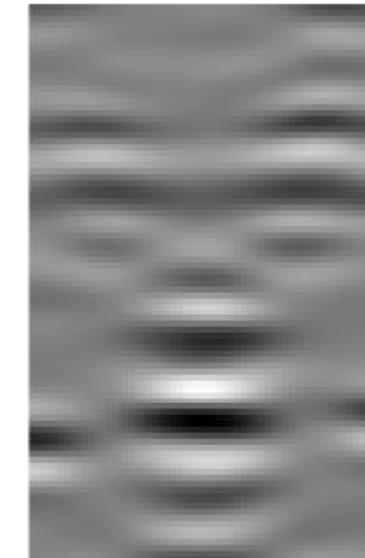
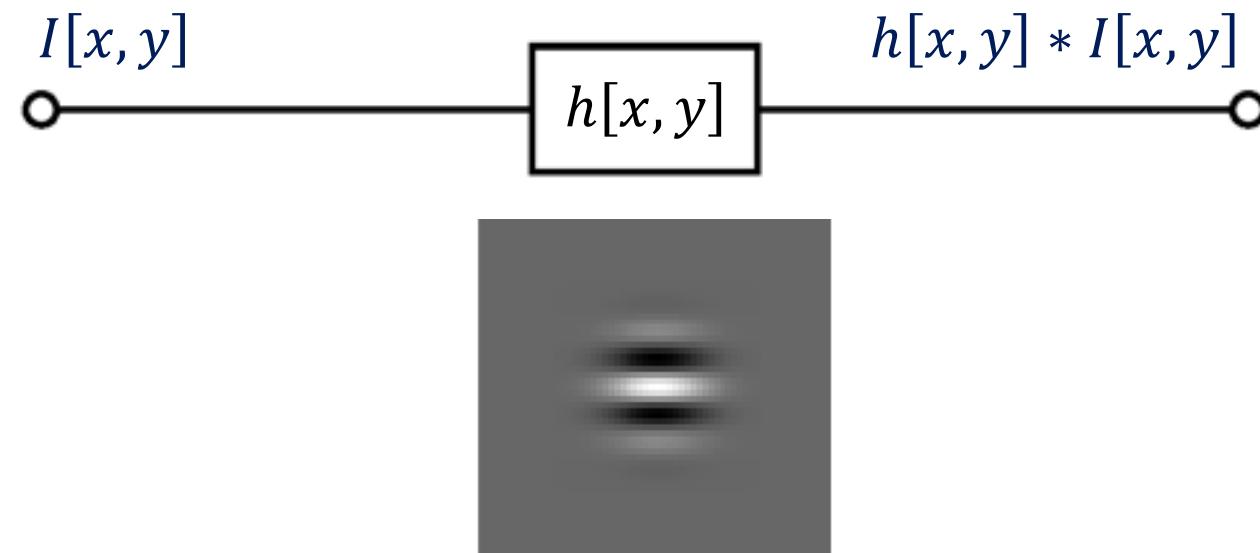
Discrete convolution

$$(f * g)[n] = \sum_{k=-\infty}^{\infty} f[k] \cdot g[n - k]$$

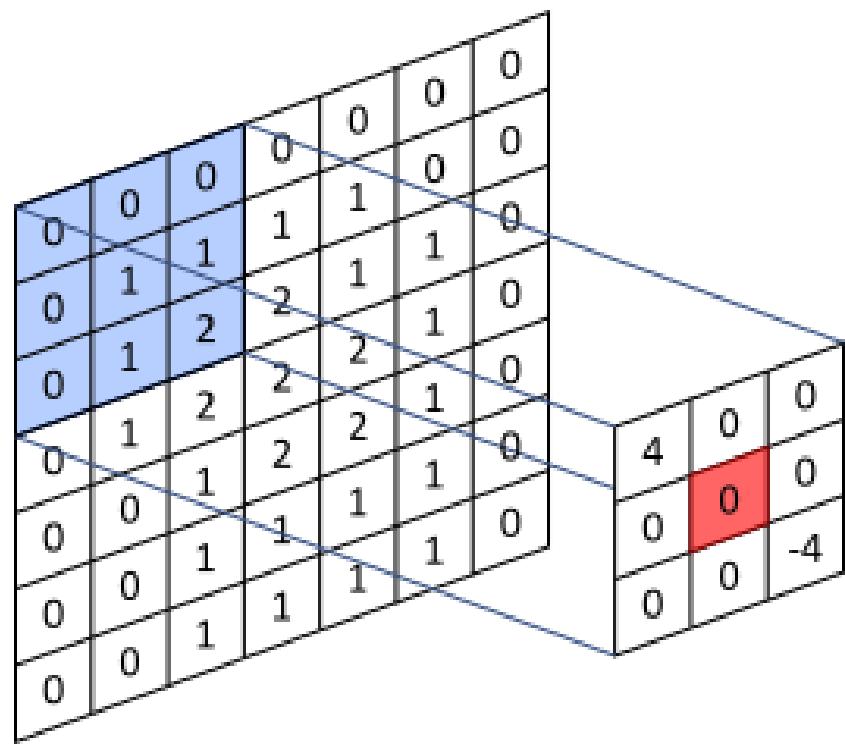


2D Discrete convolution

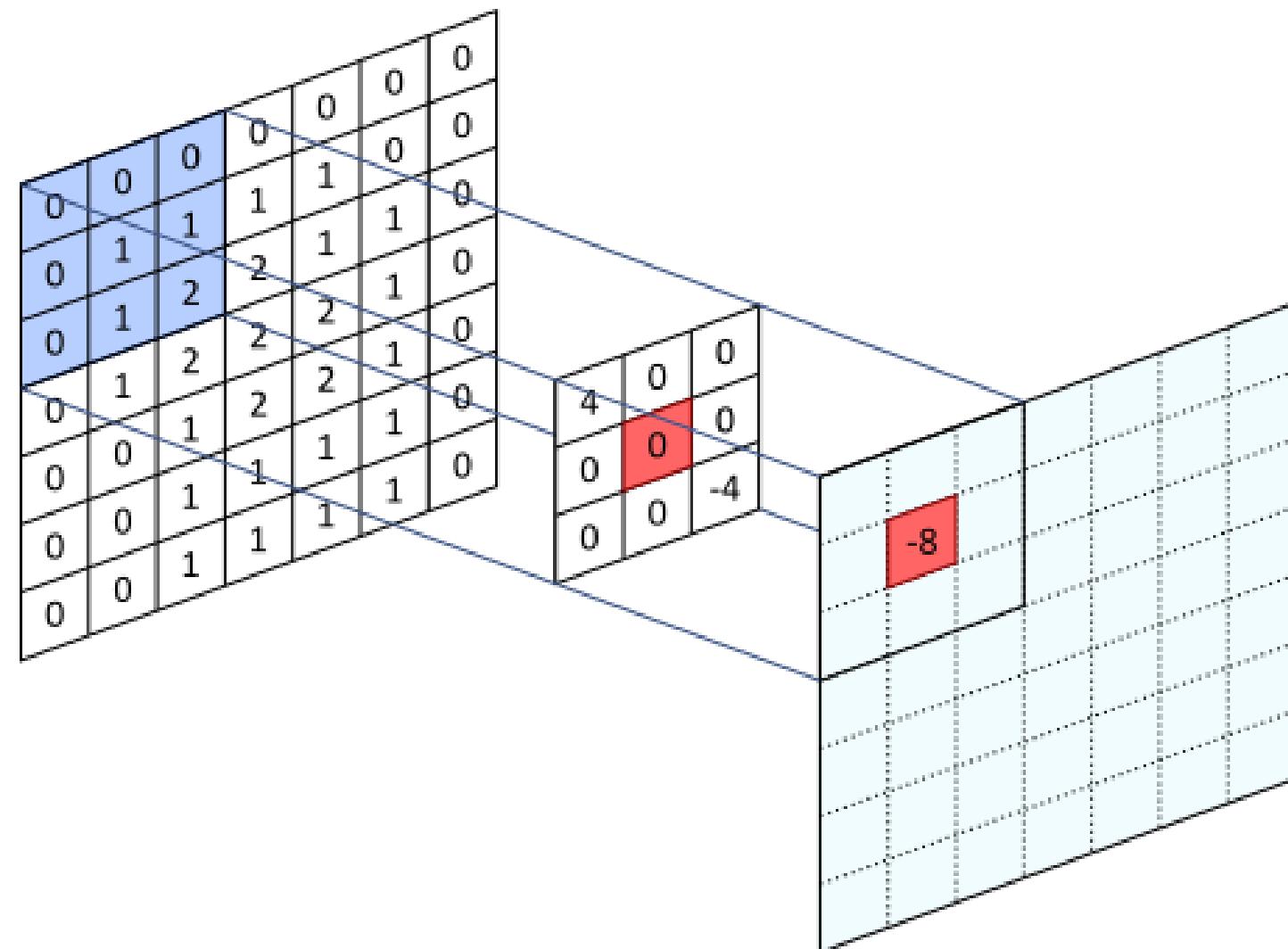
$$h[x, y] * I[x, y] = \sum_{i=-a}^a \sum_{j=-b}^b h[i, j]I[x - i, y - j]$$



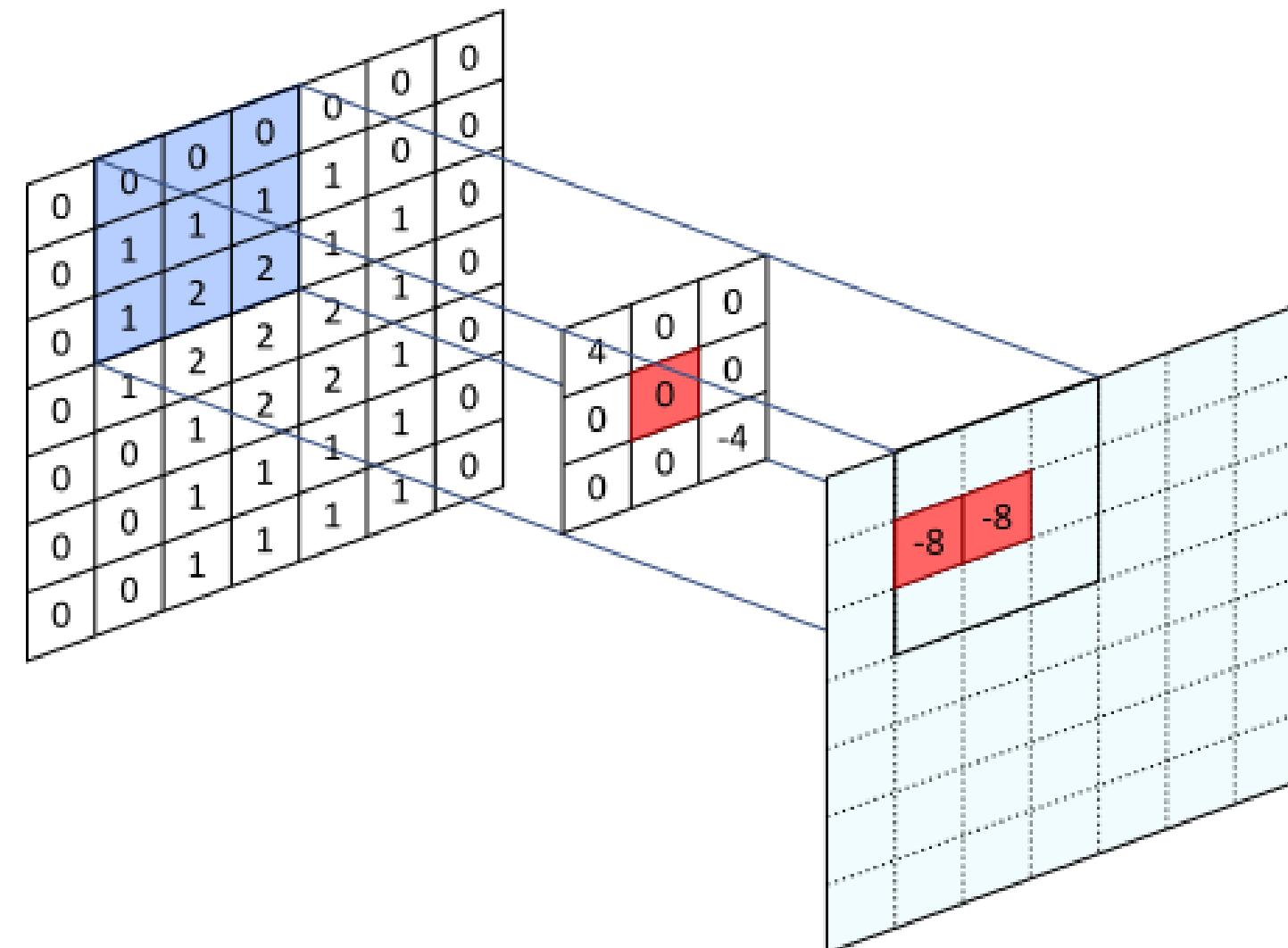
2D Discrete convolution



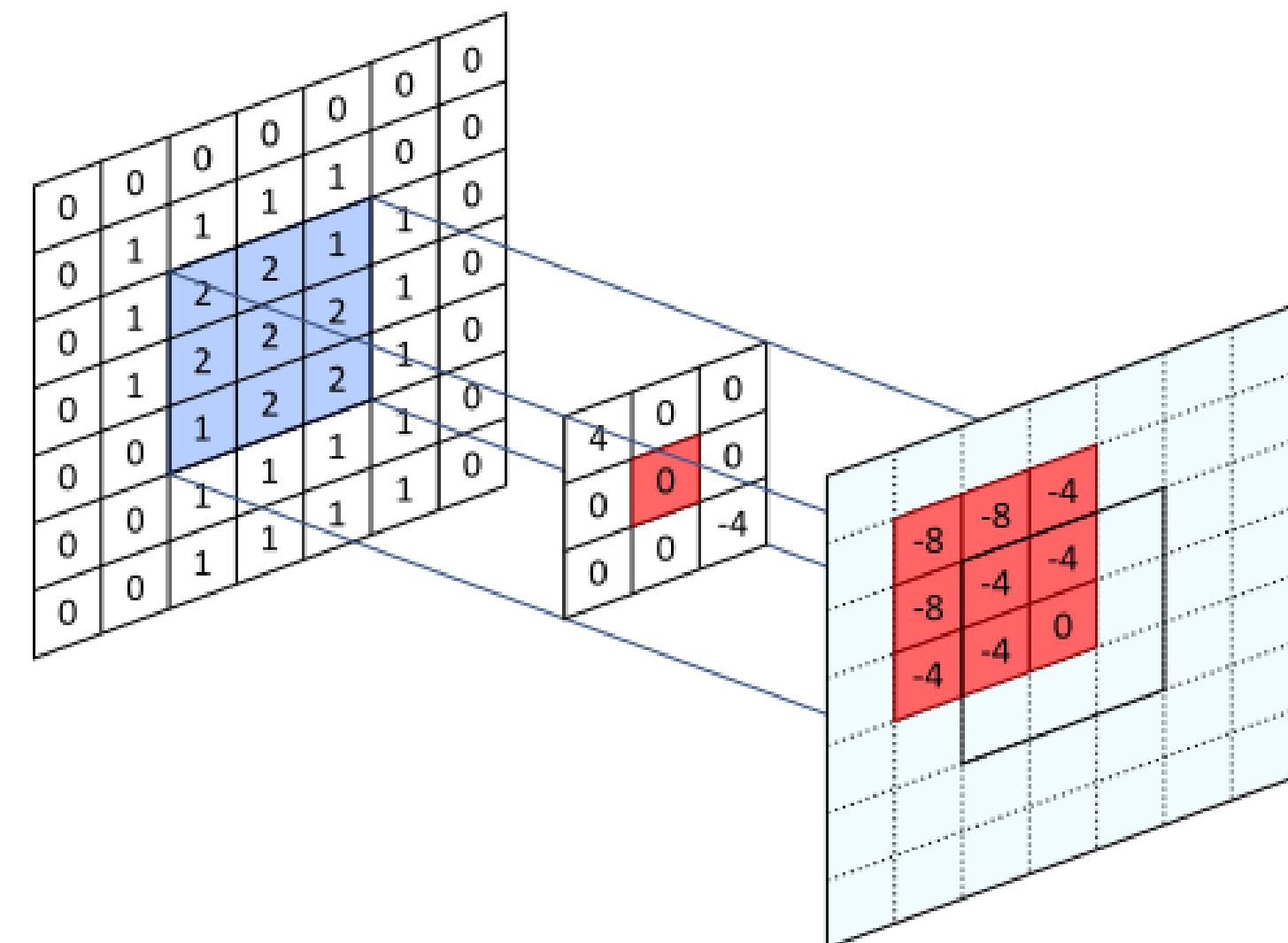
2D Discrete convolution



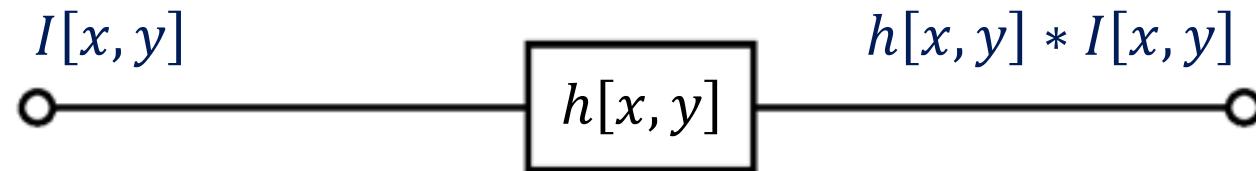
2D Discrete convolution



2D Discrete convolution



2D Discrete convolution



0	0	0	5	0	0	0
0	5	18	32	18	5	0
0	18	64	100	64	18	0
5	32	100	100	100	32	5
0	18	64	100	64	18	0
0	5	18	32	18	5	0
0	0	0	5	0	0	0

-1	-1	-1
-1	8	-1
-1	-1	-1

0	0	-1	0	0
0	-1	-2	-1	0
-1	-2	16	-2	-1
0	-1	-2	-1	0
0	0	-1	0	0



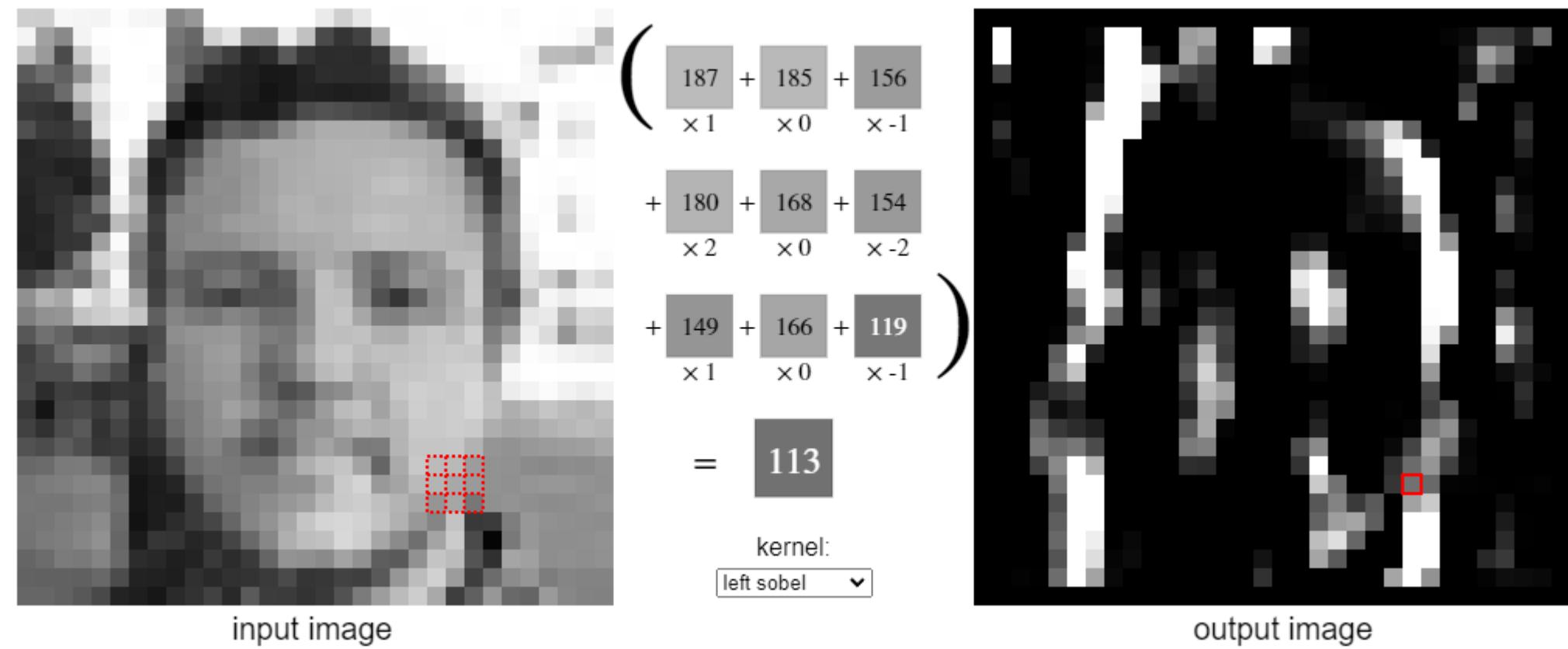
Gaussian blur

Edge detection

Laplacian filter



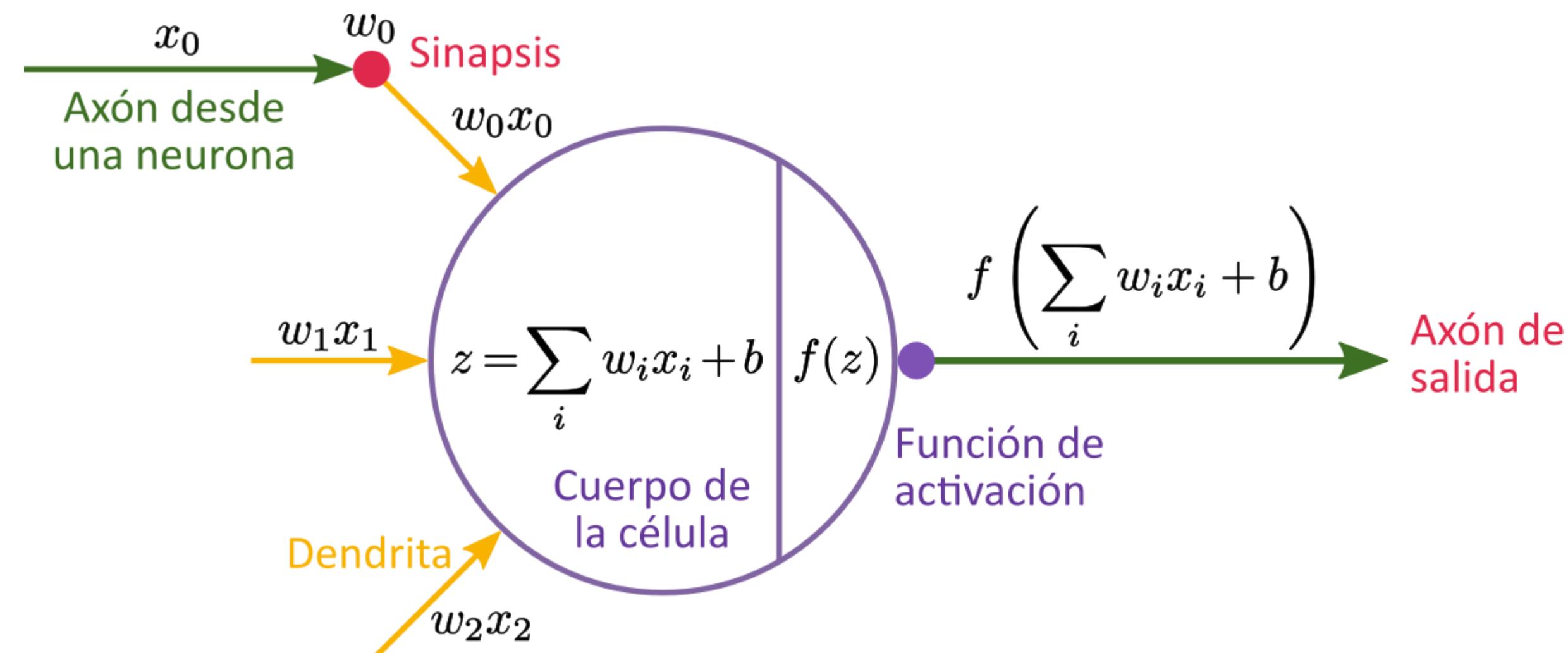
2D Discrete convolution



<https://setosa.io/ev/image-kernels/>



Red neuronal artificial

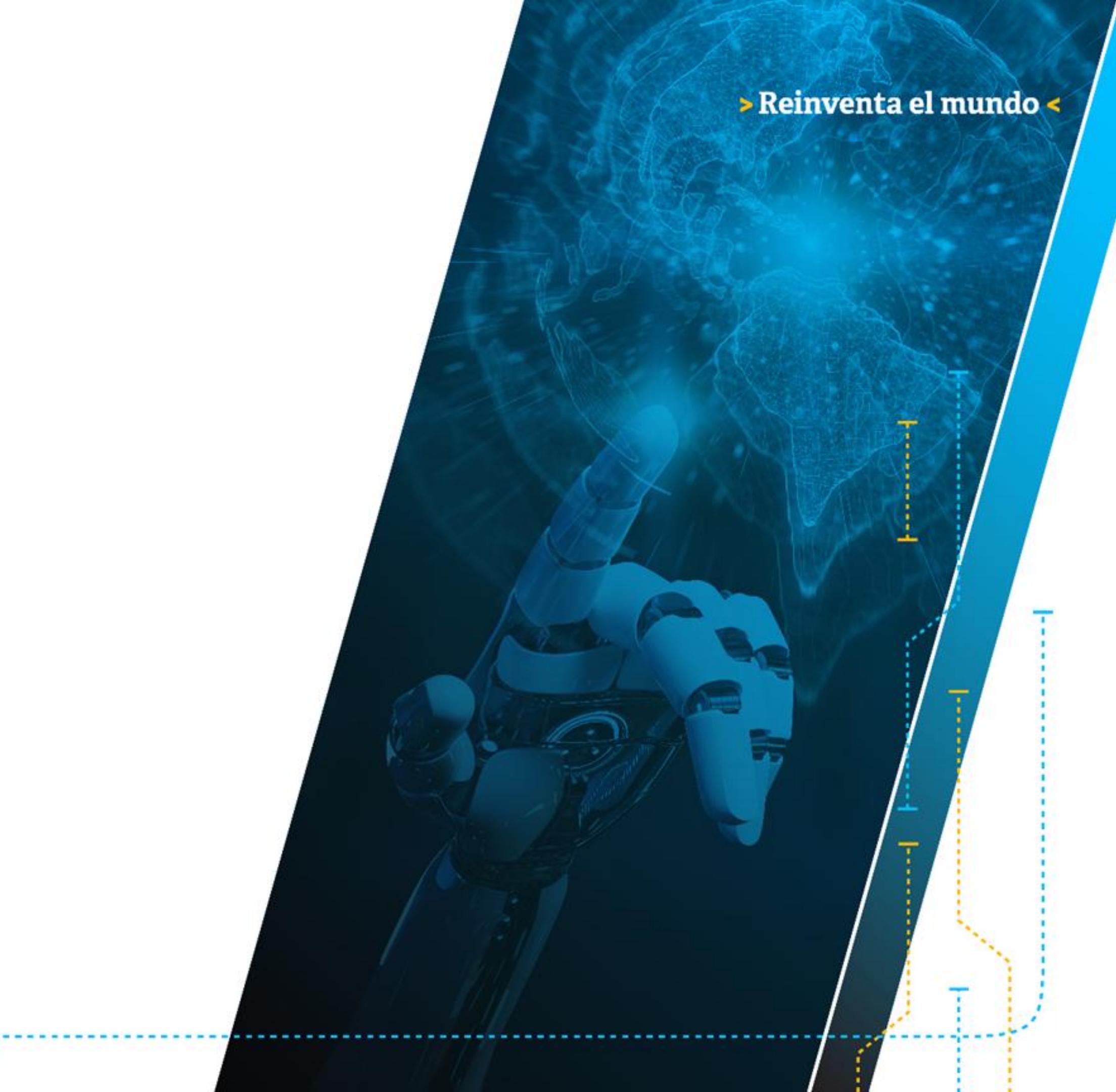


Frank Rosenblatt (1958) "The perceptron: a probabilistic model for information storage and organization in the brain".
Psychological review 65.6 (1958): 386.

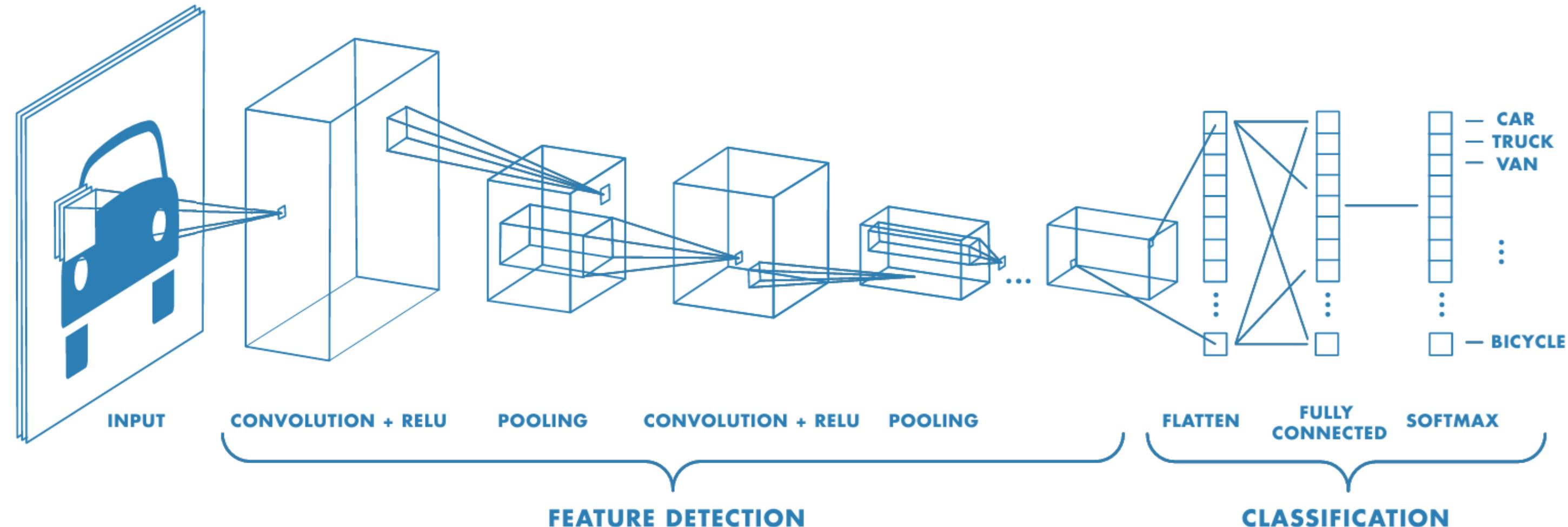
2.



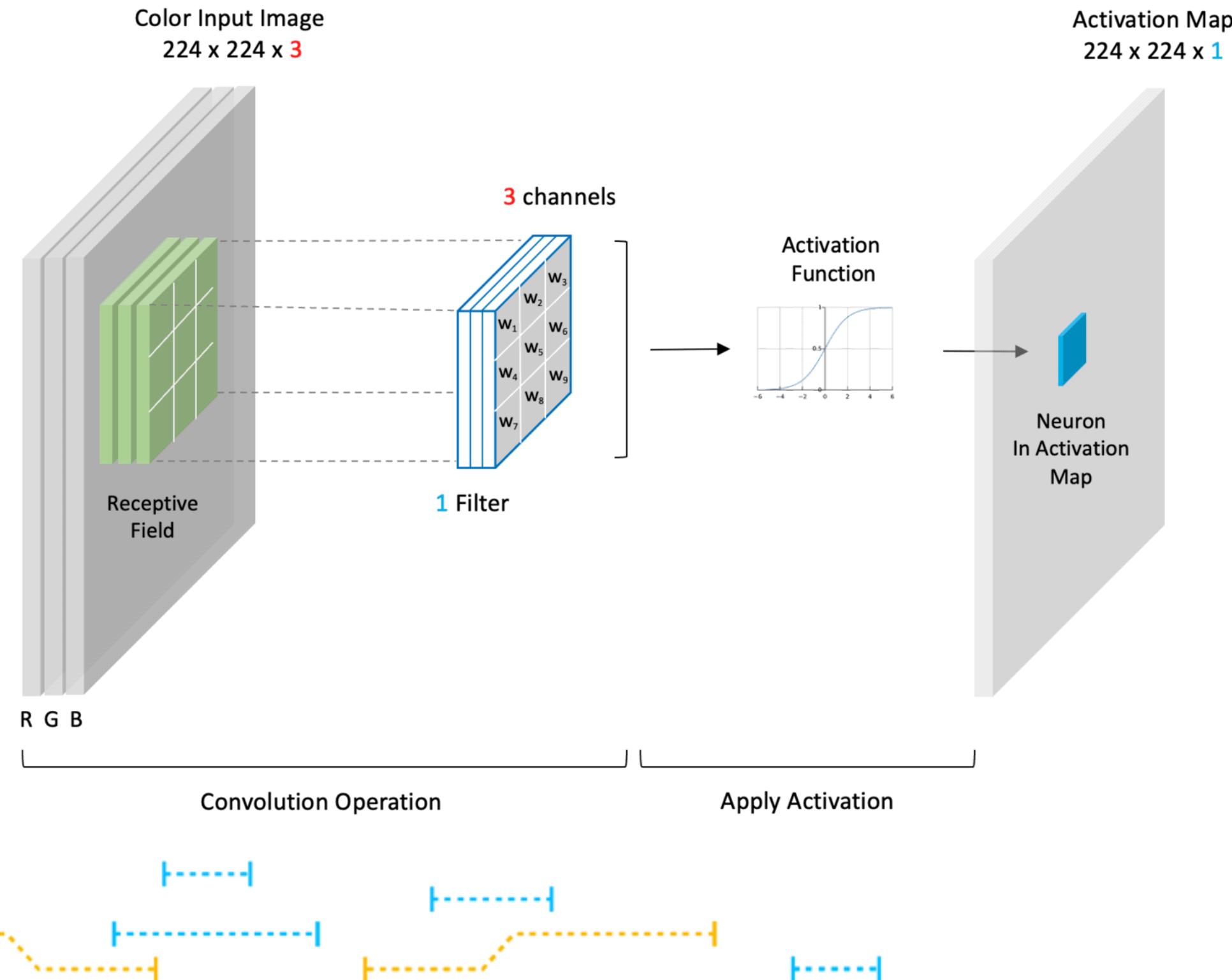
Convolutional *neural network*



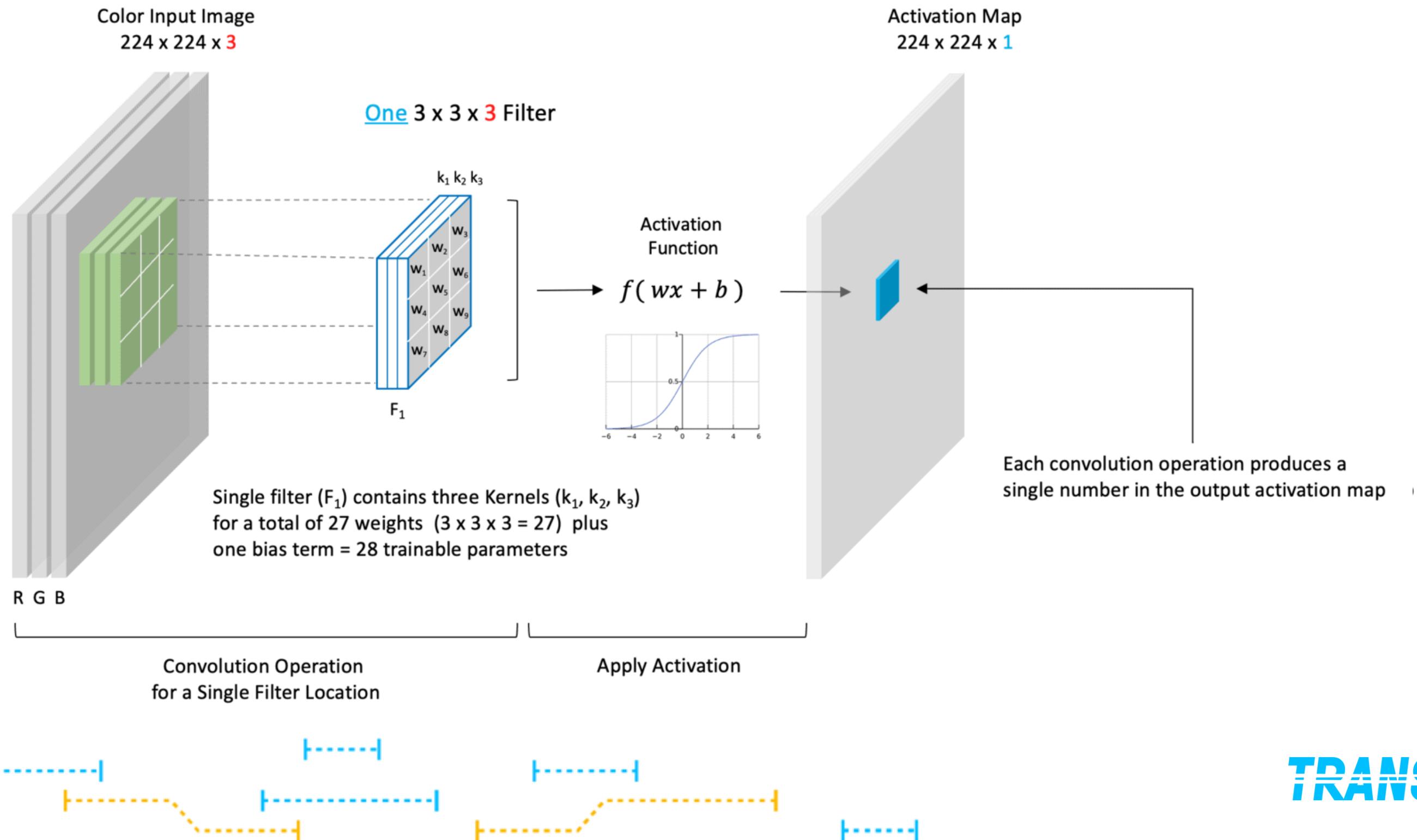
Convolutional neural network



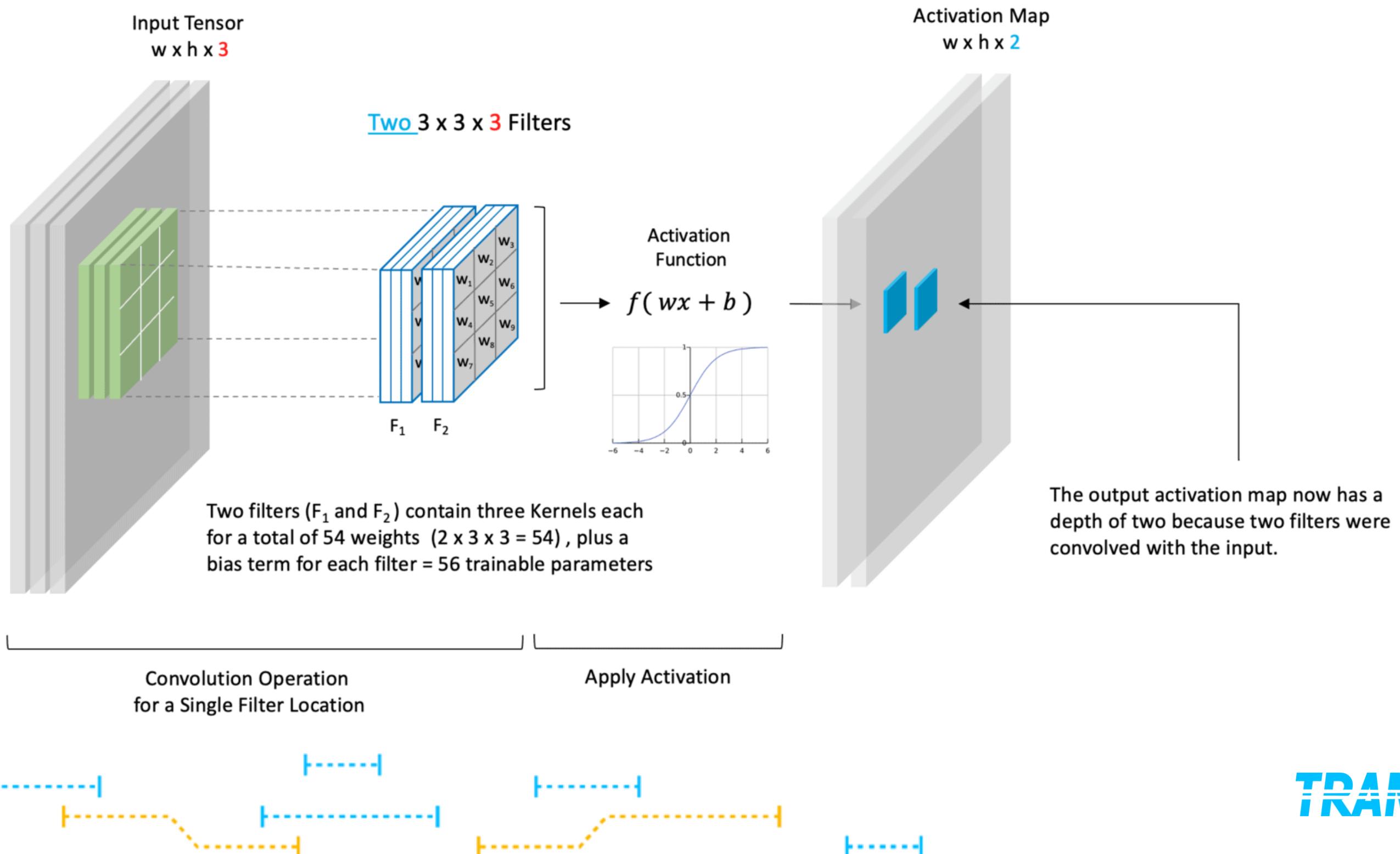
Convolutional neural network



Convolutional neural network



Convolutional neural network



Convolutional neural network

0	0	0	0	0	0	0	...
0	156	155	156	158	158	158	...
0	153	154	157	159	159	159	...
0	149	151	155	158	159	159	...
0	146	146	149	153	158	158	...
0	145	143	143	148	158	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	0	...
0	167	166	167	169	169	169	...
0	164	165	168	170	170	170	...
0	160	162	166	169	170	170	...
0	156	156	159	163	168	168	...
0	155	153	153	158	168	168	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	0	...
0	163	162	163	165	165	165	...
0	160	161	164	166	166	166	...
0	156	158	162	165	166	166	...
0	155	155	158	162	167	167	...
0	154	152	152	157	167	167	...
...

Input Channel #3 (Blue)

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1



308

1	0	0
1	-1	-1
1	0	-1

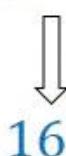
Kernel Channel #2



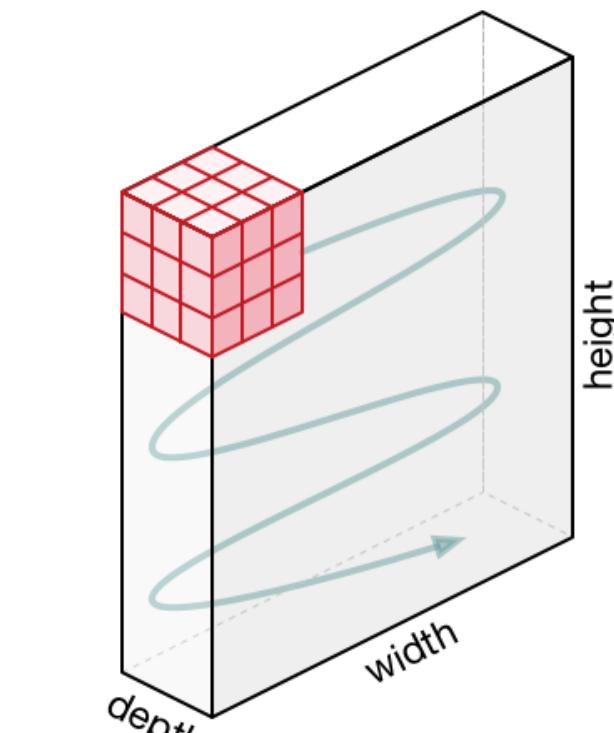
-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



164



Output

-25			...
			...
			...
			...
...

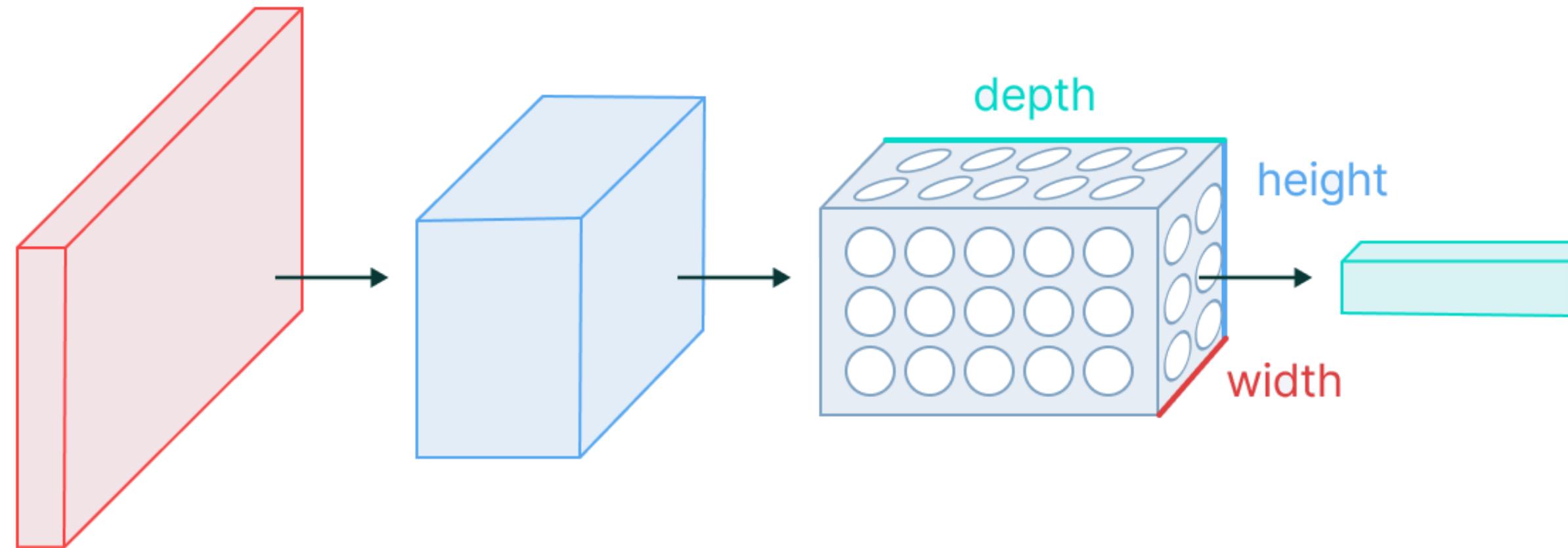
$$+ 1 = -25$$

$$\uparrow$$

Bias = 1

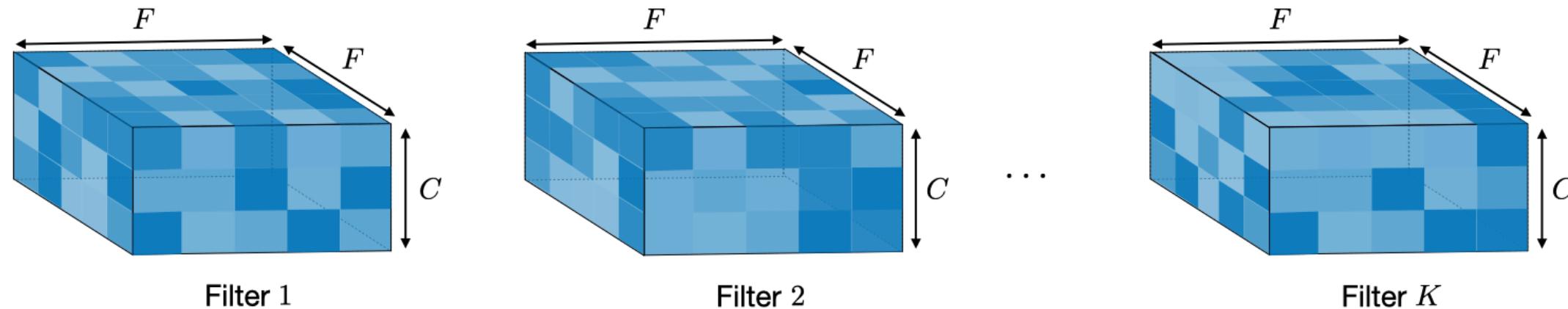


Convolutional *neural network*



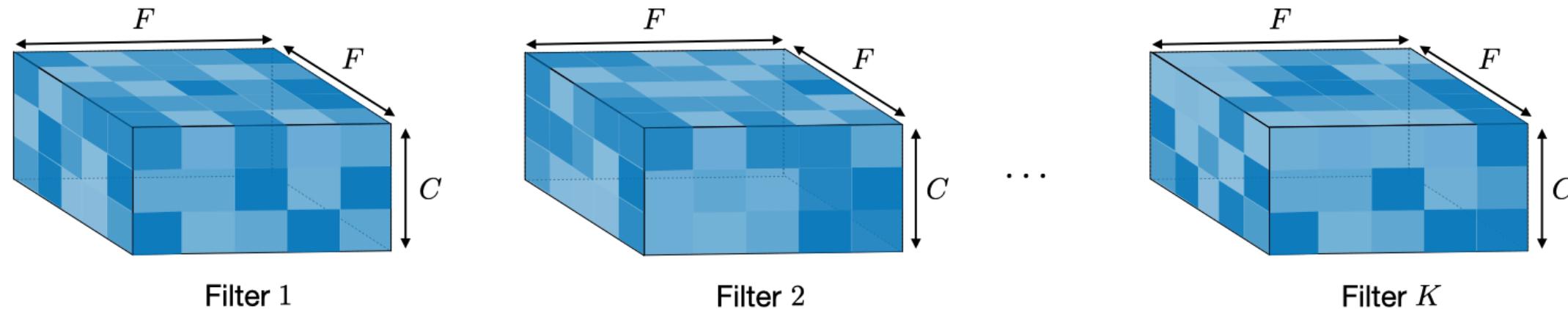
Convolutional neural network

Dimensiones del filtro

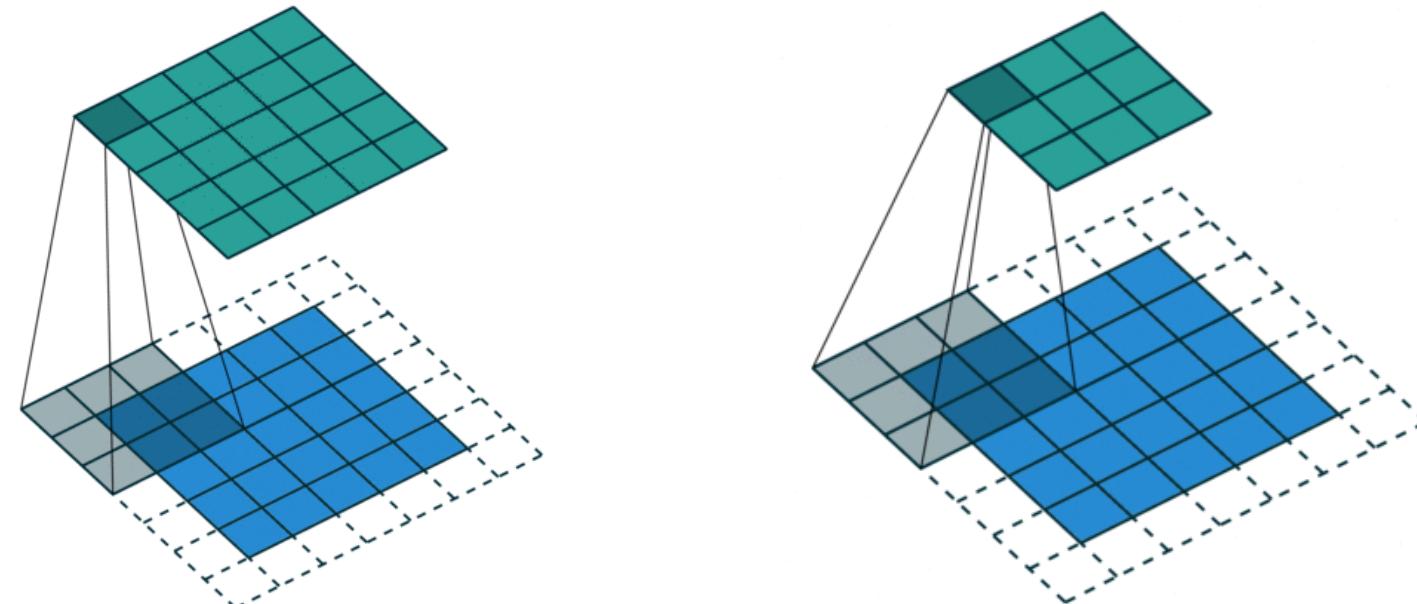


Convolutional neural network

Dimensiones del filtro

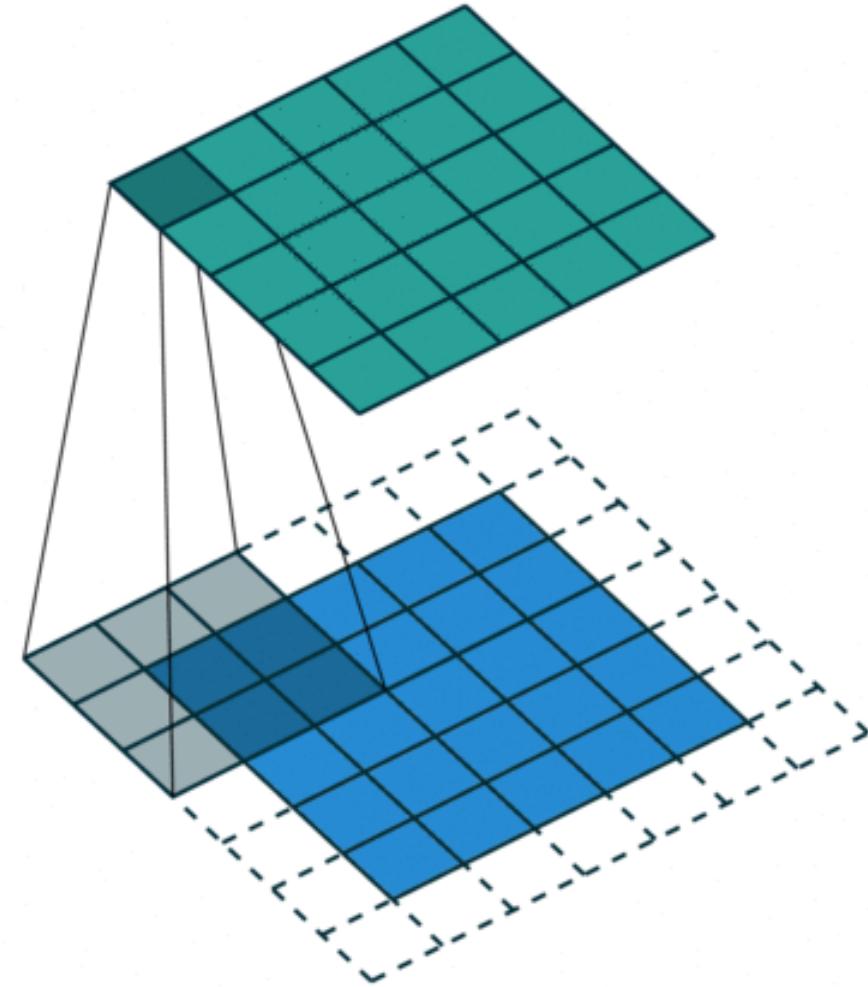


Stride



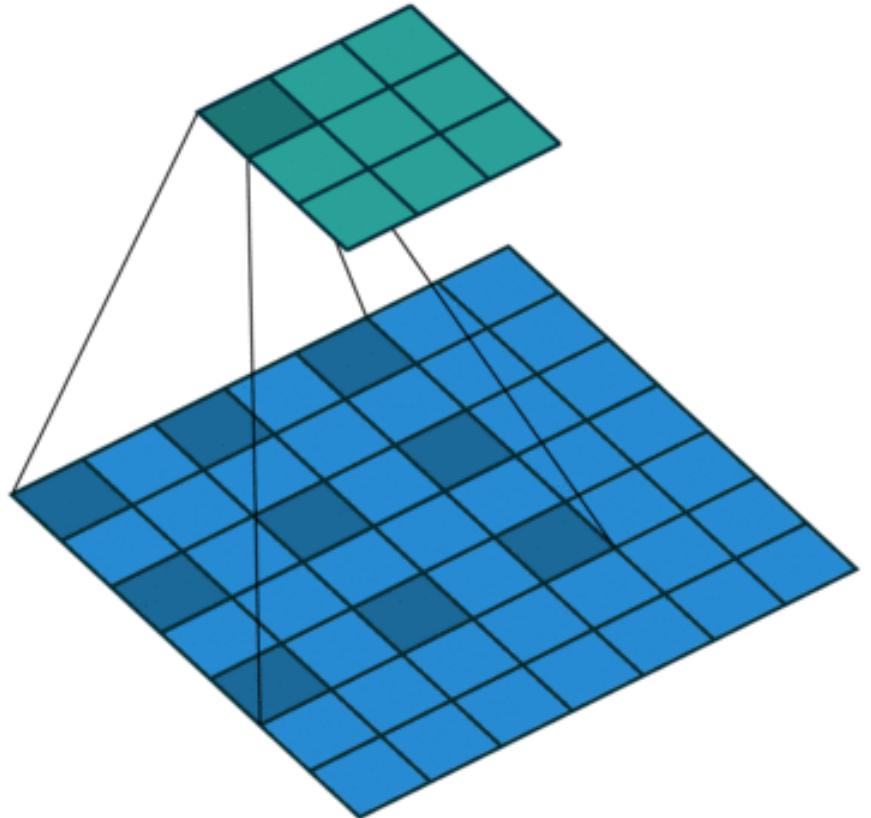
Convolutional *neural network*

Padding

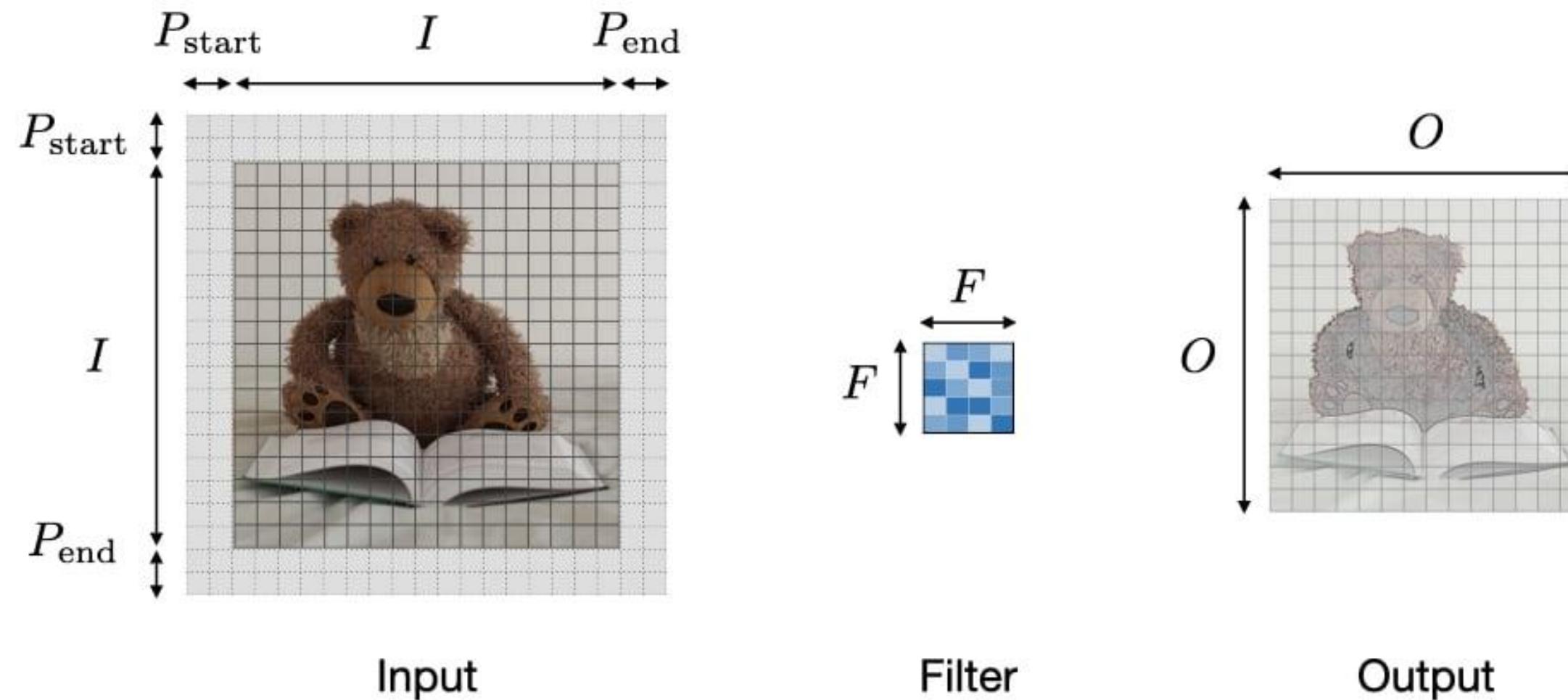


Convolutional *neural network*

Dilated Convolution



Convolutional neural network

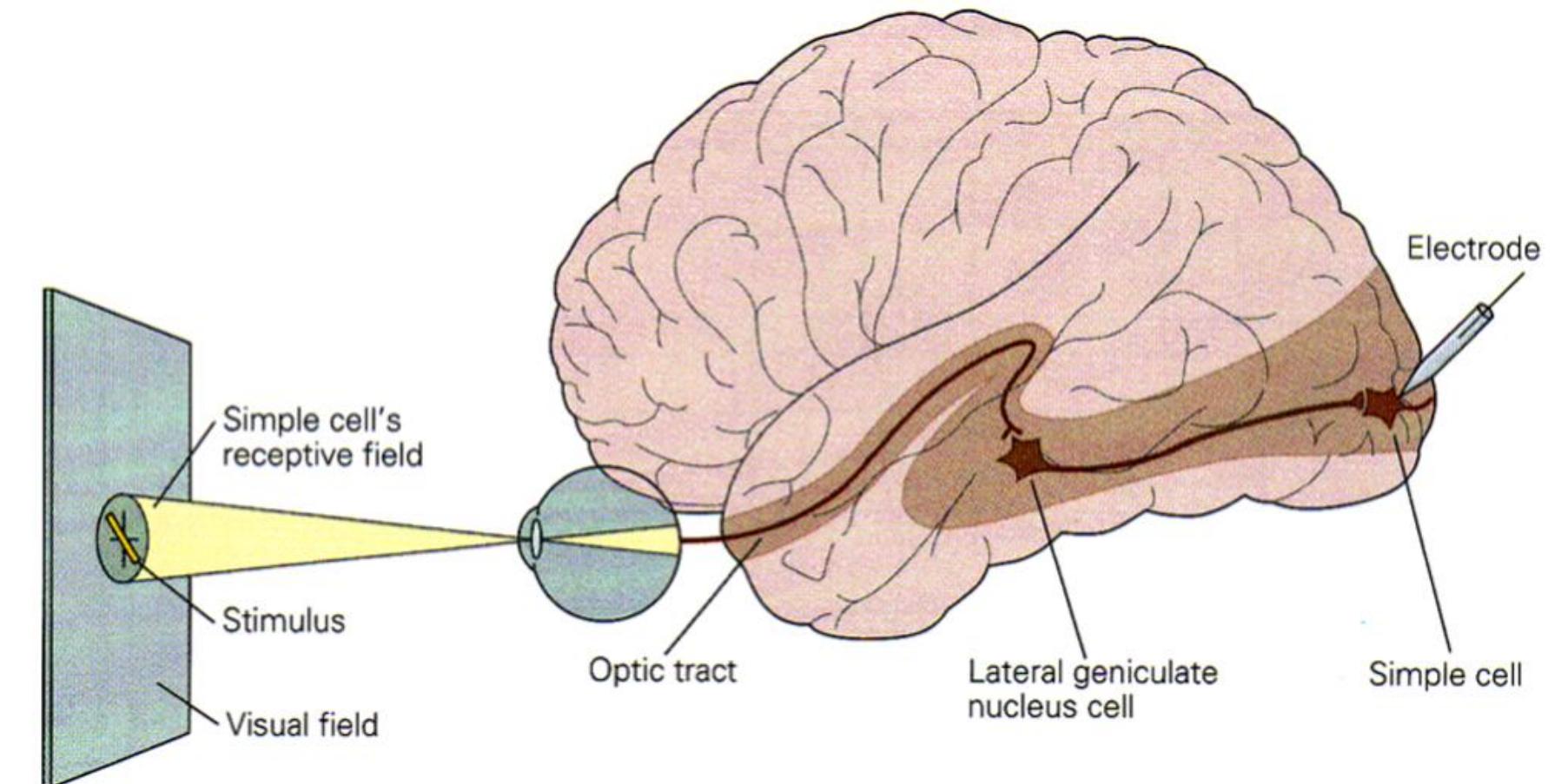
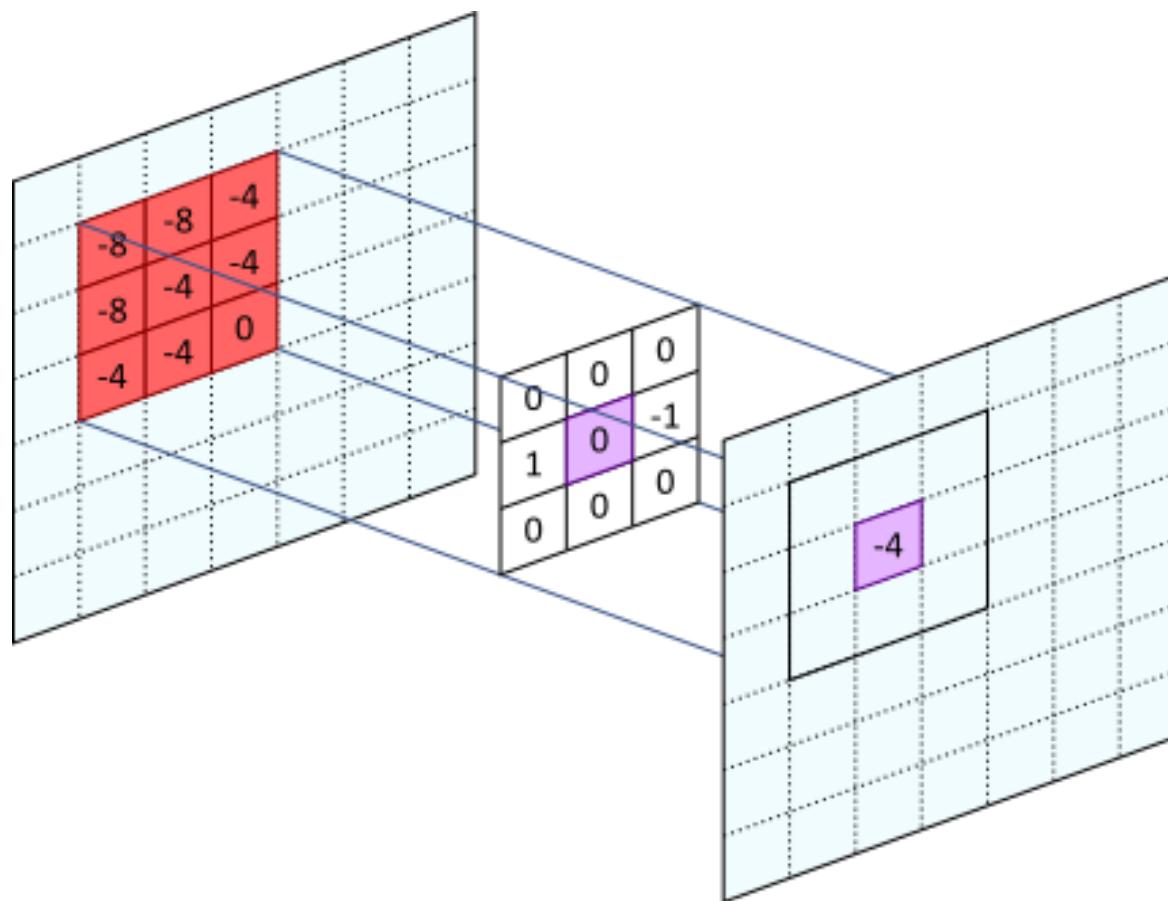


$$O = \frac{I - F + P_{start} + P_{end}}{S} + 1$$



Convolutional neural network

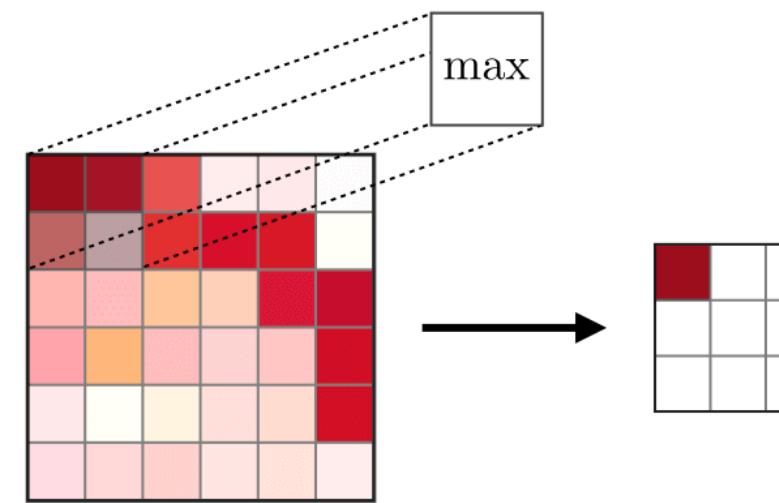
Receptive field



Convolutional neural network

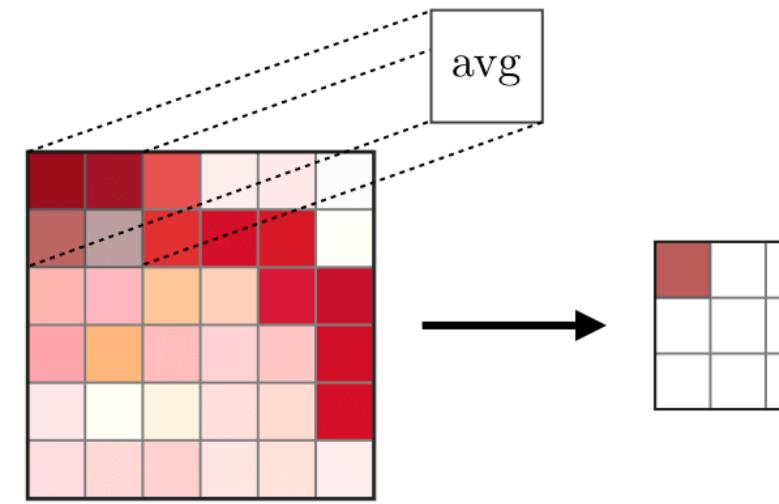
Pooling

Max Pooling



- Conserva las características detectadas.
- Más comúnmente utilizado.

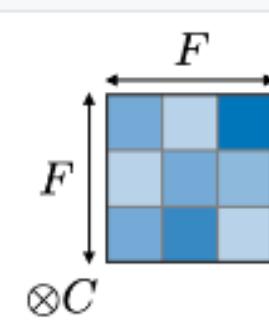
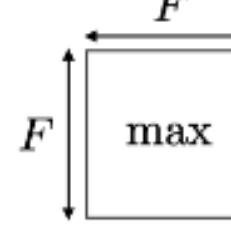
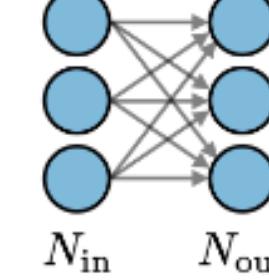
Average Pooling



- Conserva los valores promedio de las características.
- Utilizado en LeNet.

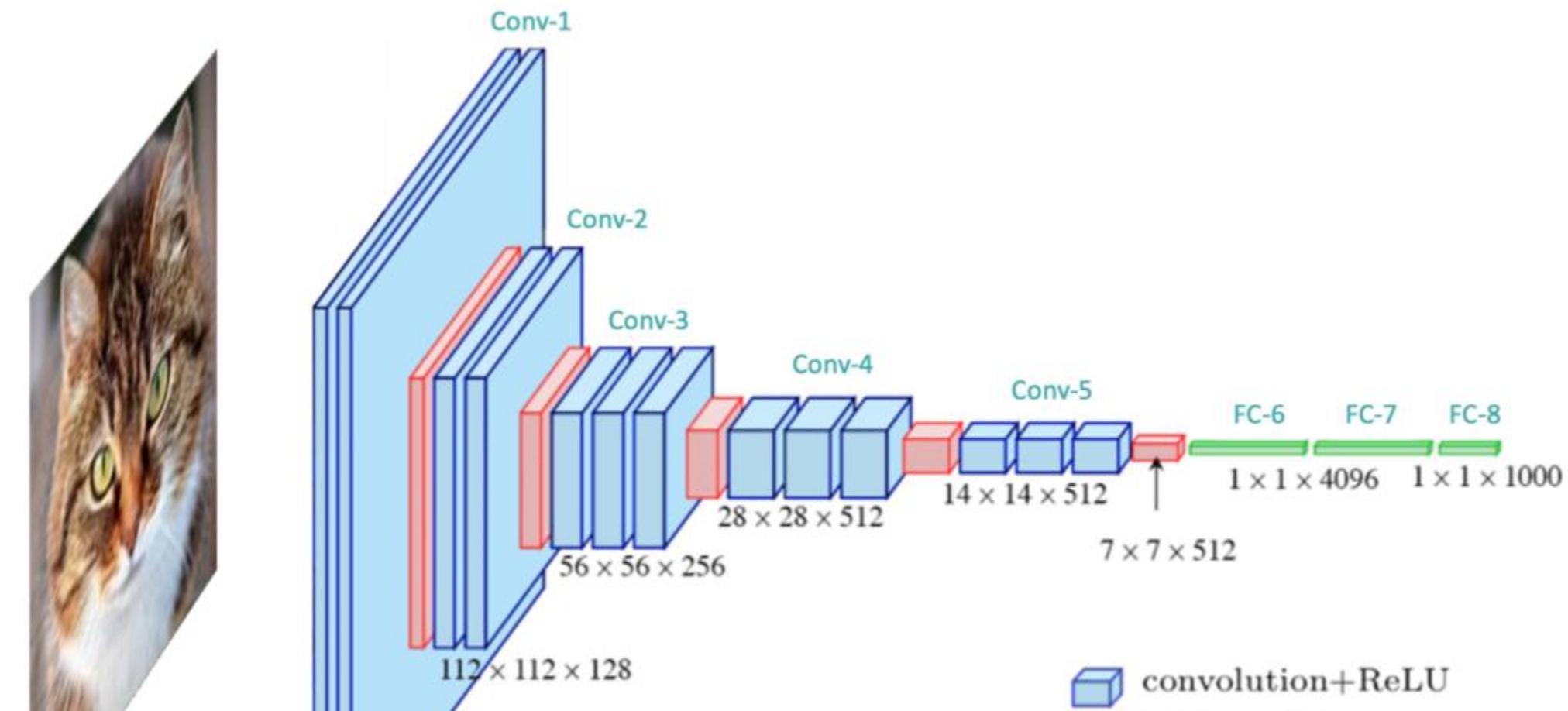


Convolutional neural network

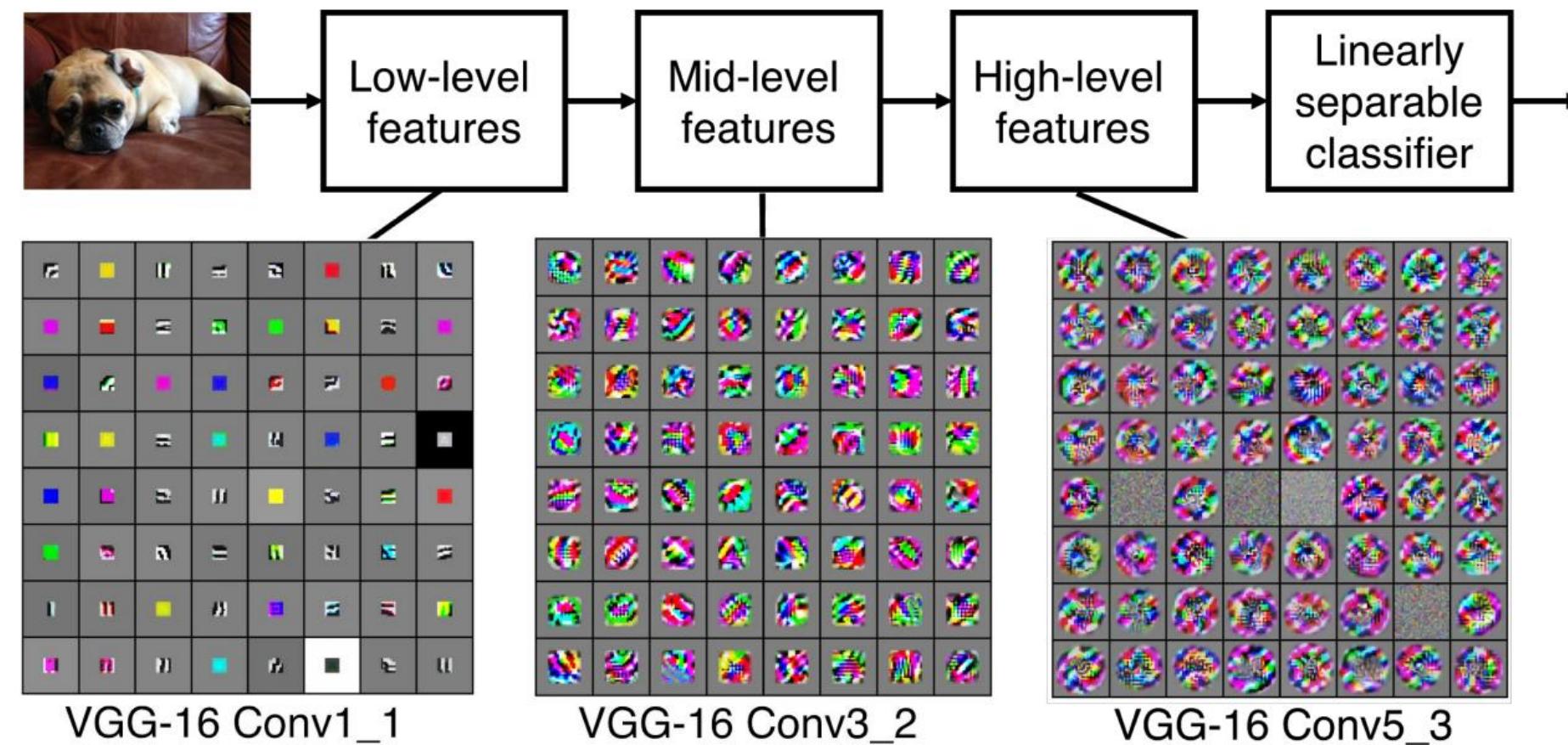
	CONV	POOL	FC
Illustration	 $F \times F \times C \times K$	 $F \times \text{max}$	 $N_{\text{in}} \times N_{\text{out}}$
Input size	$I \times I \times C$	$I \times I \times C$	N_{in}
Output size	$O \times O \times K$	$O \times O \times C$	N_{out}
Number of parameters	$(F \times F \times C + 1) \cdot K$	0	$(N_{\text{in}} + 1) \times N_{\text{out}}$
Remarks	<ul style="list-style-type: none"> One bias parameter per filter In most cases, $S < F$ A common choice for K is $2C$ 	<ul style="list-style-type: none"> Pooling operation done channel-wise In most cases, $S = F$ 	<ul style="list-style-type: none"> Input is flattened One bias parameter per neuron The number of FC neurons is free of structural constraints



Convolutional neural network

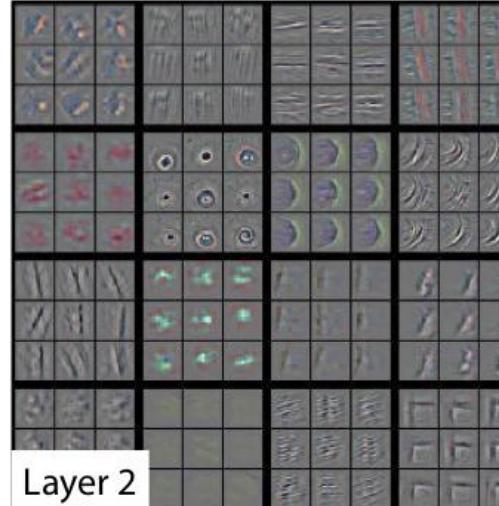


Convolutional neural network



Deep *representations*

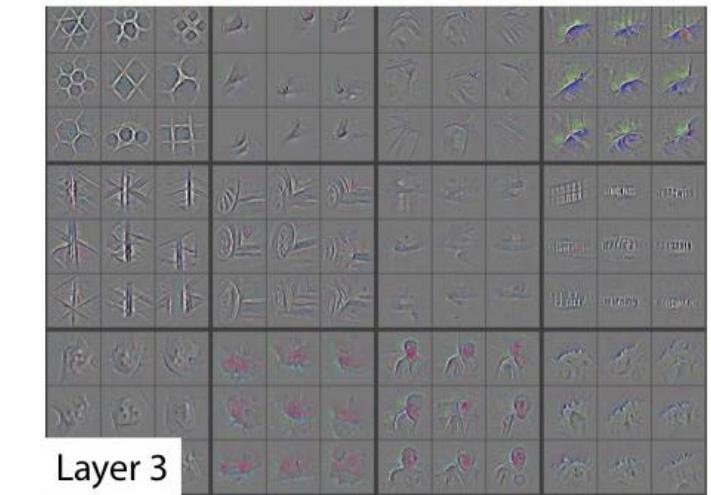
Layer 1



Layer 2

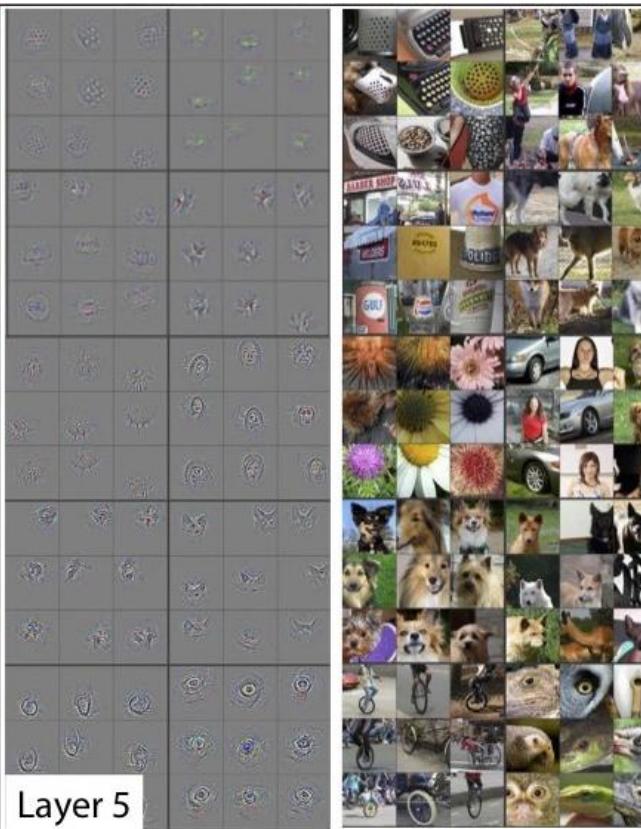


Layer 3

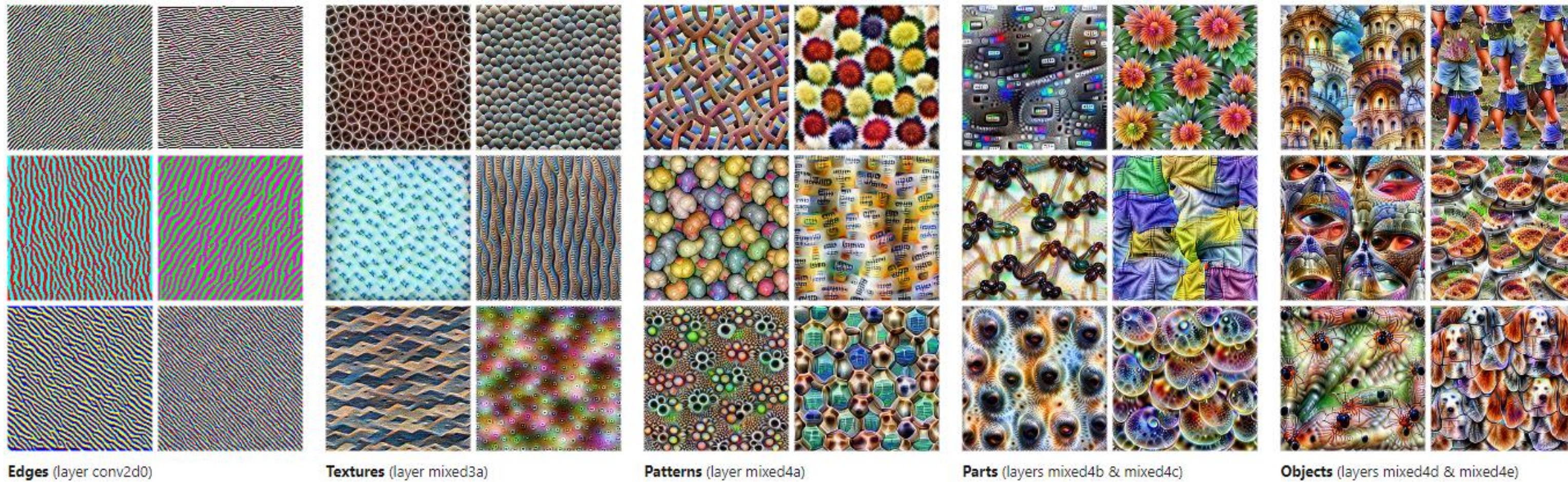


Layer 4

Layer 5



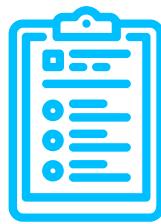
Deep *representations*



<https://distill.pub/2017/feature-visualization/>



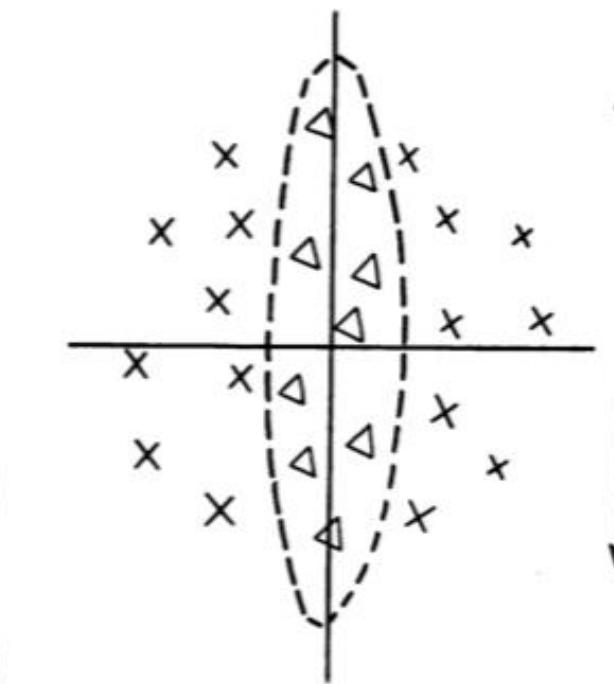
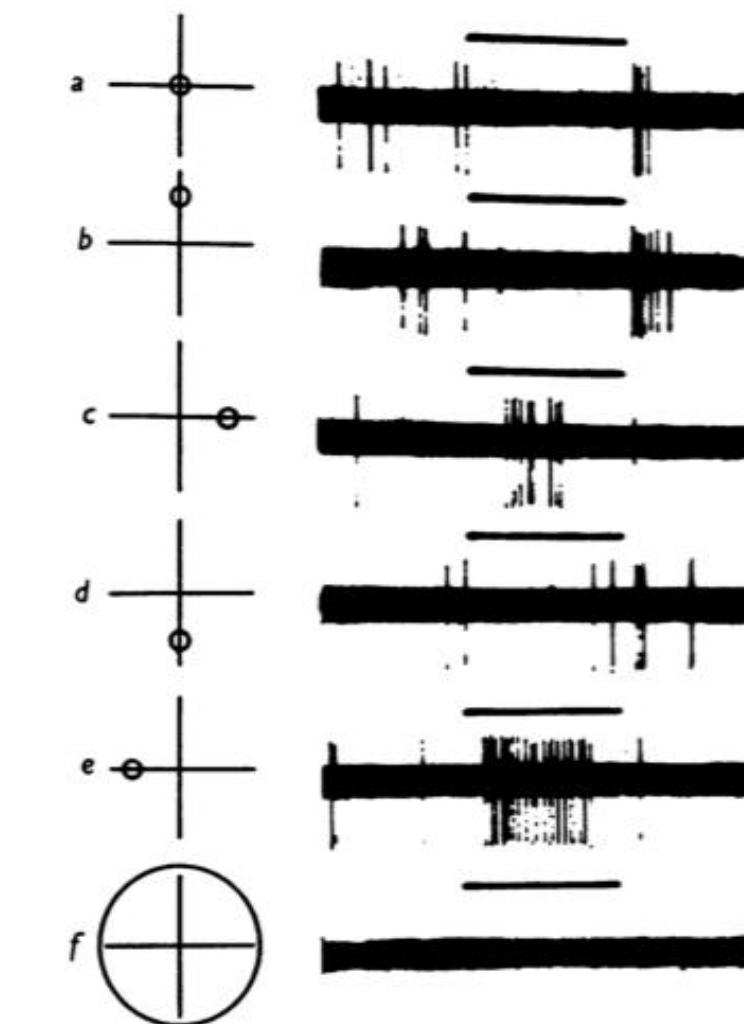
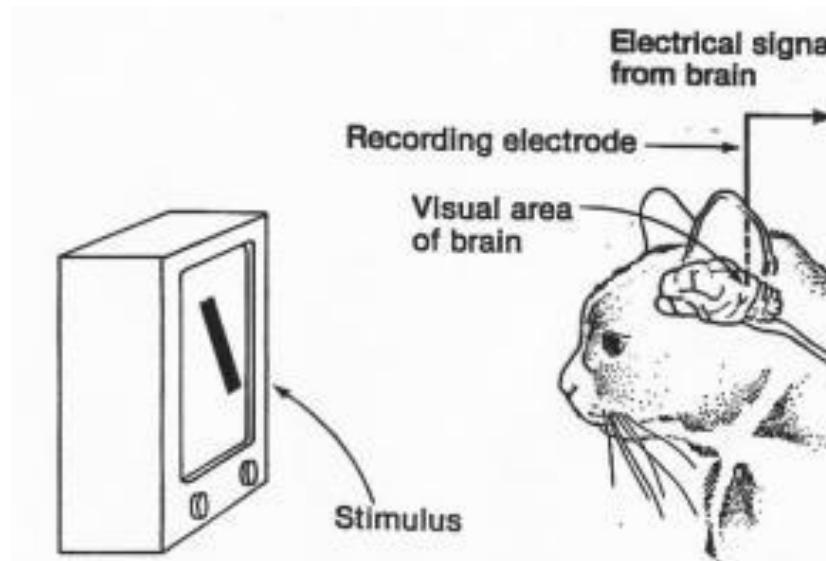
3.



Arquitecturas *clásicas*



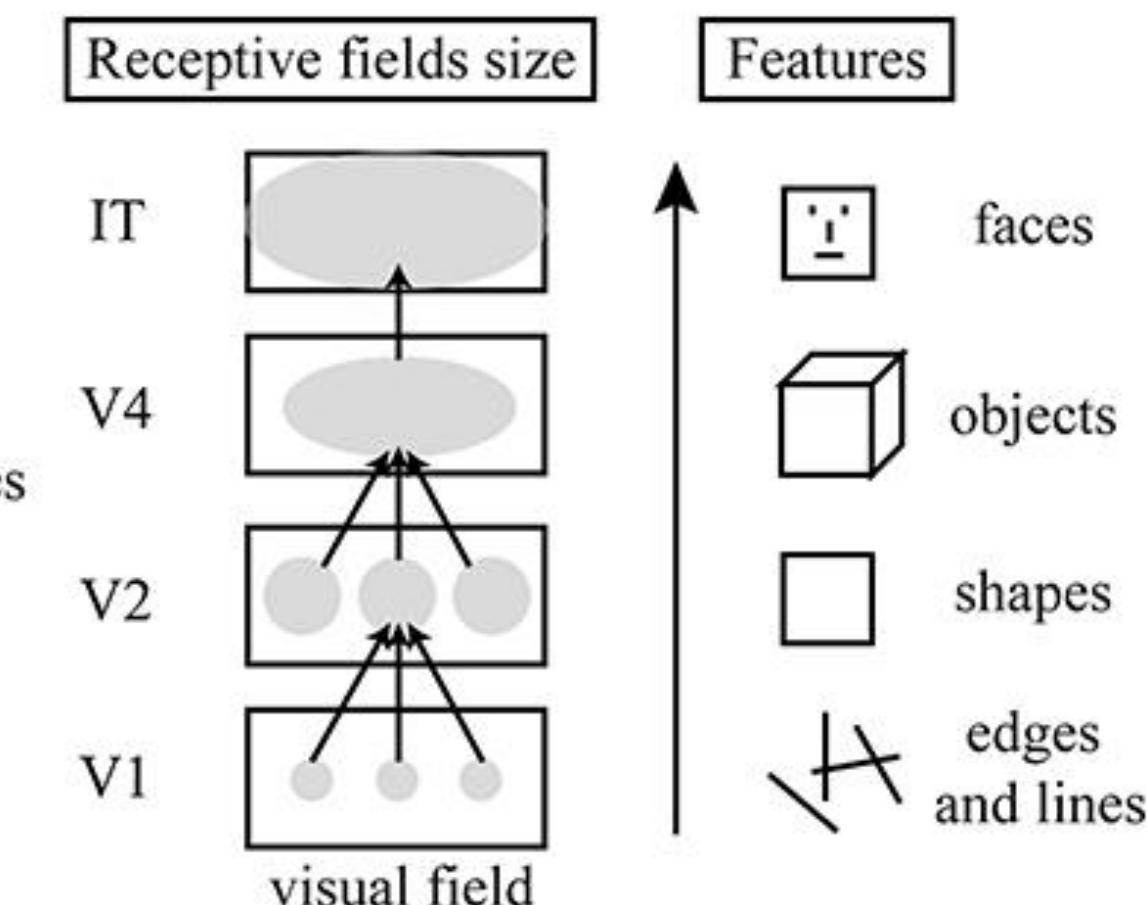
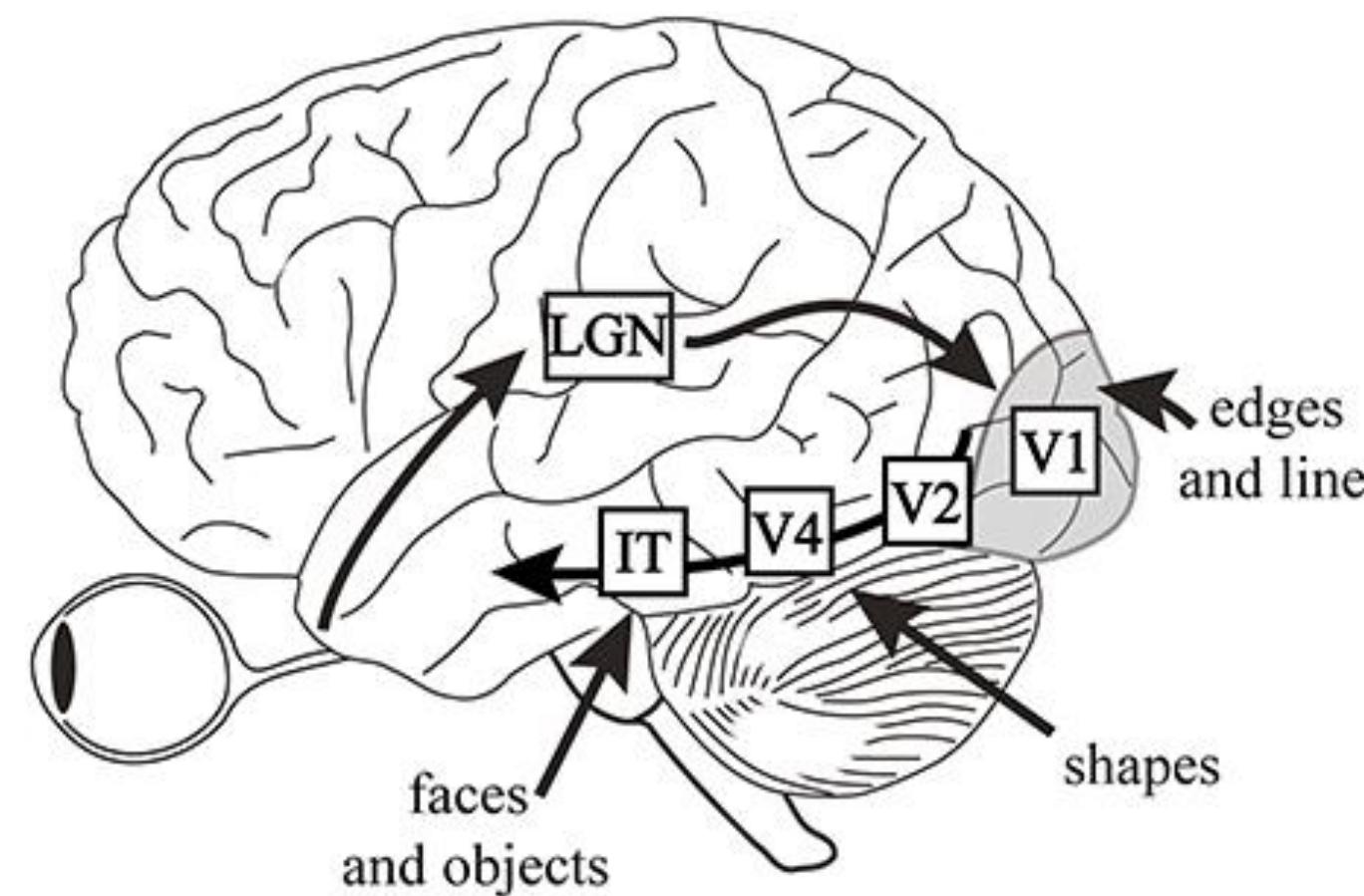
Corteza visual



TRANSFORMATEC

David Hubel et al. (1959) "Receptive fields of single neurones in the cat's striate cortex".
The Journal of physiology 148.3 (1959): 574-591.

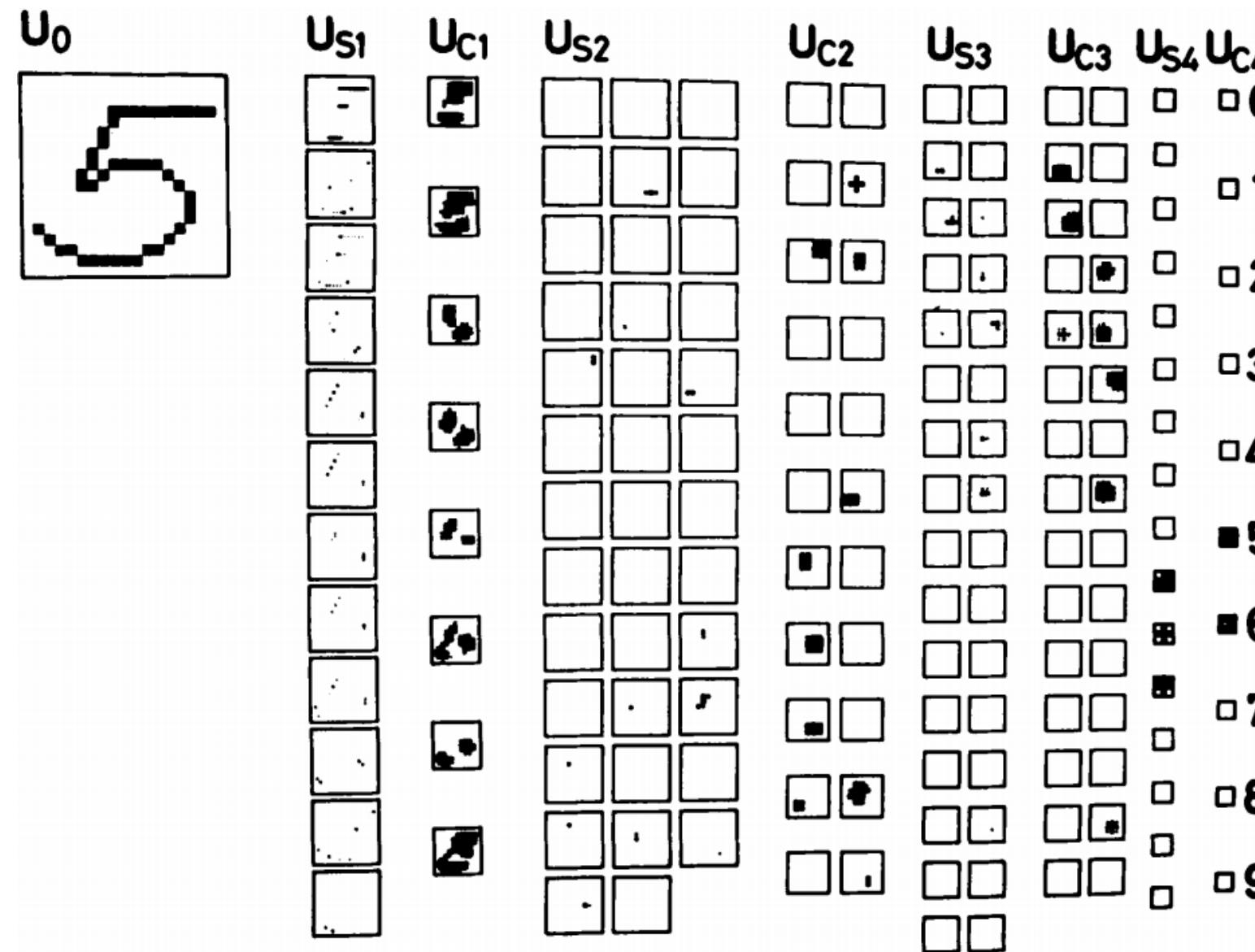
Corteza visual



TRANSFORMATEC

David Hubel et al. (1959) "Receptive fields of single neurones in the cat's striate cortex".
The Journal of physiology 148.3 (1959): 574-591.

Neocognitron



S-cell:

$$u_{Sl}(k_l, \mathbf{n}) = r_l \cdot \varphi \left(\frac{1 + \sum_{k_l=1}^{K_l-1} \sum_{v \in S_l} a_l(k_{l-1}, v, k_l) \cdot u_{Cl-1}(k_{l-1}, \mathbf{n} + v)}{1 + \frac{2r_l}{1+r_l} \cdot b_l(k_l) \cdot v_{Cl-1}(\mathbf{n})} \right)$$

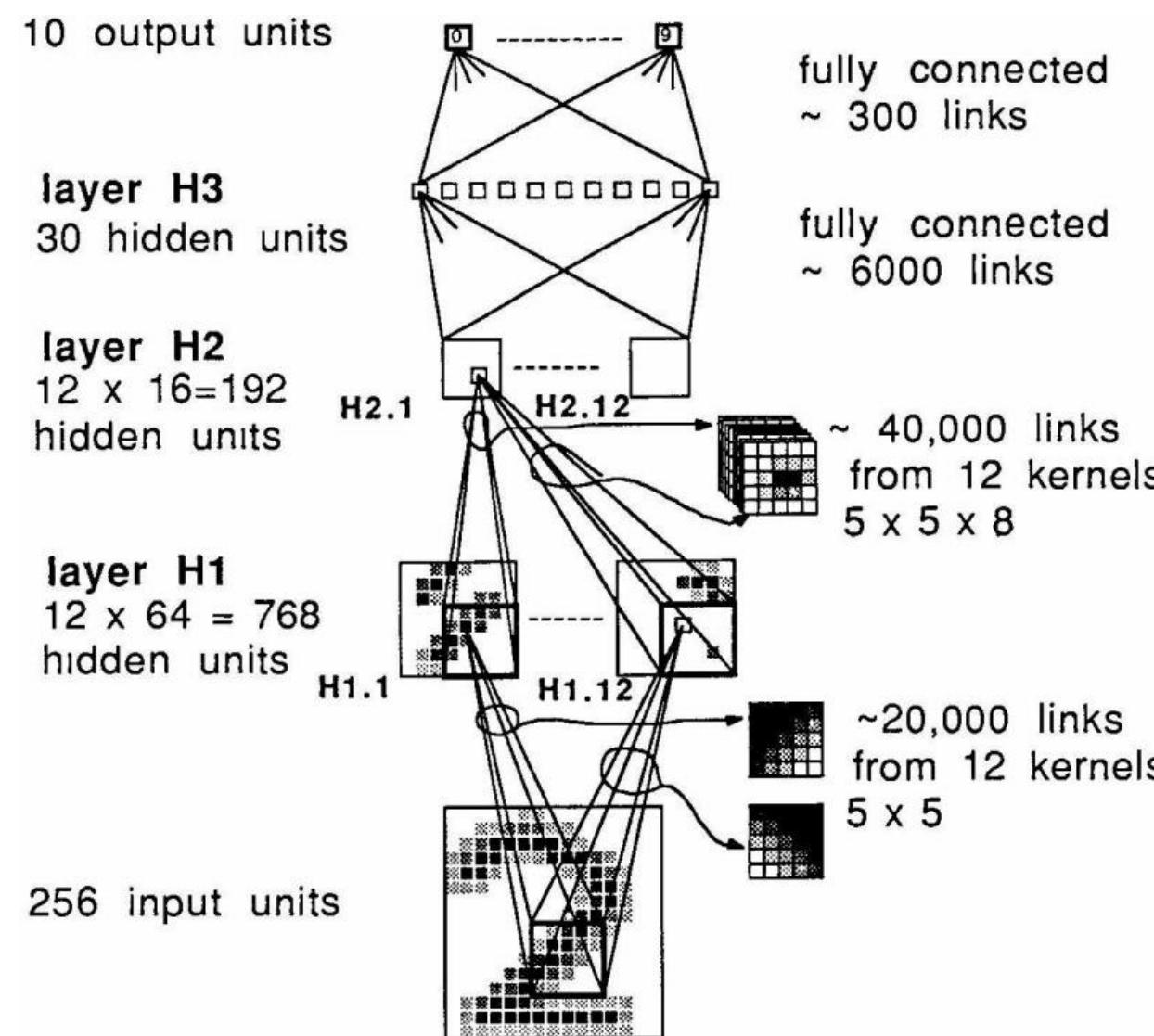
C-cell:

$$u_{Cl}(k_l, \mathbf{n}) = \psi \left(\frac{1 + \sum_{v \in D_l} d_l(v) \cdot u_{Sl}(k_l, \mathbf{n} + v)}{1 + v_{Sl}(\mathbf{n})} \right)$$

donde $\varphi(x) = \begin{cases} x, & x \geq 0 \\ 0, & x < 0 \end{cases}$ $\psi(x) = \varphi\left(\frac{x}{\alpha + x}\right)$

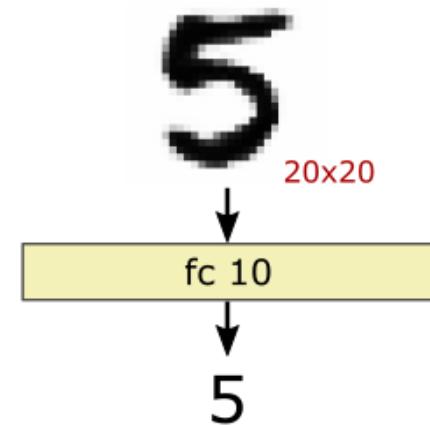


LeNet-1 (1989)



LeNet

Perceptrón

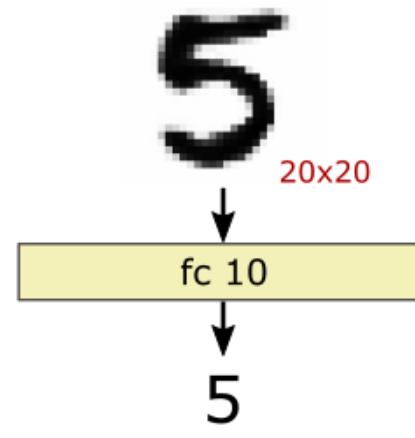


Error:
8.4%

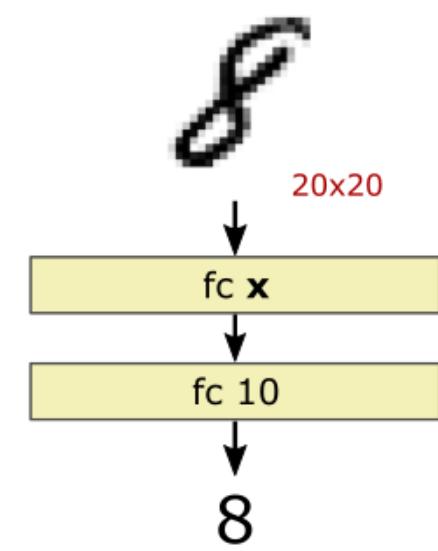


LeNet

Perceptrón



MLP



Error:
8.4%

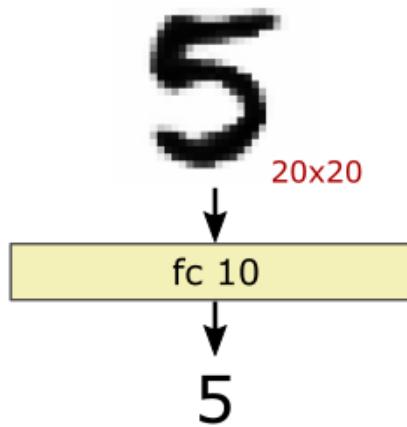
Error:
3.6% (x:300)
3.8% (x:1000)



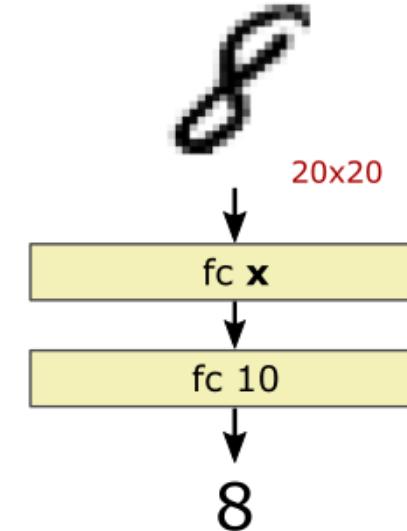
Yan LeCun et al. (1995) "Comparison of learning algorithms for handwritten digit recognition".
International conference on artificial neural networks.

LeNet

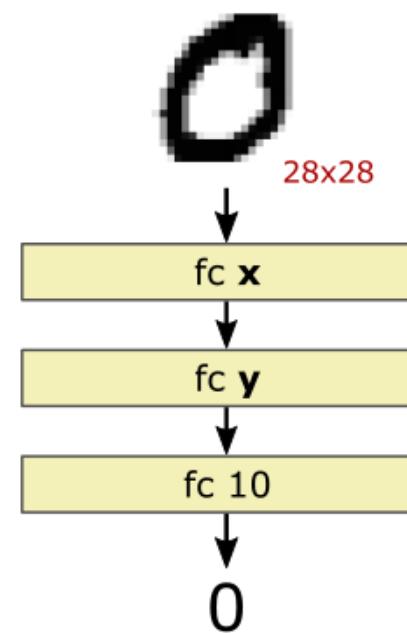
Perceptrón



MLP



MLP



Error:
8.4%

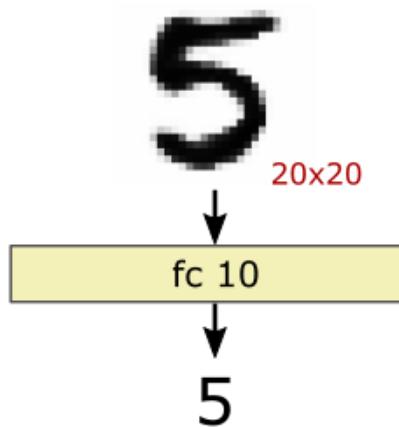
Error:
3.6% (x:300)
3.8% (x:1000)

Error:
3.05% (x:300, y:100)
2.95% (x:1000,y:150)

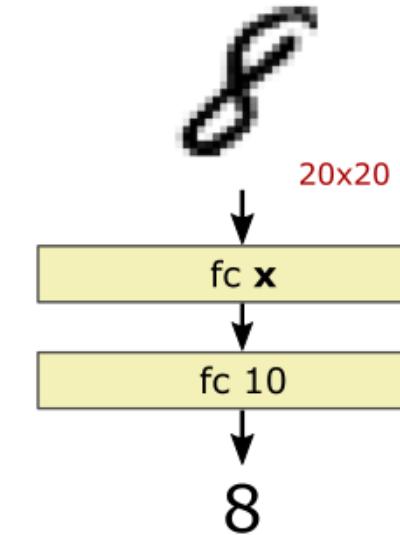


LeNet

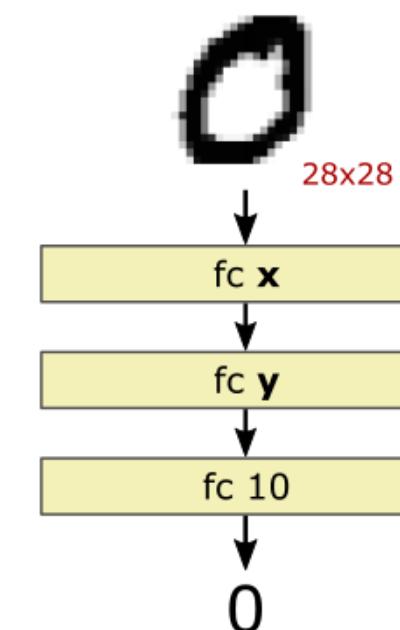
Perceptrón



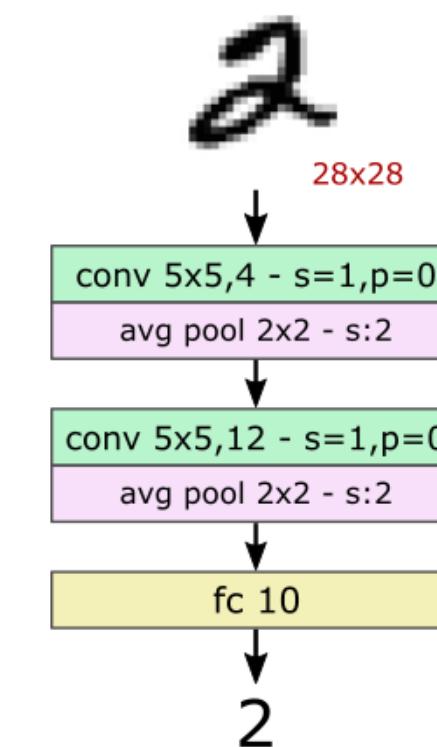
MLP



MLP



LeNet-1



Error:
8.4%

Error:
3.6% (x:300)
3.8% (x:1000)

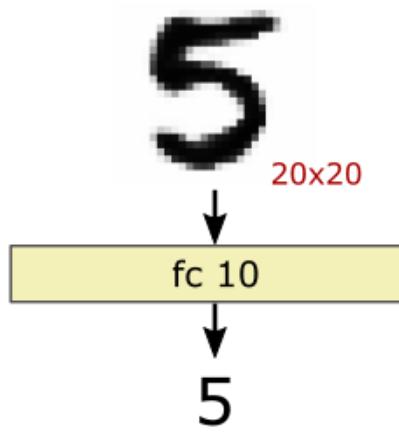
Error:
3.05% (x:300, y:100)
2.95% (x:1000,y:150)

Error:
1.7%



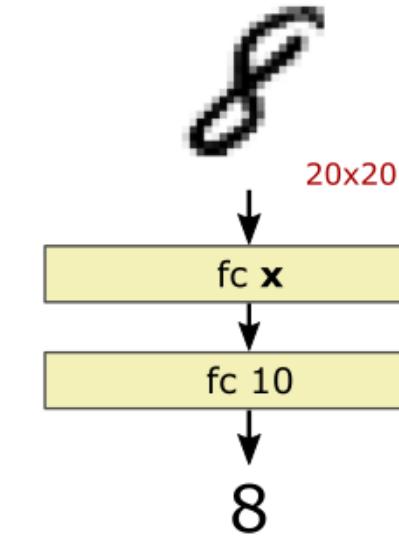
LeNet

Perceptrón



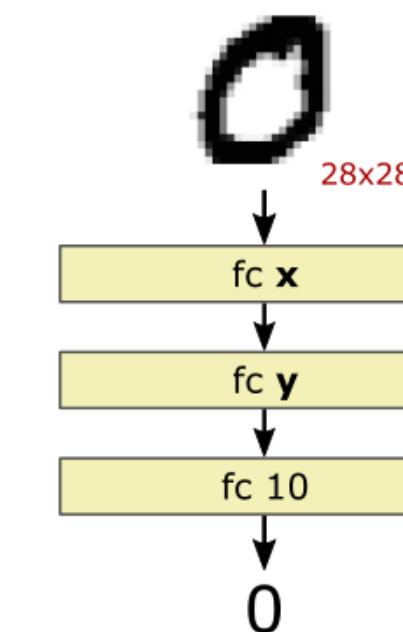
Error:
8.4%

MLP



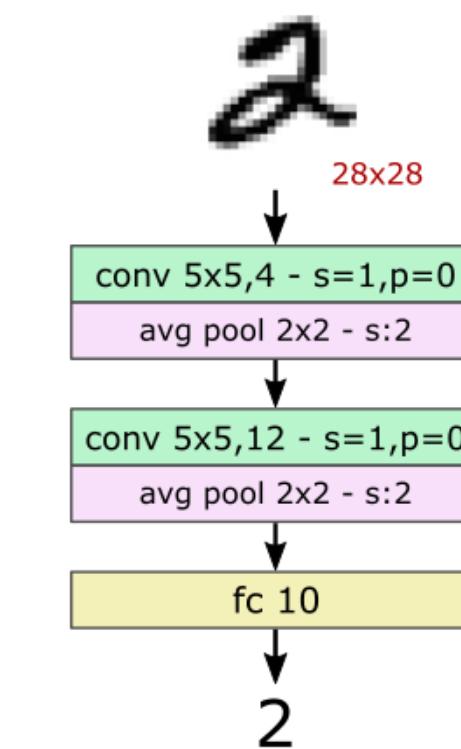
Error:
3.6% (x:300)
3.8% (x:1000)

MLP



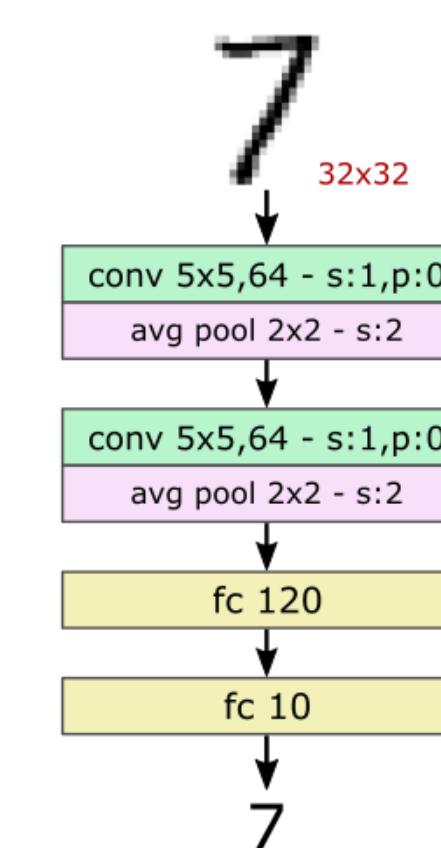
Error:
3.05% (x:300, y:100)
2.95% (x:1000,y:150)

LeNet-1



Error:
1.7%

LeNet-4

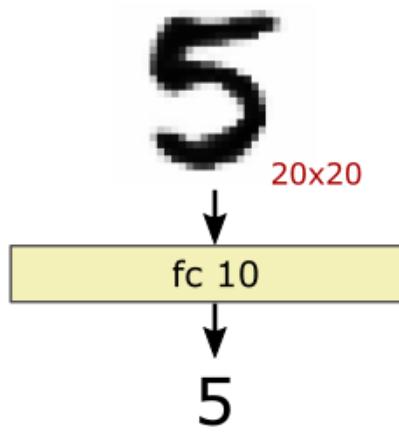


Error:
1.1%



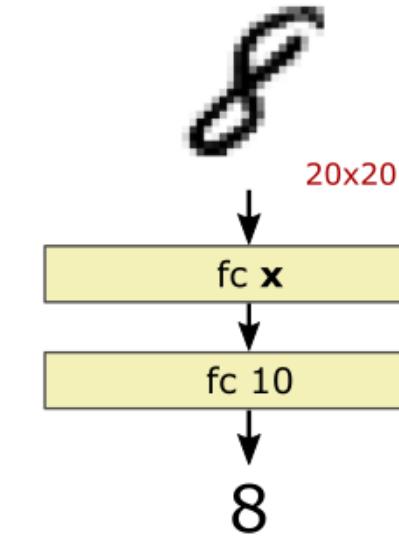
LeNet

Perceptrón



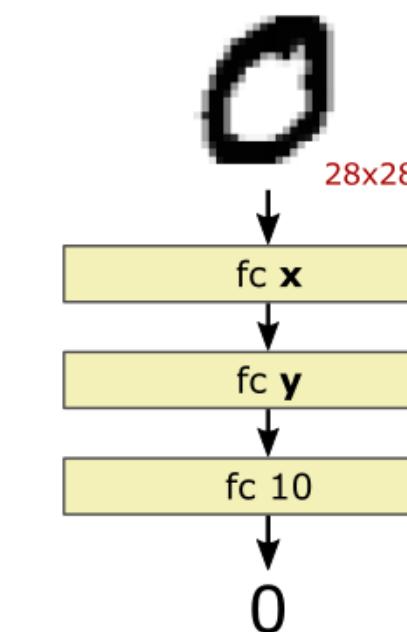
Error:
8.4%

MLP



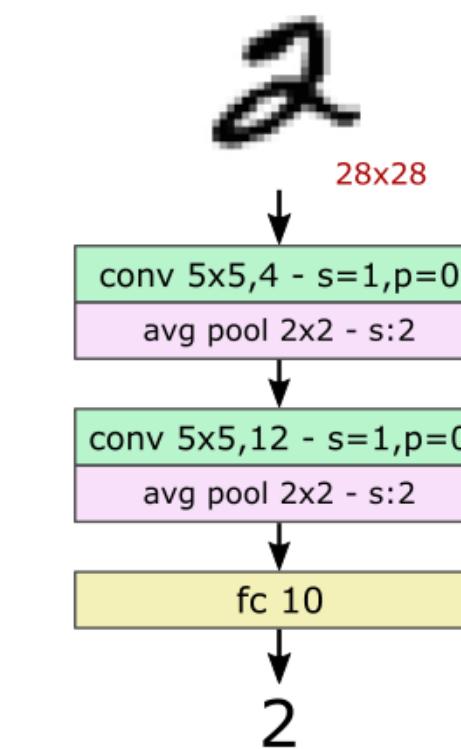
Error:
3.6% (x:300)
3.8% (x:1000)

MLP



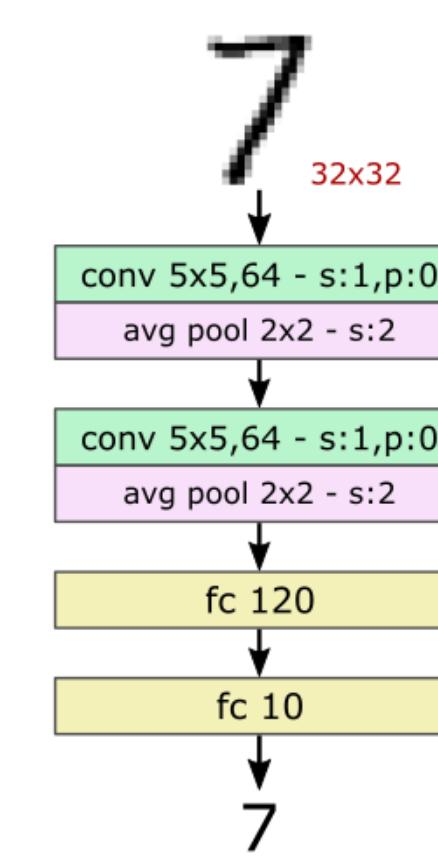
Error:
3.05% (x:300, y:100)
2.95% (x:1000,y:150)

LeNet-1



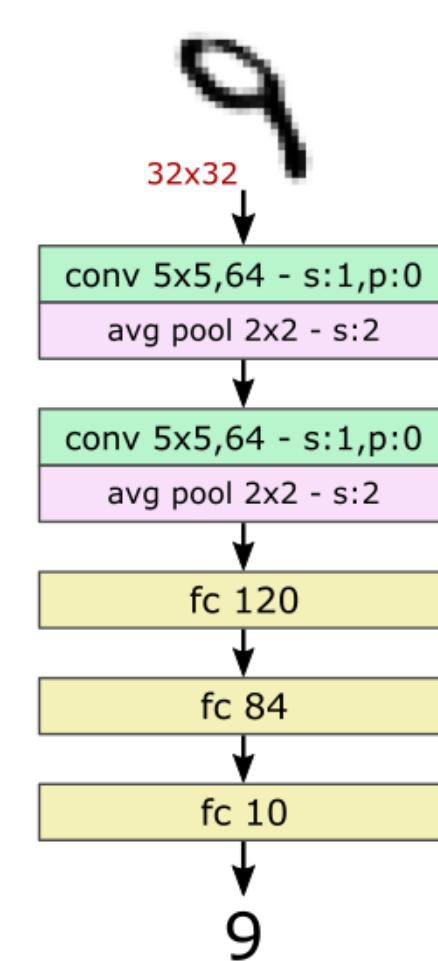
Error:
1.7%

LeNet-4



Error:
1.1%

LeNet-5

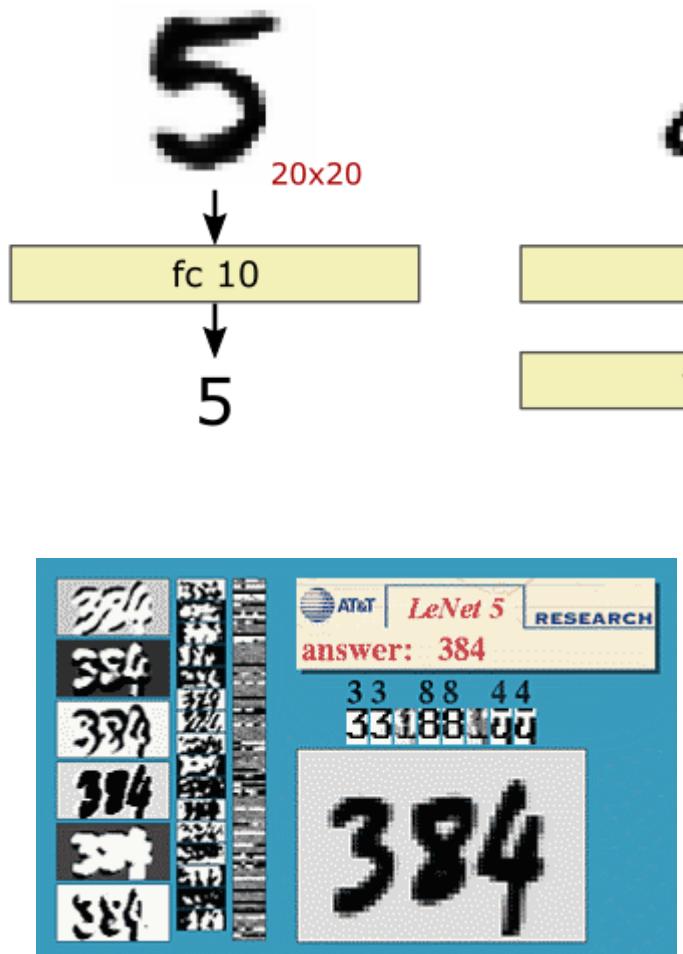


Error:
0.95%

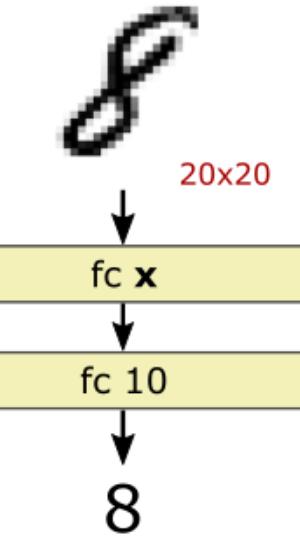


LeNet

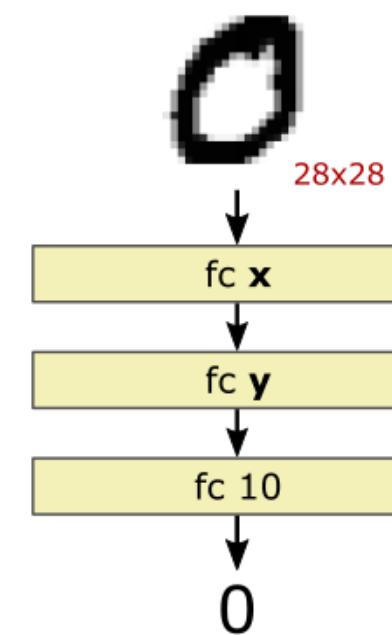
Perceptrón



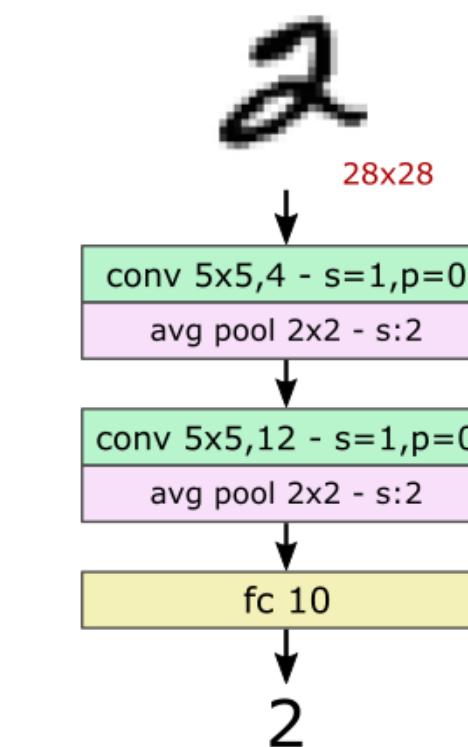
MLP



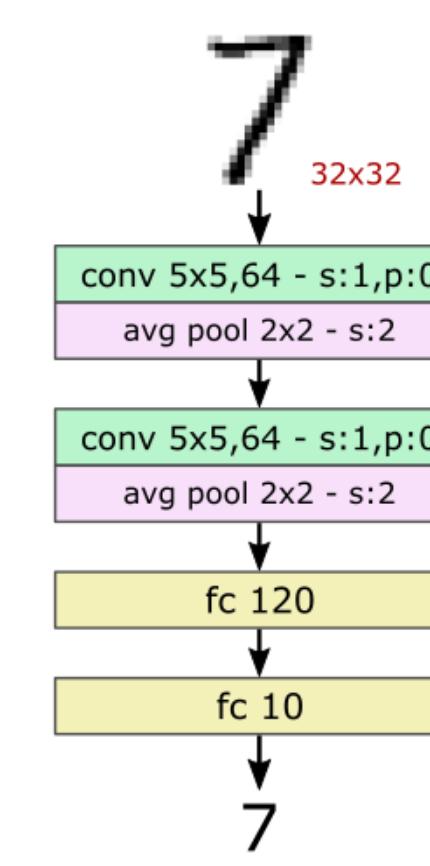
MLP



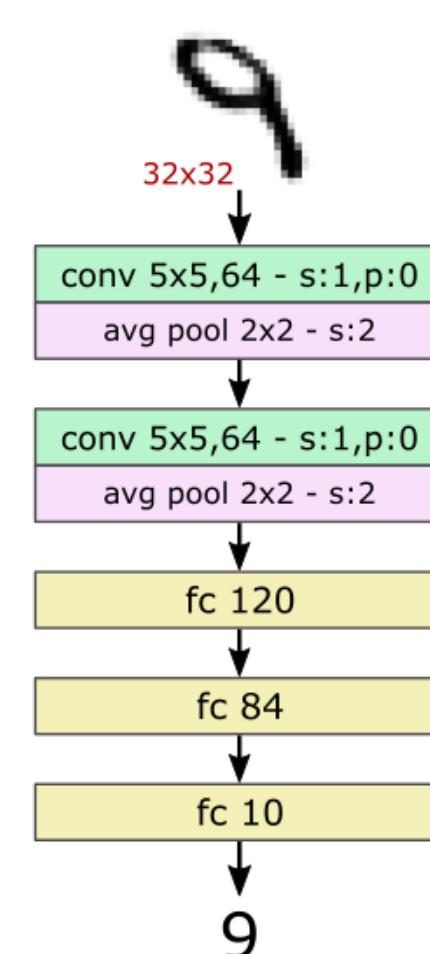
LeNet-1



LeNet-4



LeNet-5



Error:
8.4%

Error:
3.6% (x:300)
3.8% (x:1000)

Error:
3.05% (x:300, y:100)
2.95% (x:1000,y:150)

Error:
1.7%

Error:
1.1%

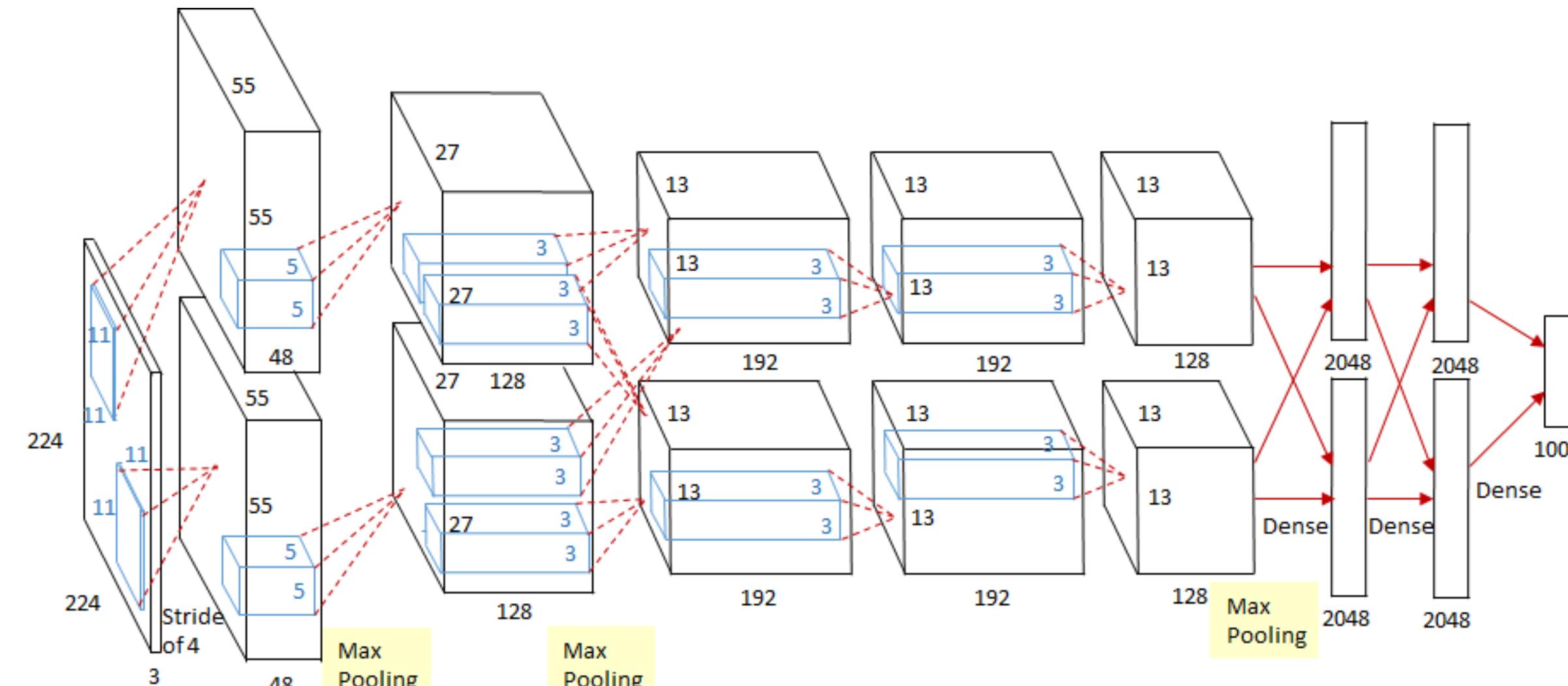
Error:
0.95%



TRANSFORMATEC

Yan LeCun et al. (1995) "Comparison of learning algorithms for handwritten digit recognition".
International conference on artificial neural networks.

AlexNet



TRANSFORMATEC

Alex Krizhevsky et al. (2012) "ImageNet Classification with Deep Convolutional Neural Networks".
Advances in neural information processing systems.

AlexNet



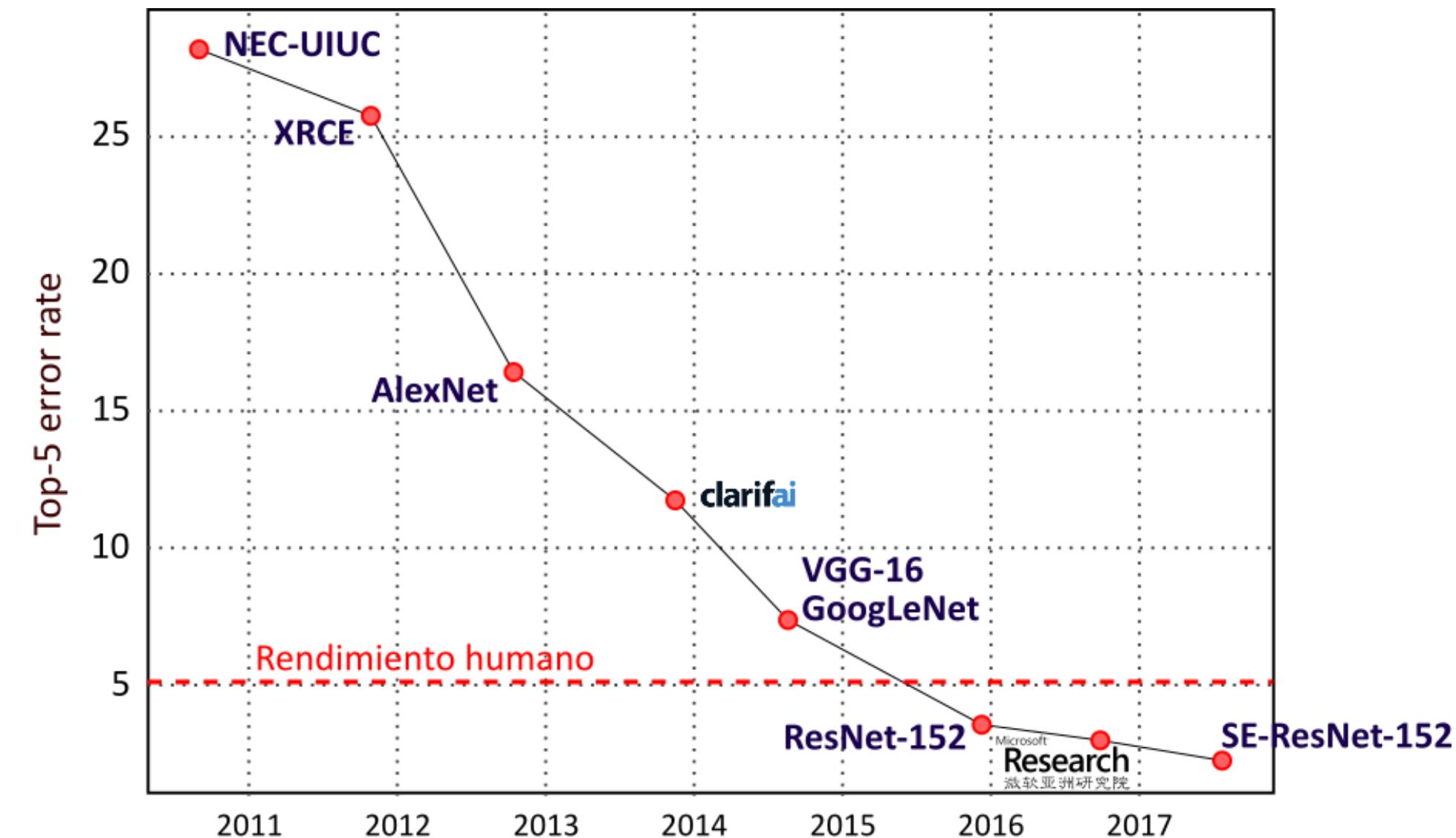
- 14 197 122 imágenes
- 21 841 categorías
- 1 280 000 imágenes de entrenamiento
- 50 000 imágenes de evaluación
- 10 000 imágenes de prueba
- 1000 categorías



TRANSFORMATEC

Jia Deng et al. (2009) "Imagenet: A large-scale hierarchical image database".
2009 IEEE Conference on Computer Vision and Pattern Recognition

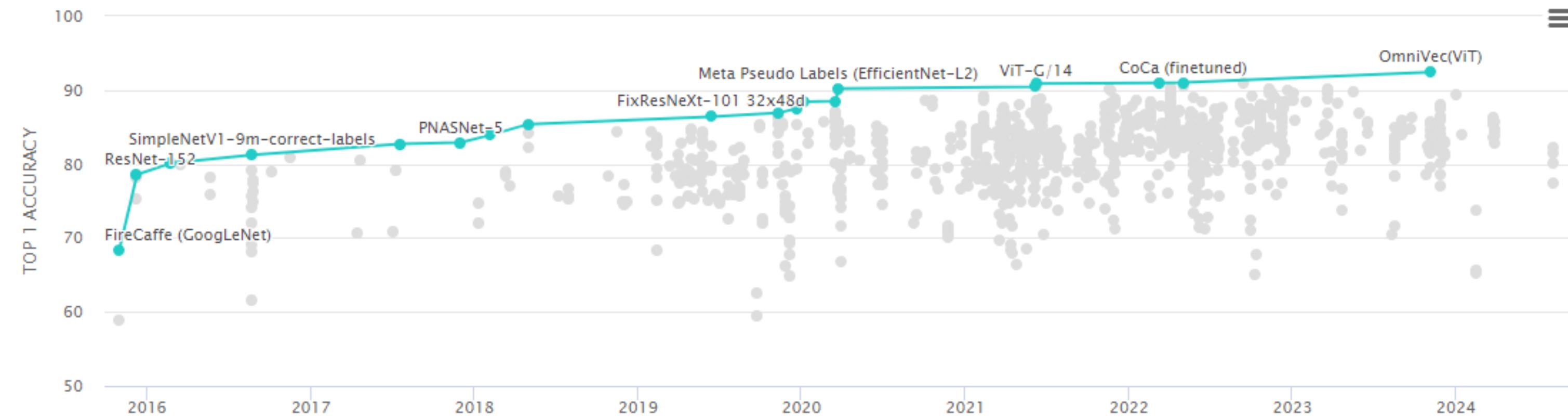
AlexNet



TRANSFORMATEC

Jia Deng et al. (2009) "Imagenet: A large-scale hierarchical image database".
2009 IEEE Conference on Computer Vision and Pattern Recognition.

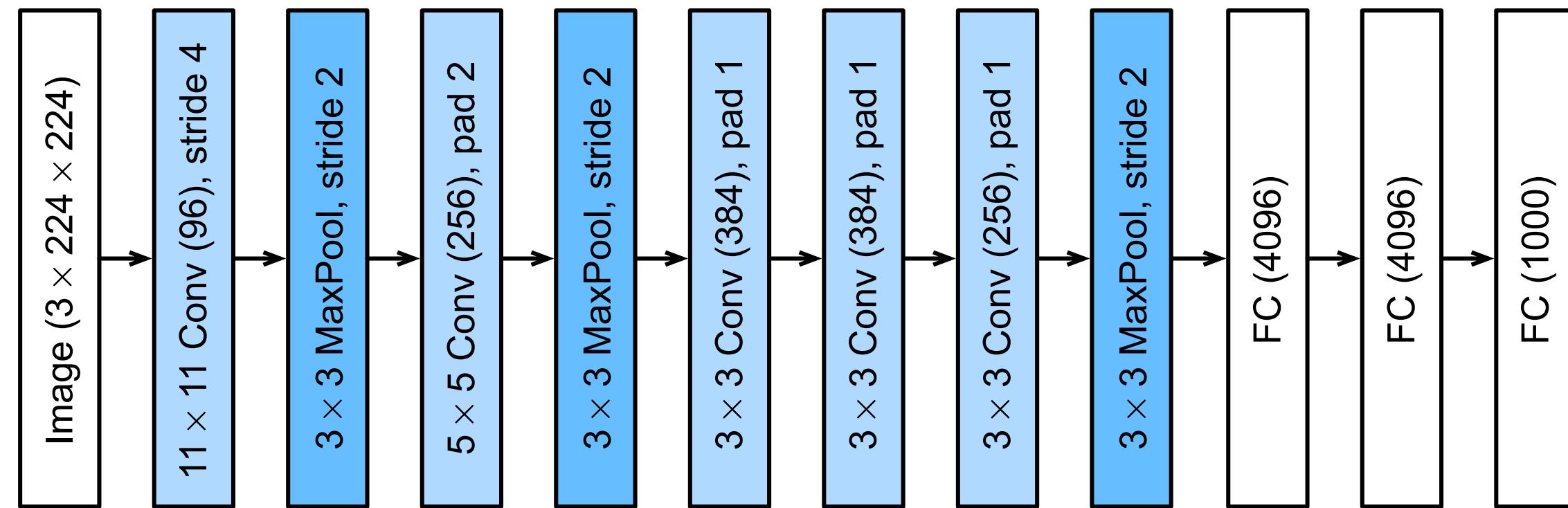
AlexNet



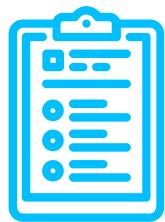
TRANSFORMATEC

Jia Deng et al. (2009) "Imagenet: A large-scale hierarchical image database".
2009 IEEE Conference on Computer Vision and Pattern Recognition.

AlexNet



4.



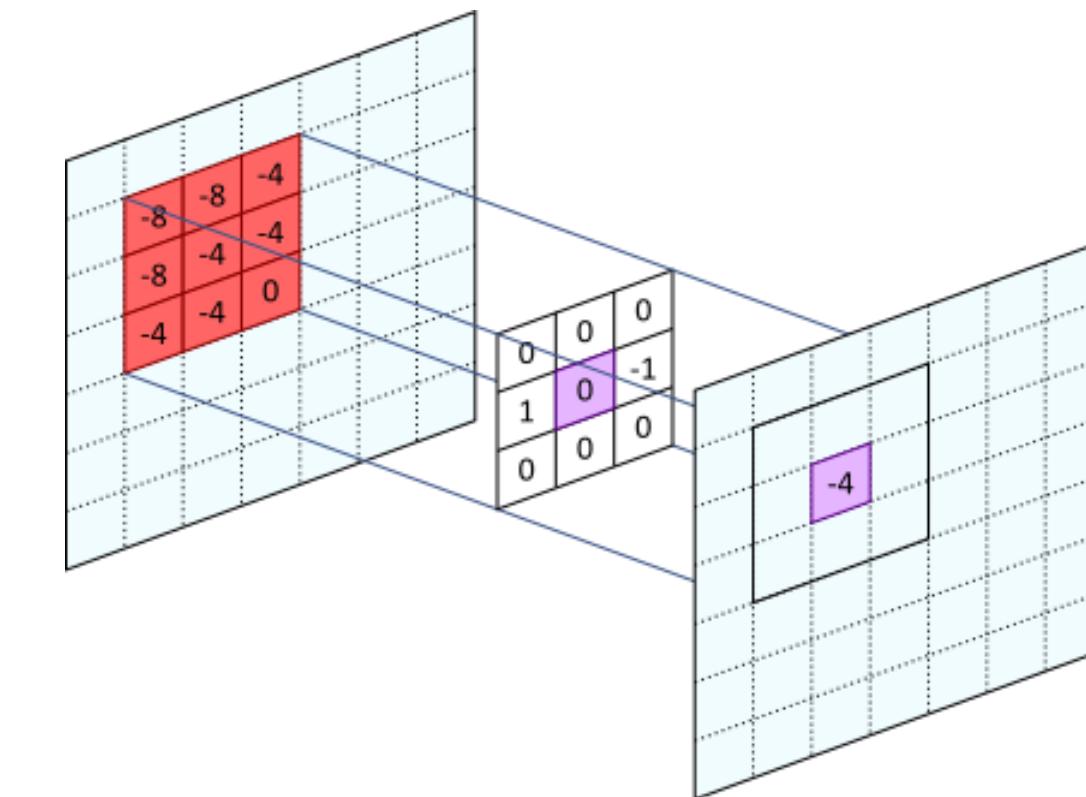
vGG

TRANSFORMATEC

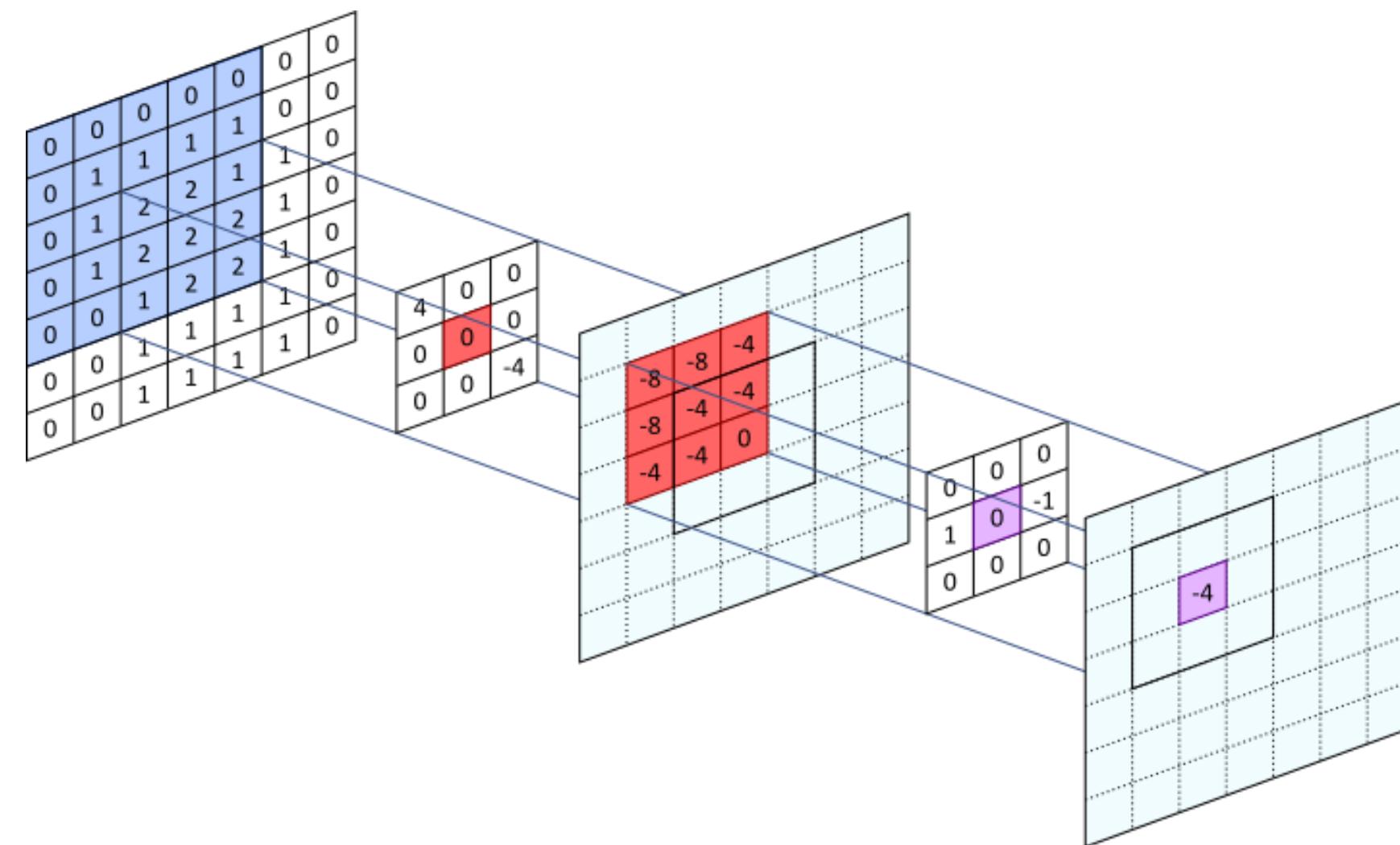
> Reinventa el mundo <



VGG

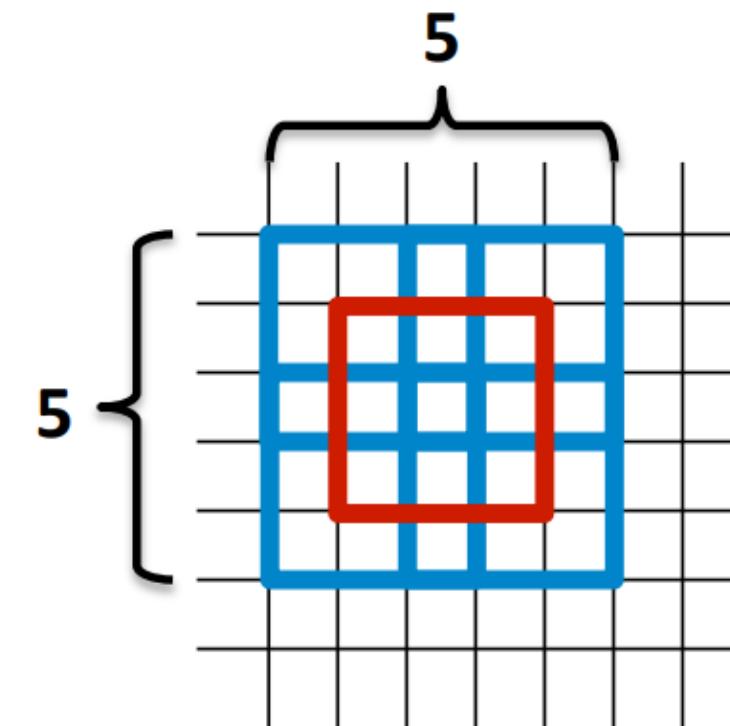


VGG



VGG

Stacks of 3x3 Filters



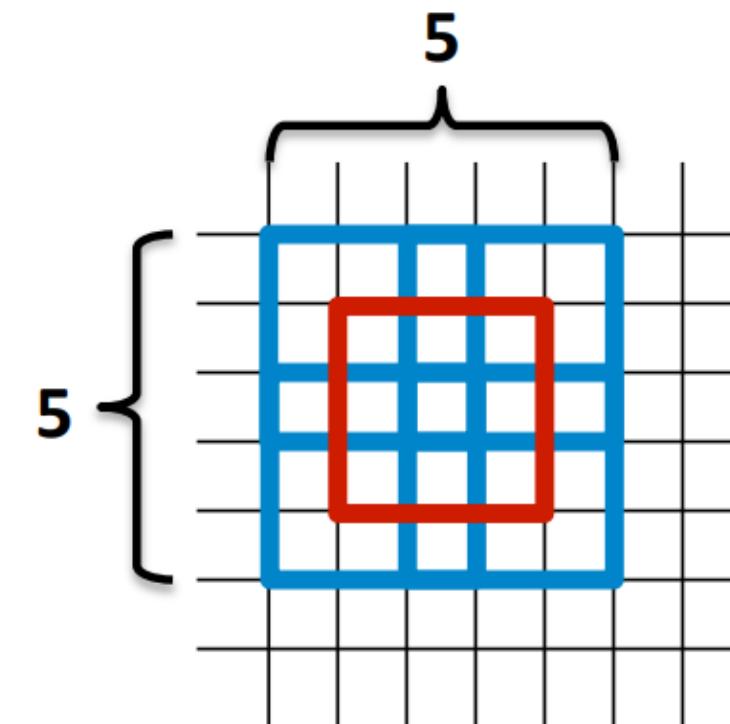
-  1st 3x3 conv. layer
-  2nd 3x3 conv. layer



1 filtro 5×5, num. de parámetros $5 \times 5 = 25$
 2 filtros 3×3, num. de parámetros $2 \times 3 \times 3 = 18$
 El número de parámetros se reduce en **28%**

VGG

Stacks of 3x3 Filters



-  1st 3x3 conv. layer
-  2nd 3x3 conv. layer

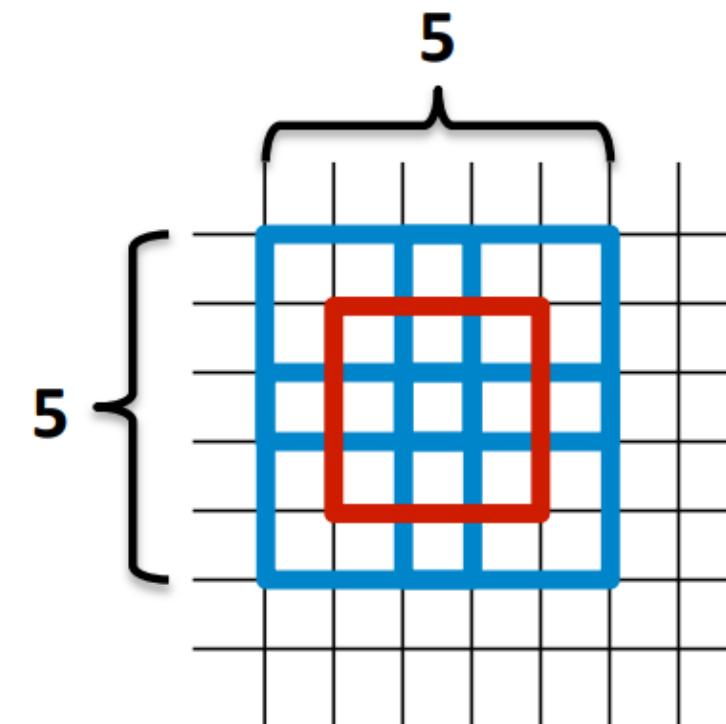


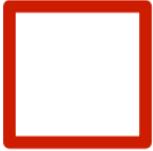
1 filtro 5×5, num. de parámetros $5 \times 5 = 25$
 2 filtros 3×3, num. de parámetros $2 \times 3 \times 3 = 18$
 El número de parámetros se reduce en **28%**

ZFNet 1 filtro 7×7, num. de parámetros: $7 \times 7 = 49$
 3 filtros 3×3, num. de parámetros: $3 \times 3 \times 3 = 27$
 El número de parámetros se reduce en **45%**

VGG

Stacks of 3x3 Filters



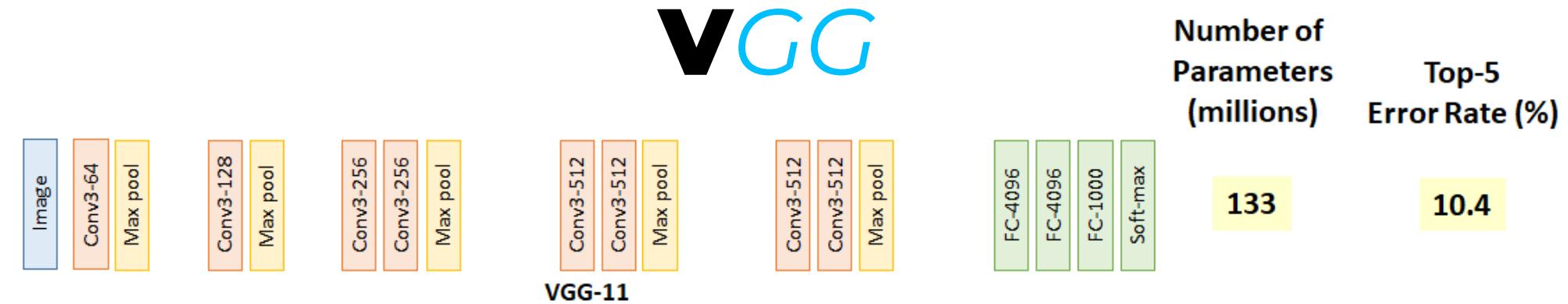
-  1st 3x3 conv. layer
-  2nd 3x3 conv. layer

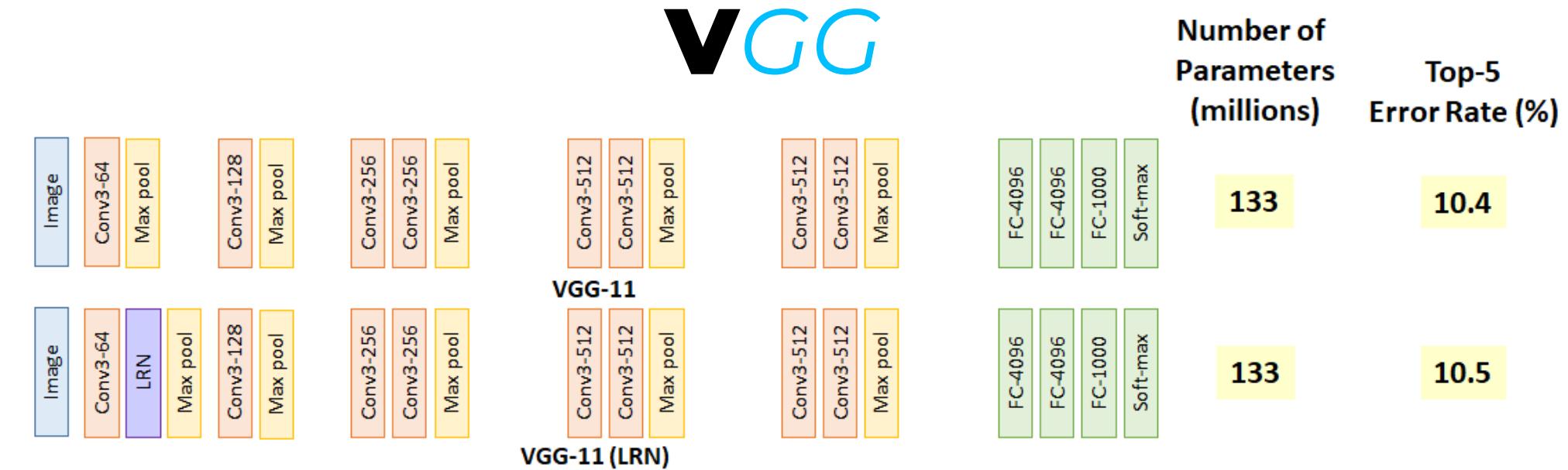


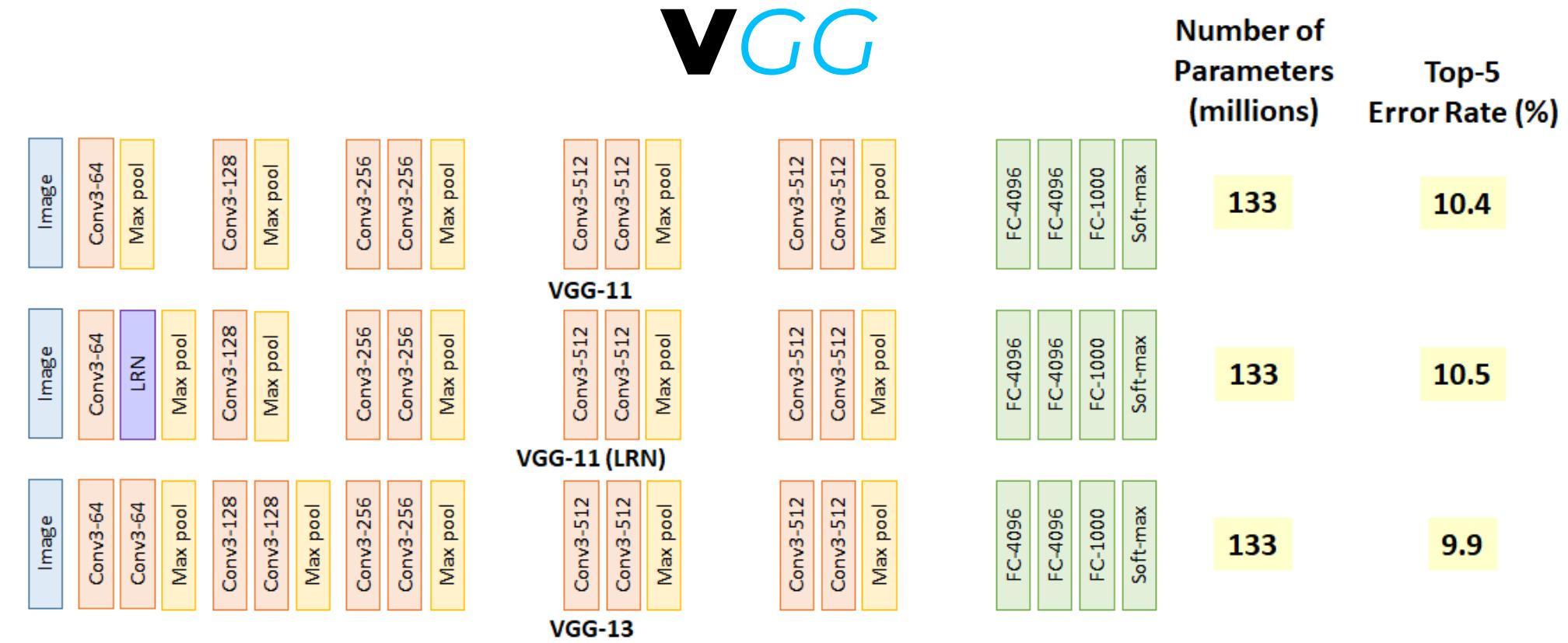
1 filtro 5×5, num. de parámetros $5 \times 5 = 25$
2 filtros 3×3, num. de parámetros $2 \times 3 \times 3 = 18$
El número de parámetros se reduce en **28%**

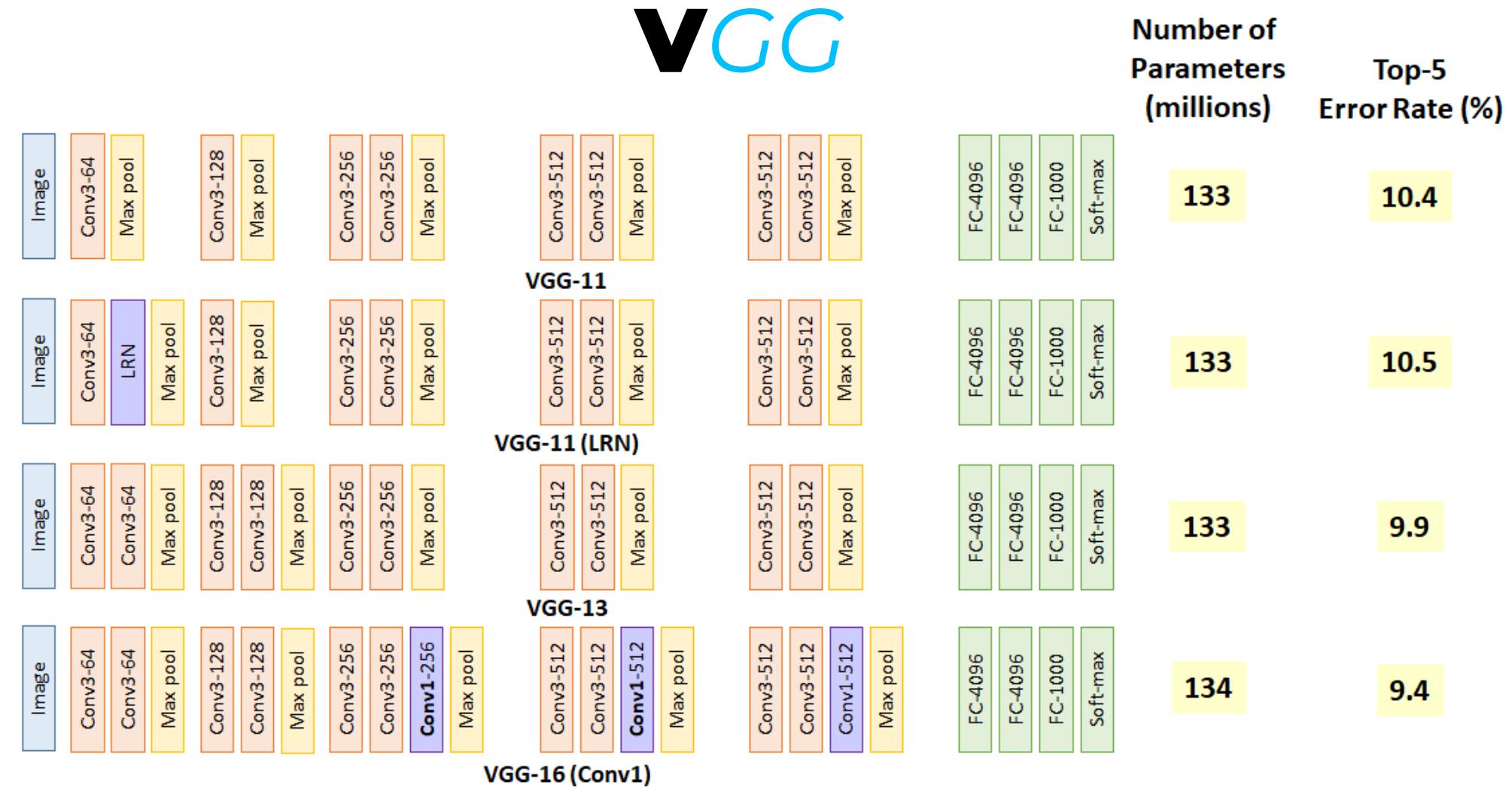
ZFNet 1 filtro 7×7, num. de parámetros: $7 \times 7 = 49$
3 filtros 3×3, num. de parámetros: $3 \times 3 \times 3 = 27$
El número de parámetros se reduce en **45%**

AlexNet 1 filtro 11×11, num. de parámetros = $11 \times 11 = 121$
5 filtros 3×3, num. de parámetros = $3 \times 3 \times 5 = 45$
El número de parámetros se reduce en **63%**

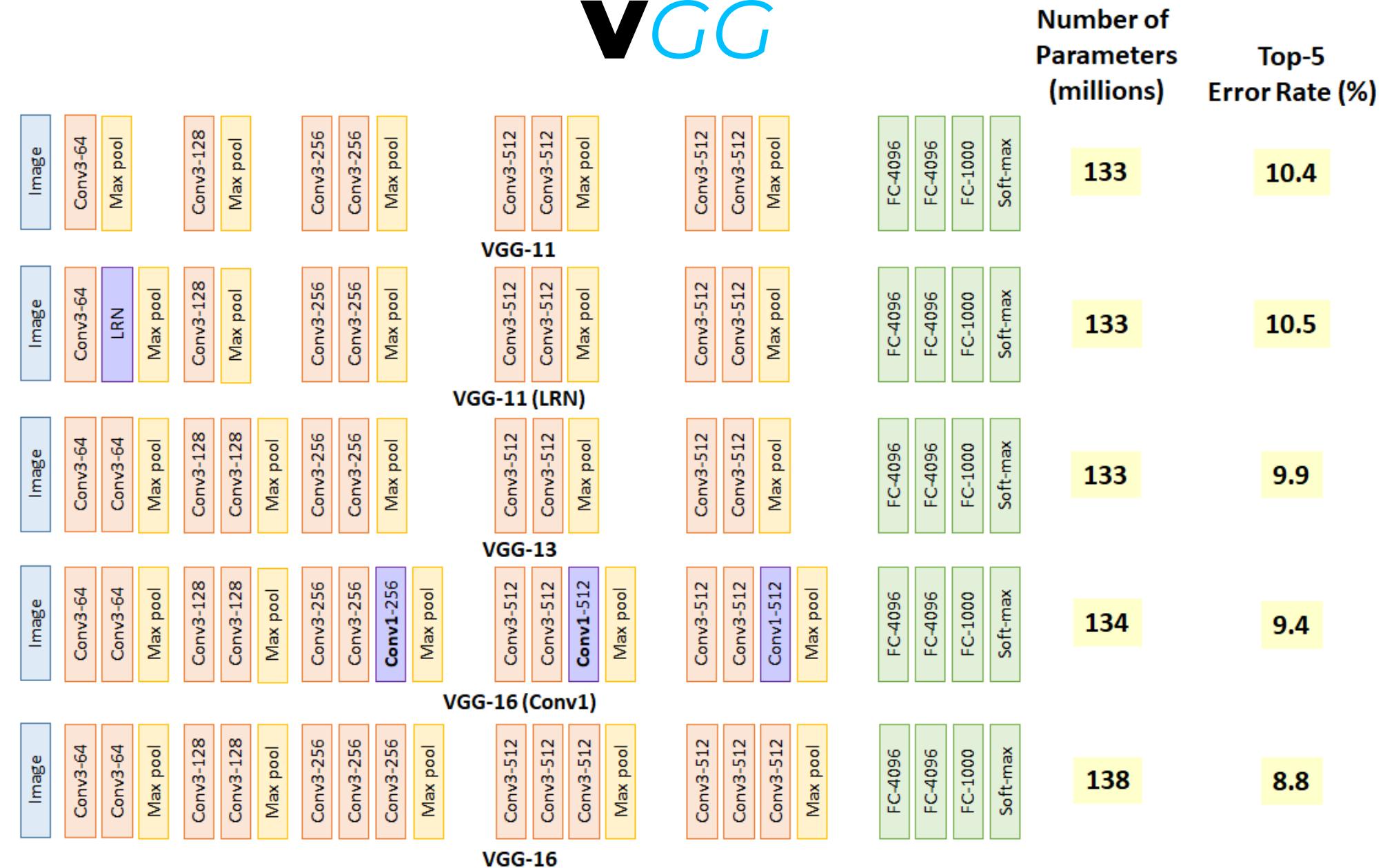




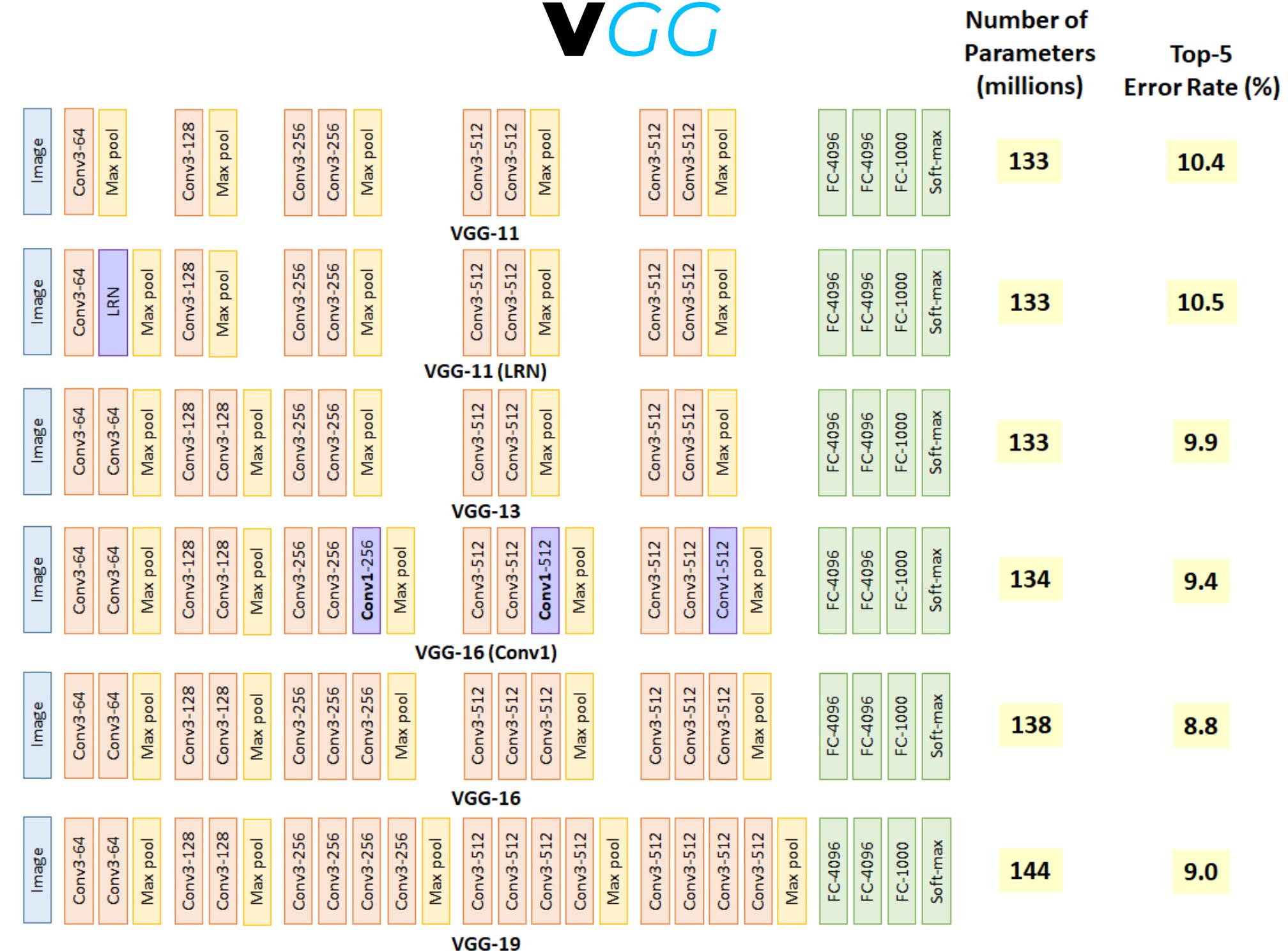




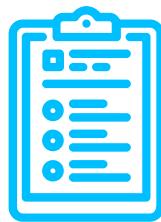
VGG



VGG



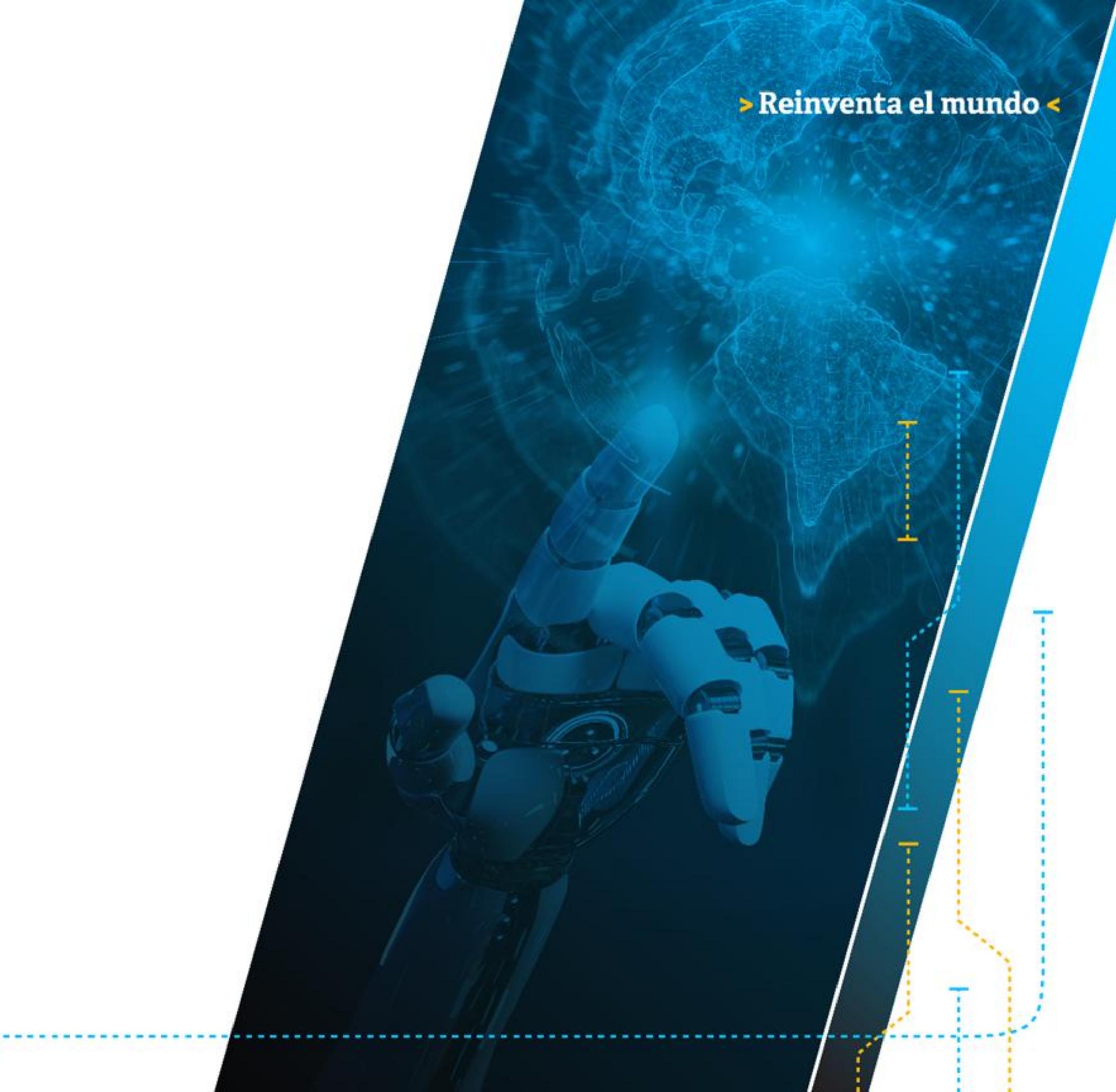
5.



Incep*tion*

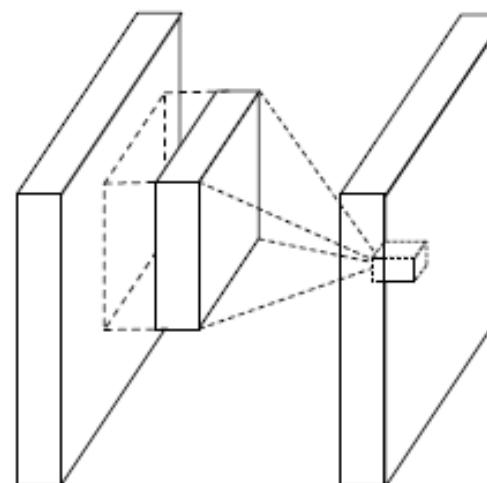
TRANSFORMA*TEC*

>Reinventa el mundo<

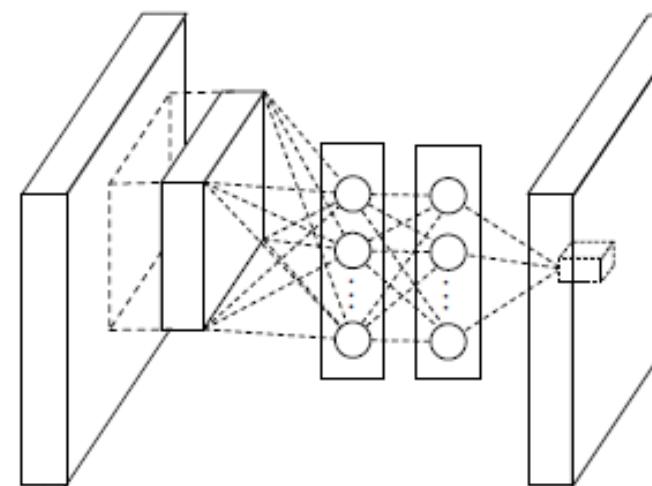


Network *in Network*

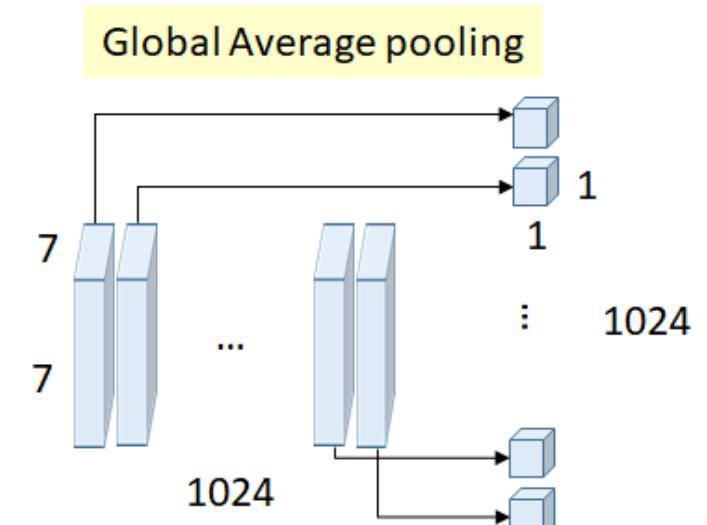
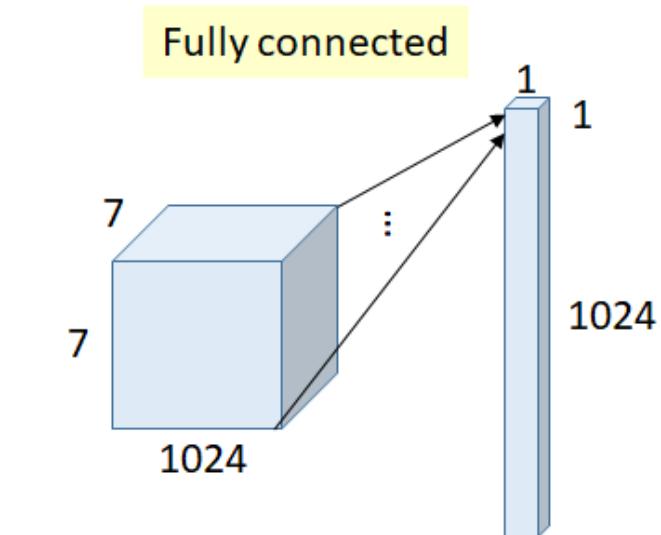
Linear Convolutional Layer



Mlpconv Layer



Fully Connected Layer VS Global Average Pooling Layer

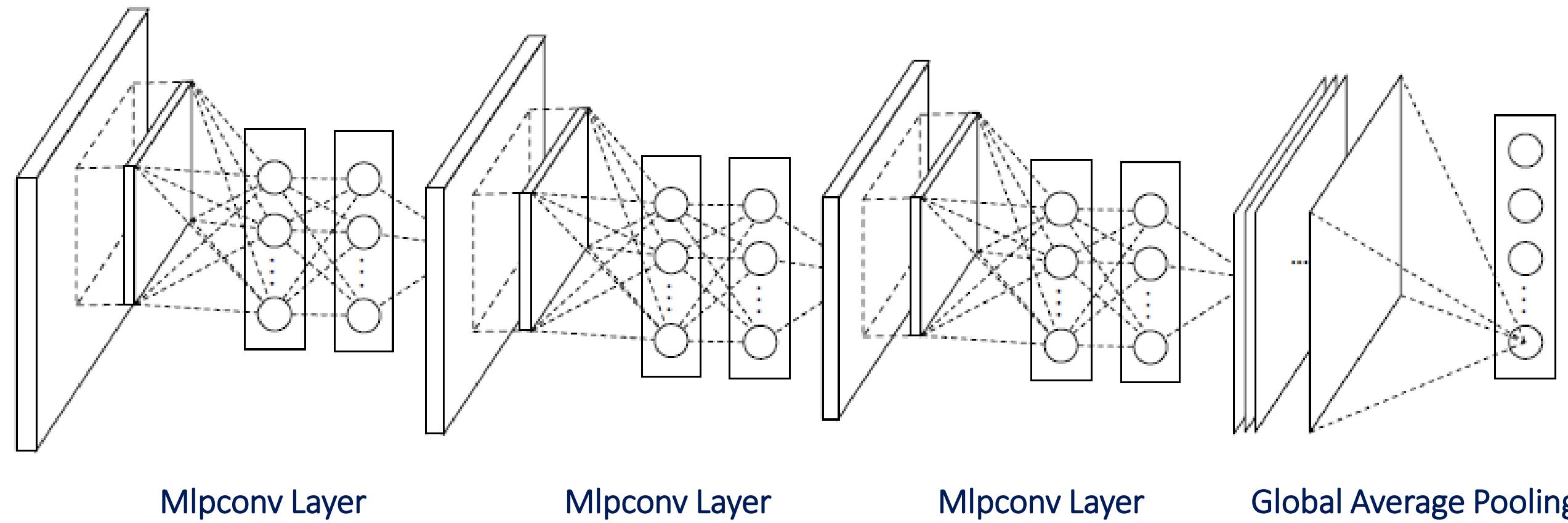


$$f_{i,j,k} = \max(w_k^T x_{i,j}, 0).$$

$$\begin{aligned} f_{i,j,k_1}^1 &= \max(w_{k_1}^1 {}^T x_{i,j} + b_{k_1}, 0). \\ &\vdots \\ f_{i,j,k_n}^n &= \max(w_{k_n}^n {}^T f_{i,j}^{n-1} + b_{k_n}, 0). \end{aligned}$$



Network *in Network*



Mlpconv Layer

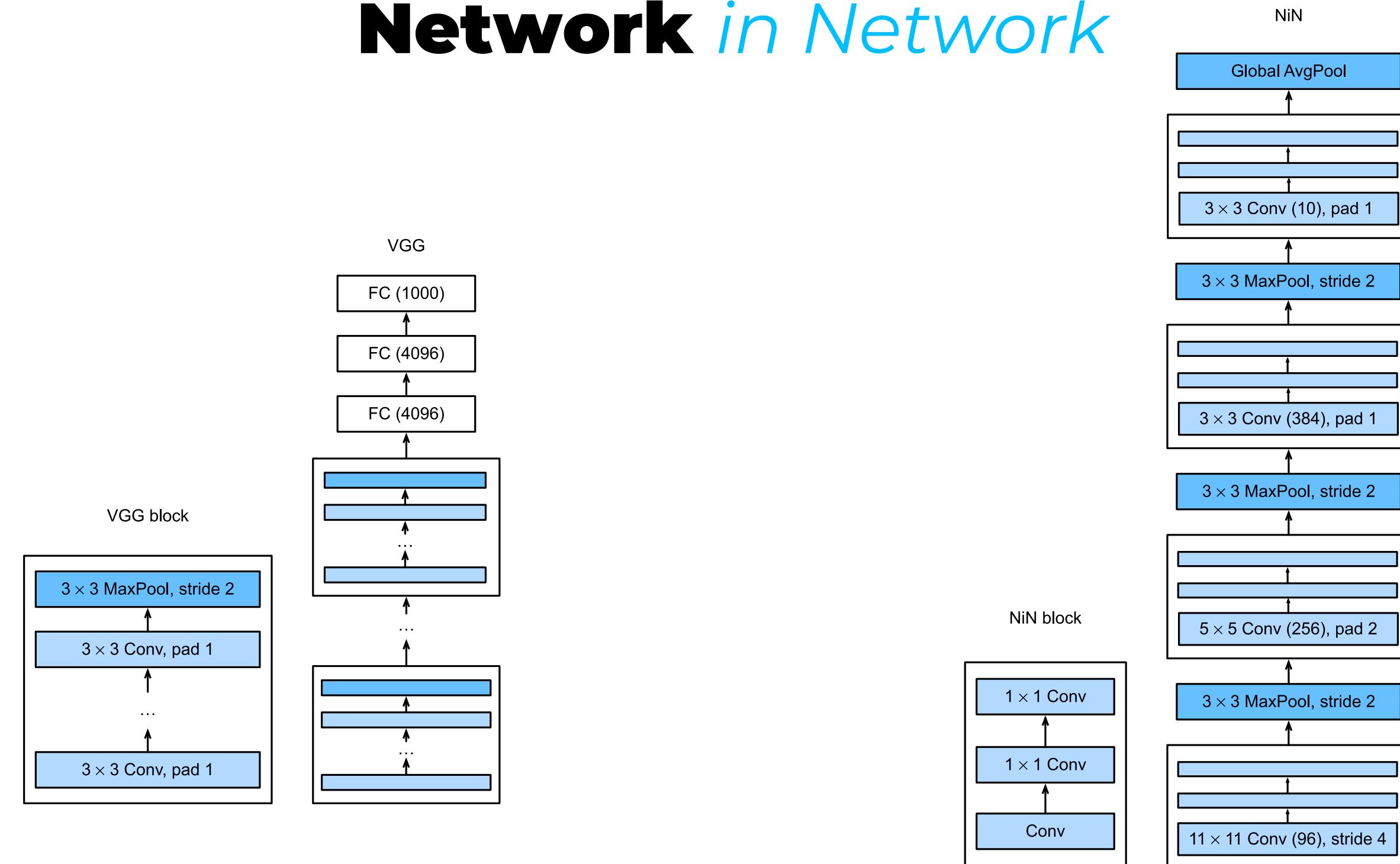
Mlpconv Layer

Mlpconv Layer

Global Average Pooling

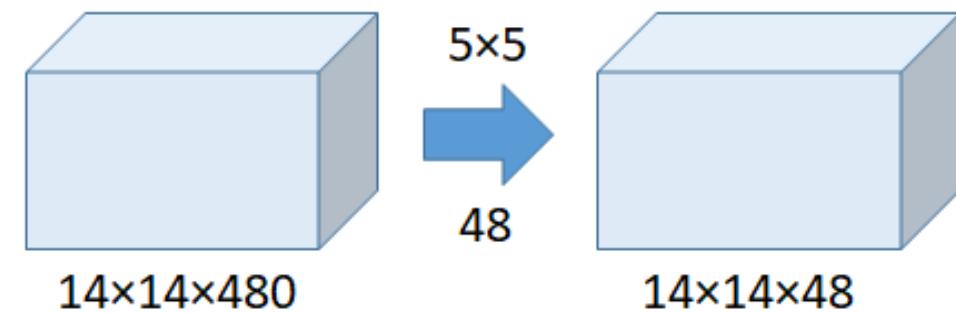


Network *in Network*



GoogLeNet

Bottleneck

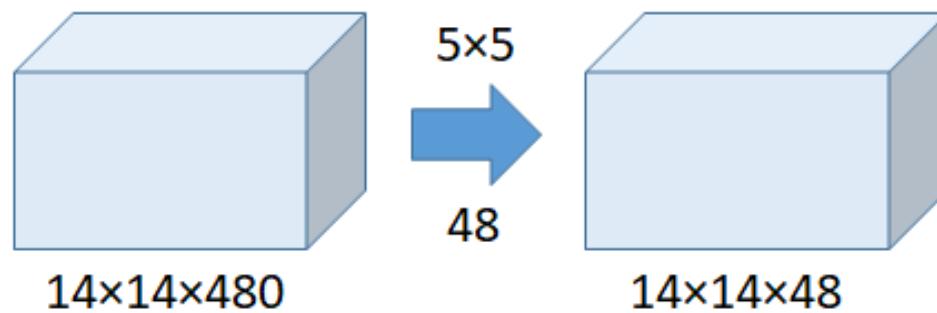


Número de operaciones = $(14 \times 14 \times 48) \times (5 \times 5 \times 480) = 112.9M$

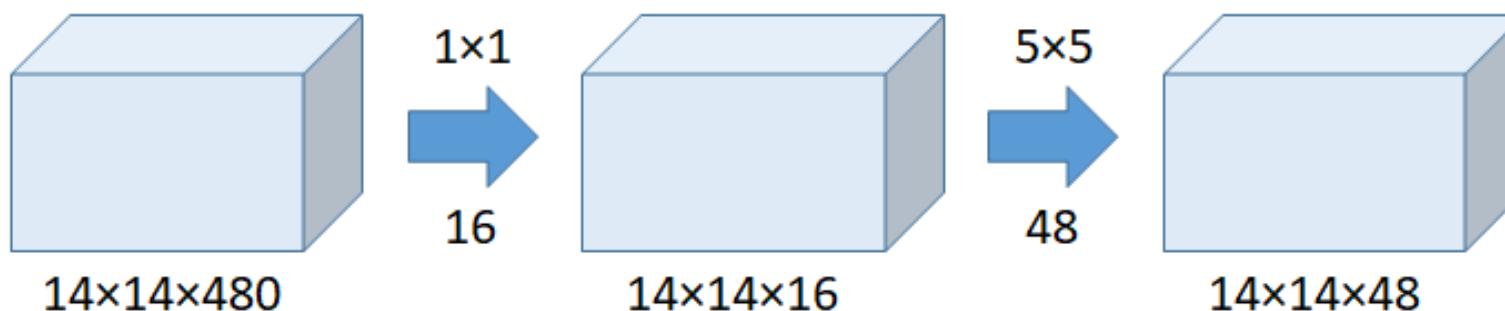


GoogLeNet

Bottleneck



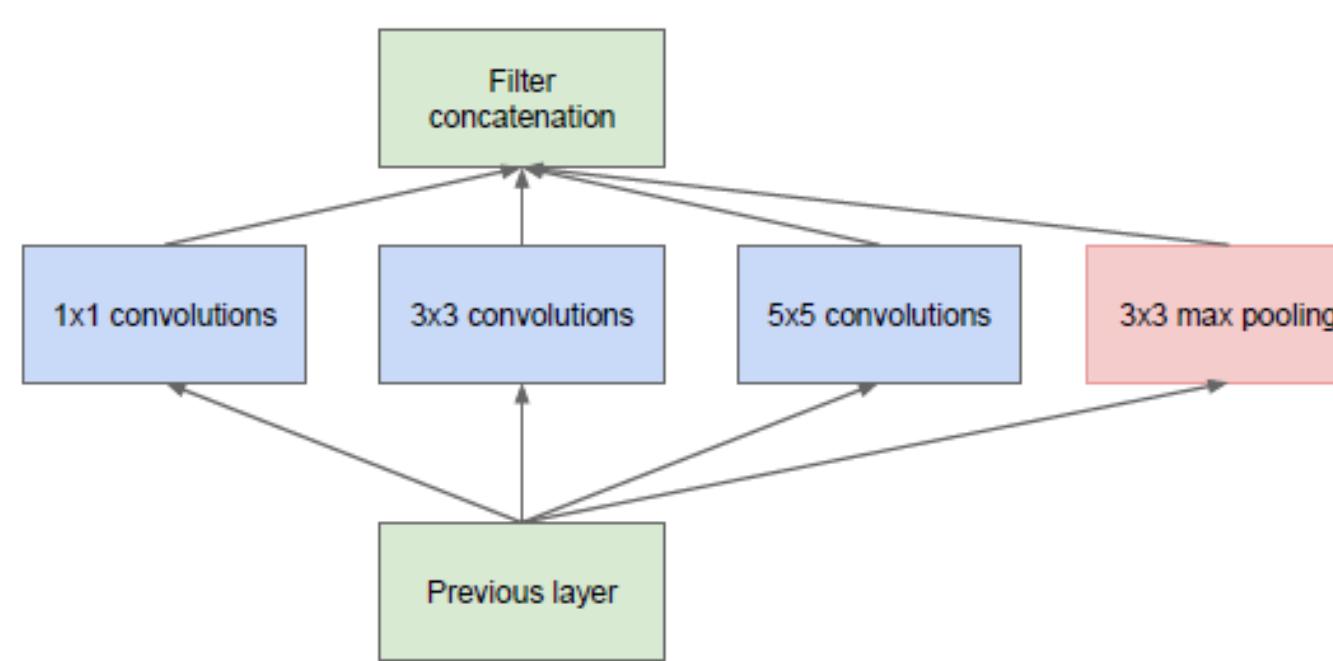
Número de operaciones = $(14 \times 14 \times 48) \times (5 \times 5 \times 480) = 112.9M$



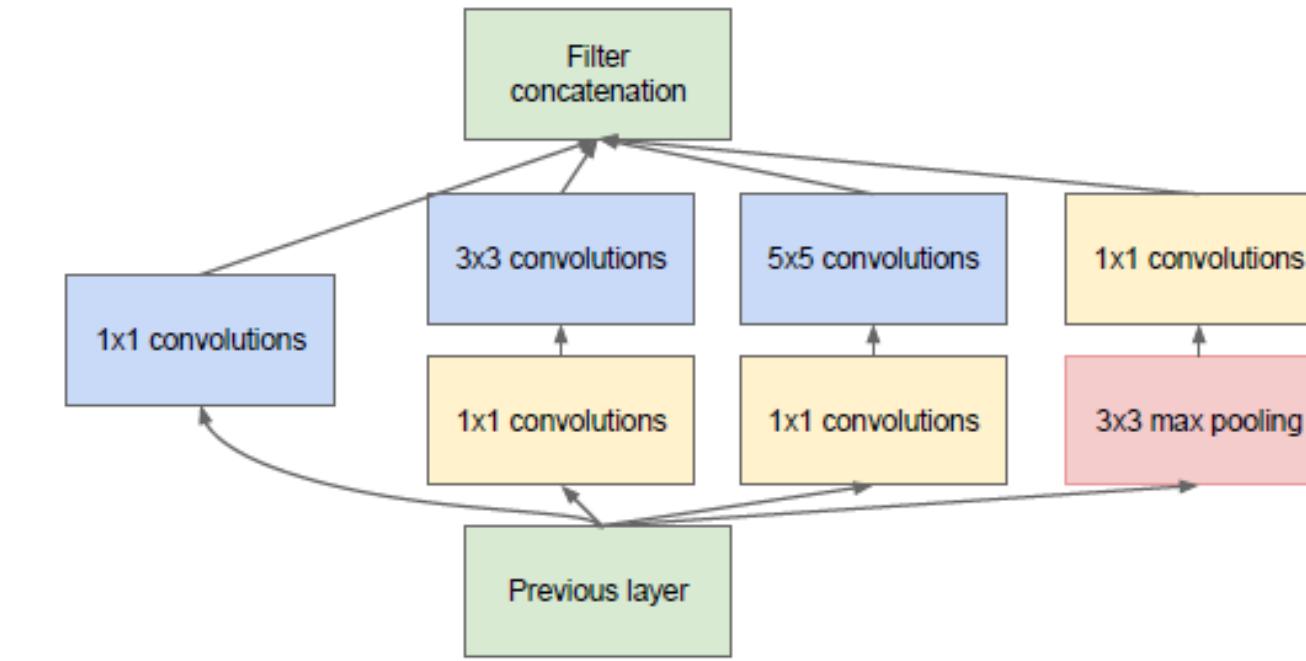
Número de operaciones $1 \times 1 = (14 \times 14 \times 16) \times (1 \times 1 \times 480) = 1.5M$
 Número de operaciones $5 \times 5 = (14 \times 14 \times 48) \times (5 \times 5 \times 16) = 3.8M$
 Número total de operaciones = $1.5M + 3.8M = 5.3M$



GoogLeNet



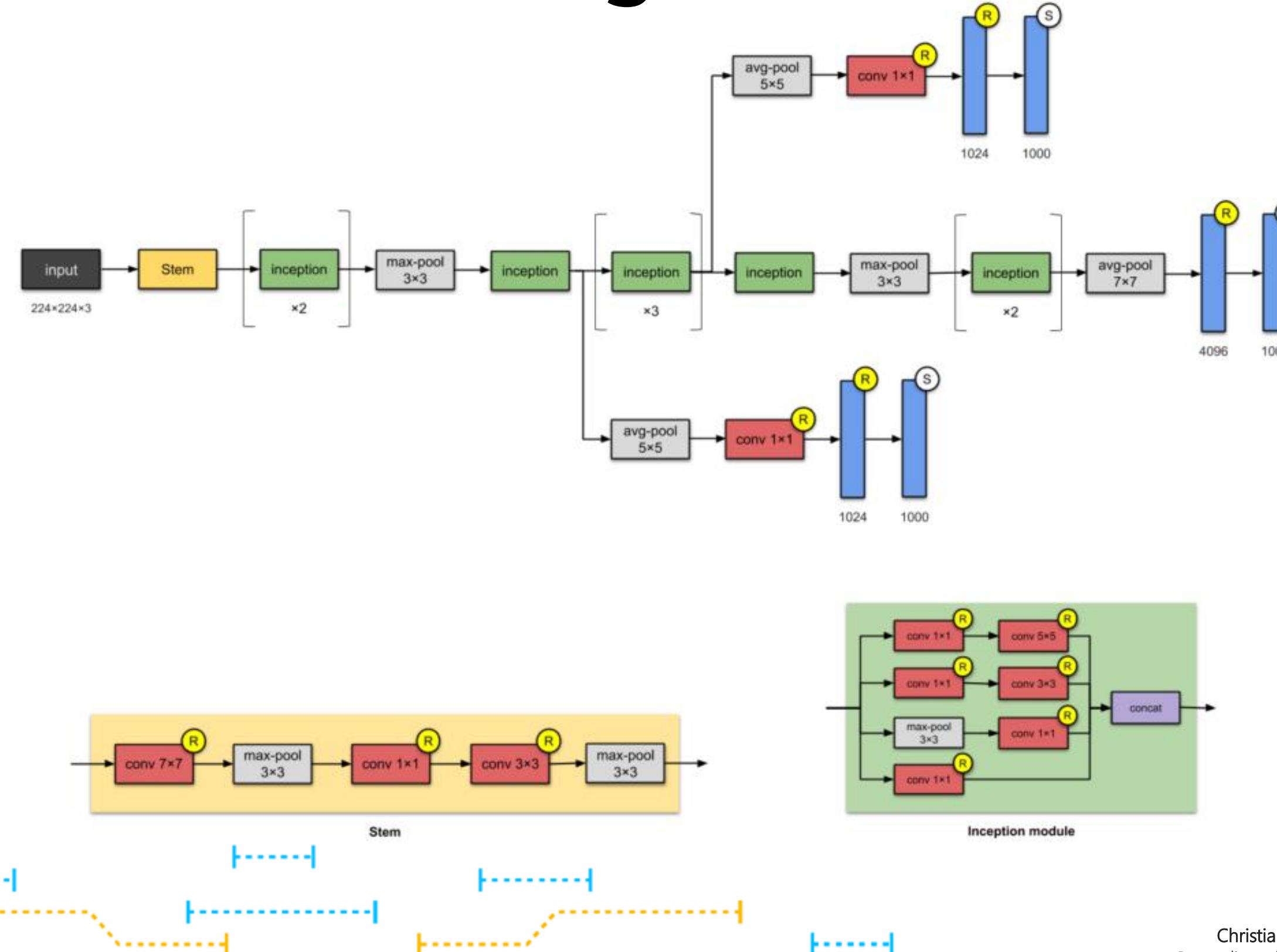
(a) Inception module, naïve version



(b) Inception module with dimensionality reduction



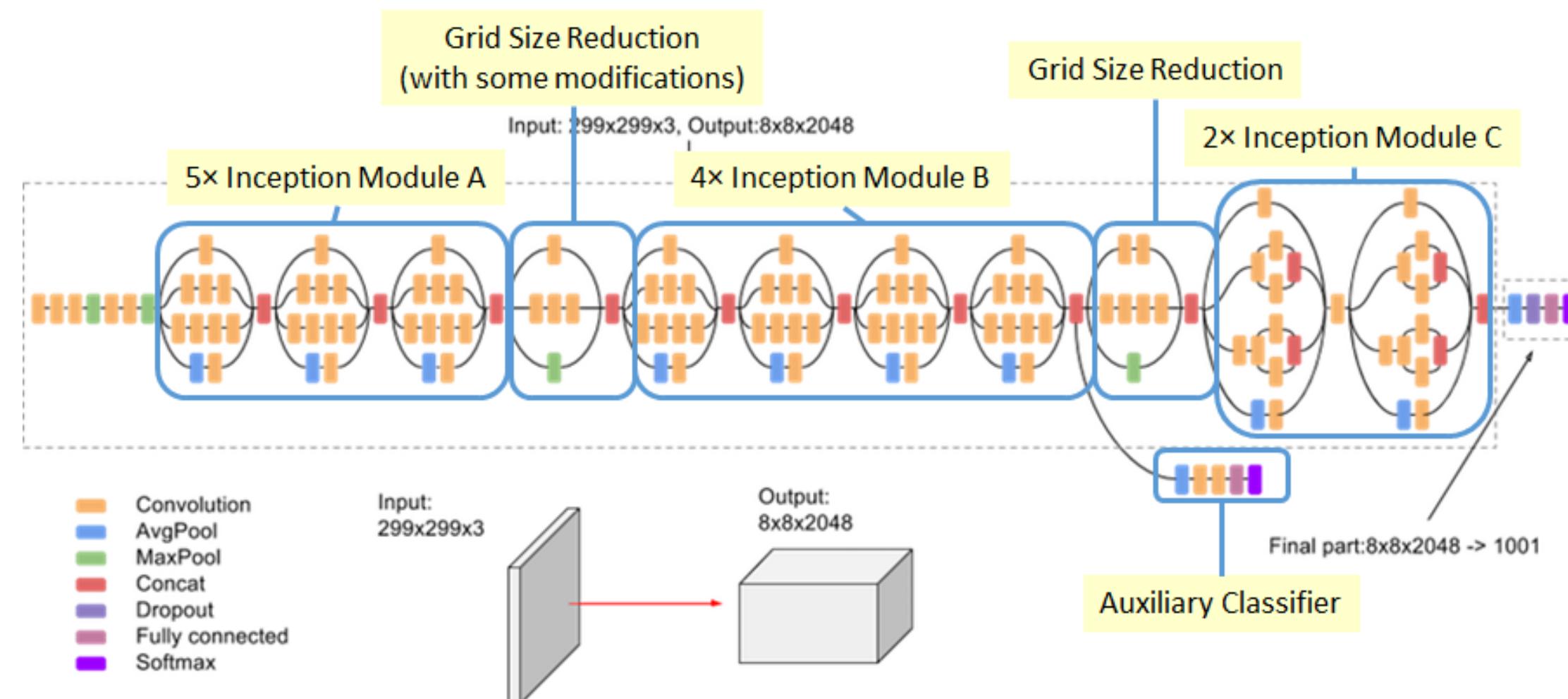
GoogLeNet



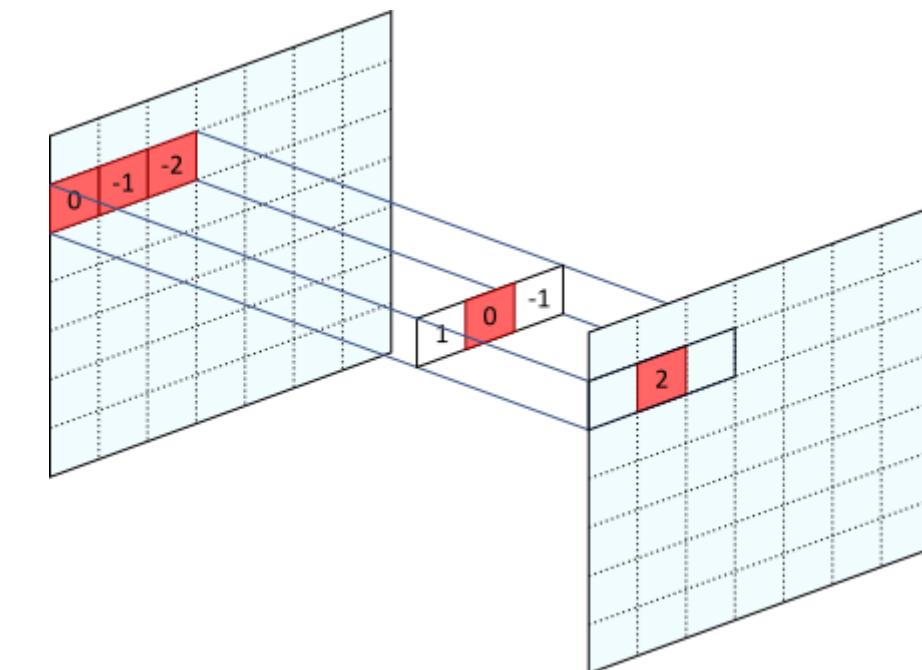
TRANSFORMATEC

Christian Szegedy et al. (2015) "Going deeper with convolutions".
Proceedings of the IEEE conference on computer vision and pattern recognition.

Inception v3



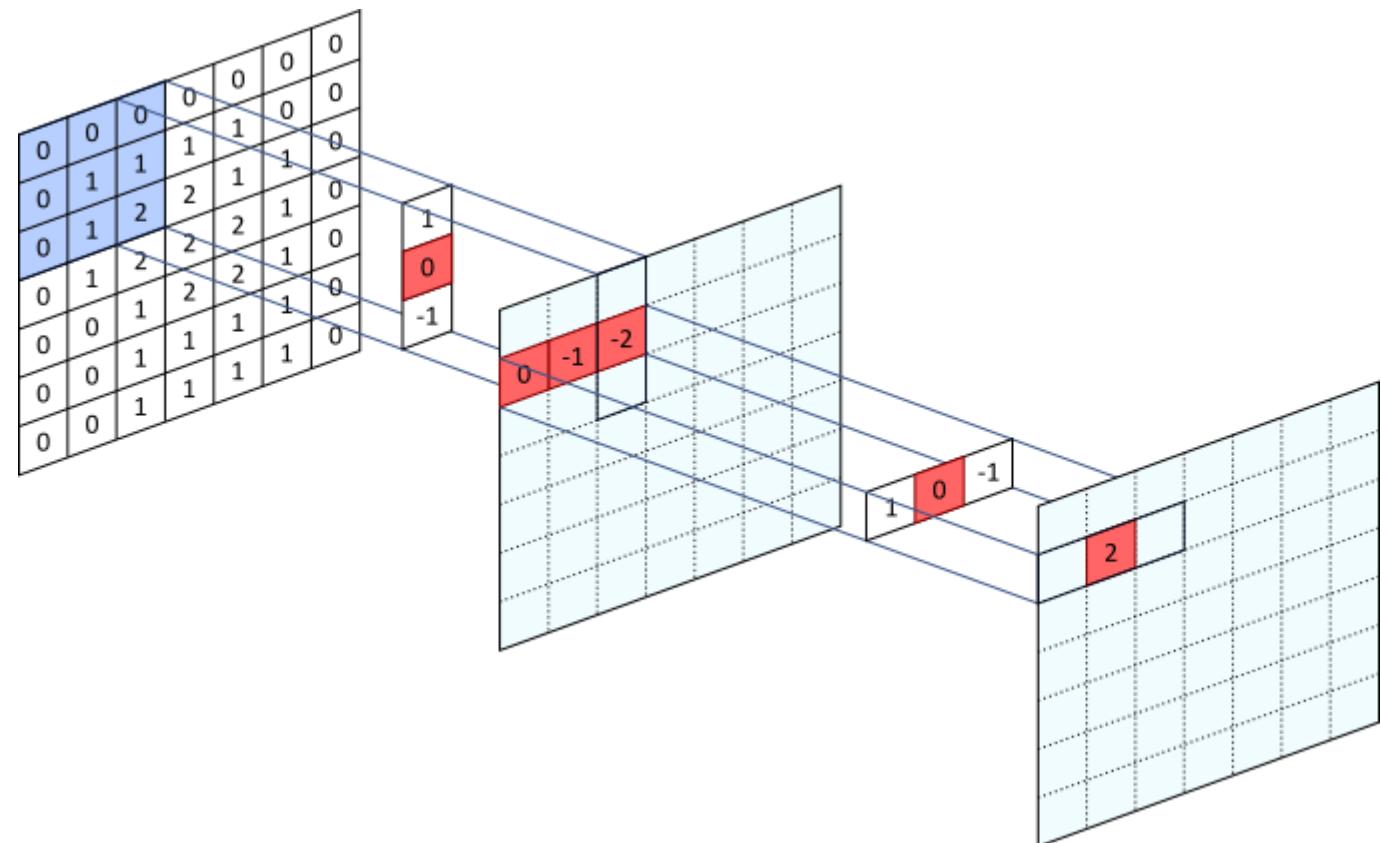
Inception v3



TRANSFORMATEC

Christian Szegedy et al. (2015) "Rethinking the Inception Architecture for Computer Vision".
Proceedings of the IEEE conference on computer vision and pattern recognition

Inception v3

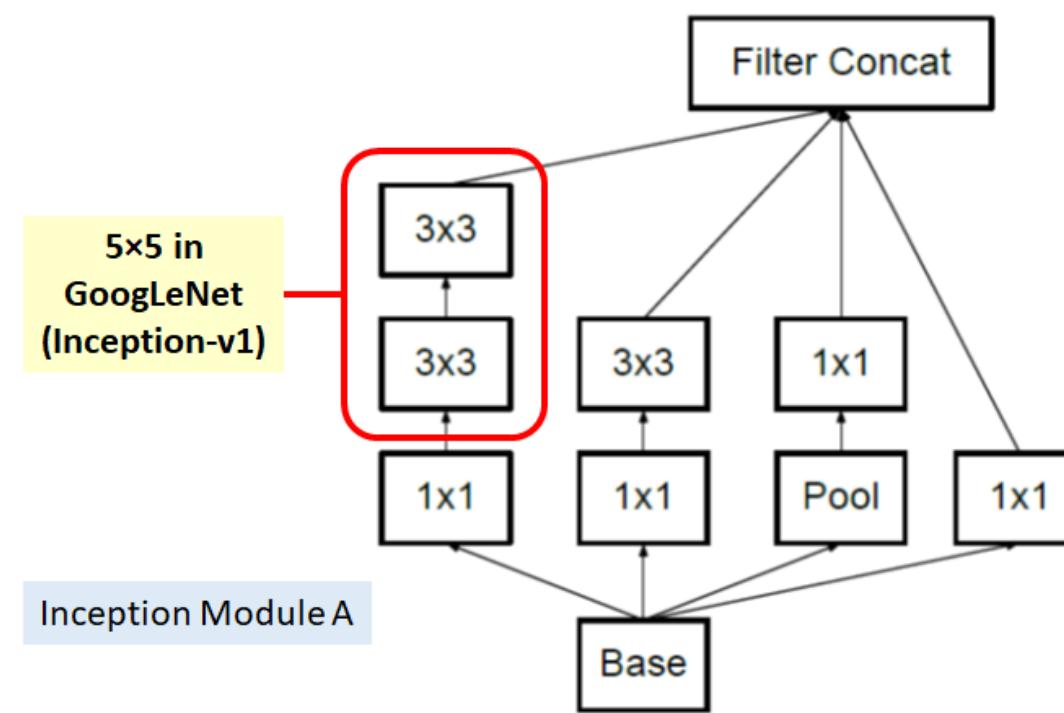


TRANSFORMATEC

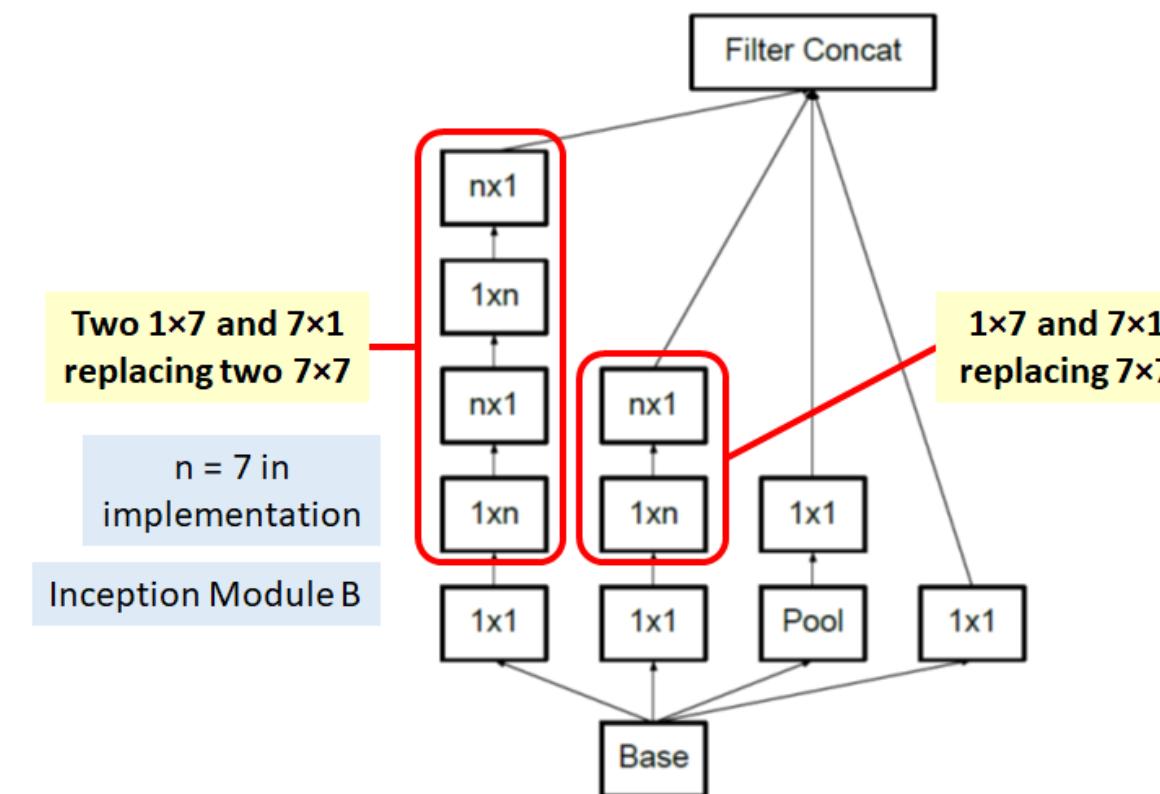
Christian Szegedy et al. (2015) "Rethinking the Inception Architecture for Computer Vision".
Proceedings of the IEEE conference on computer vision and pattern recognition

Inception v3

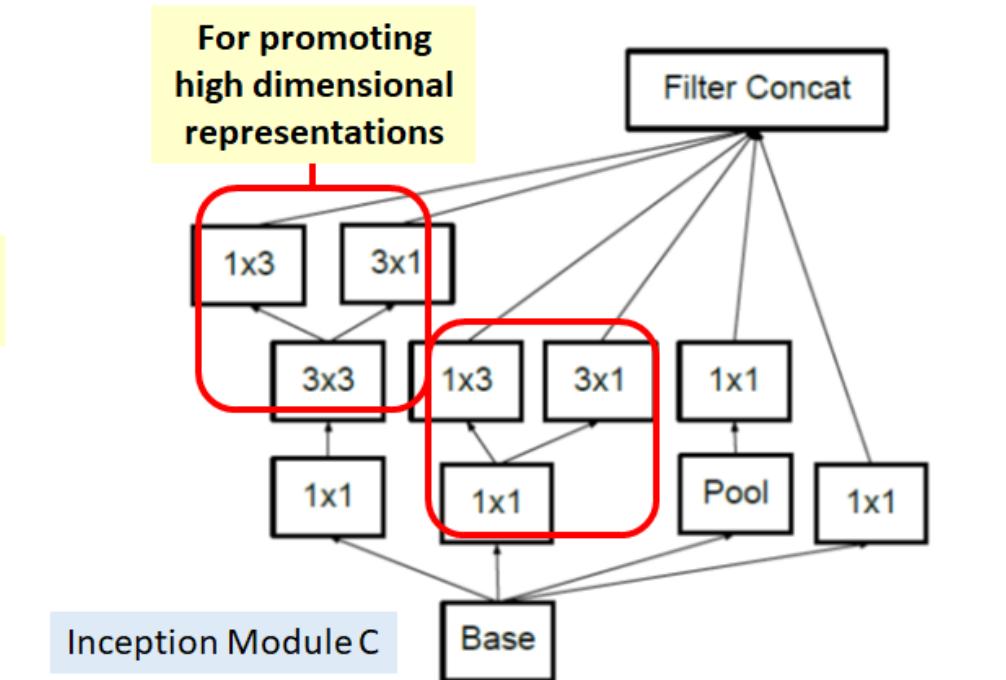
Inception Module A



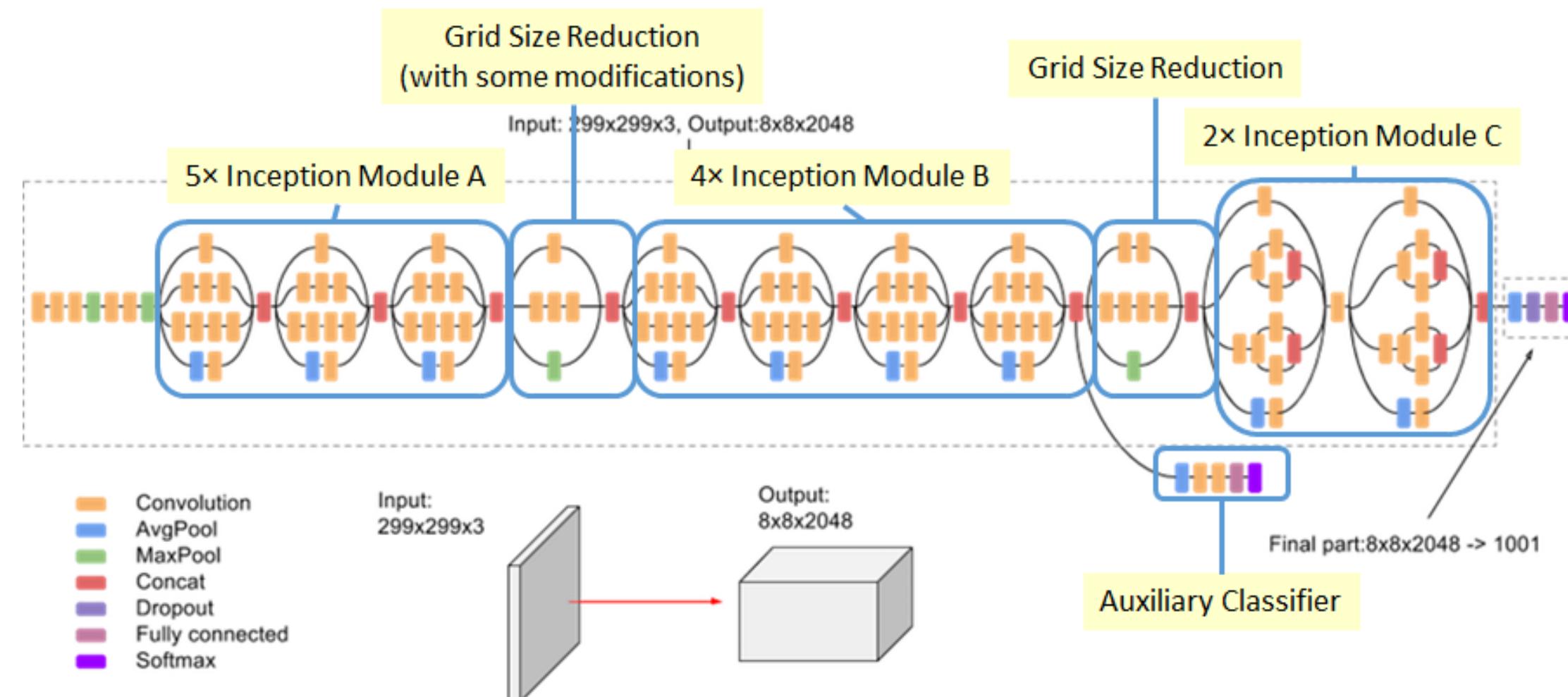
Inception Module B



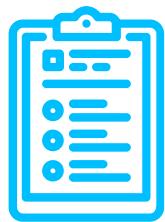
Inception Module C



Inception v3



6.



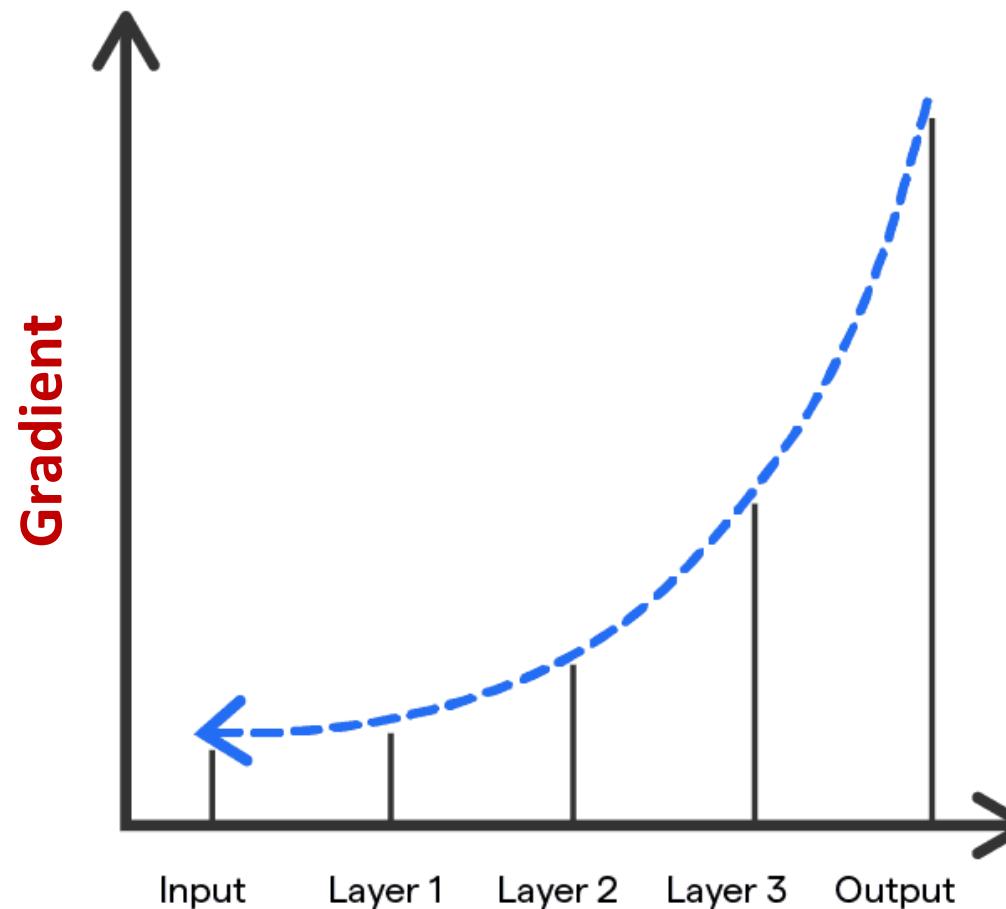
ResNet

TRANSFORMATEC

> Reinventa el mundo <



Vanishing gradient problem



Supongamos una red con L capas, y cada capa l tiene un conjunto de pesos W_l y un conjunto de activaciones A_l . En un forward pass, para una capa l , tenemos:

$$Z_l = W_l A_{l-1} + b_l \quad A_l = \sigma(Z_l)$$

Consideremos la derivada de la función de pérdida \mathcal{L} con respecto a los pesos en la capa l :

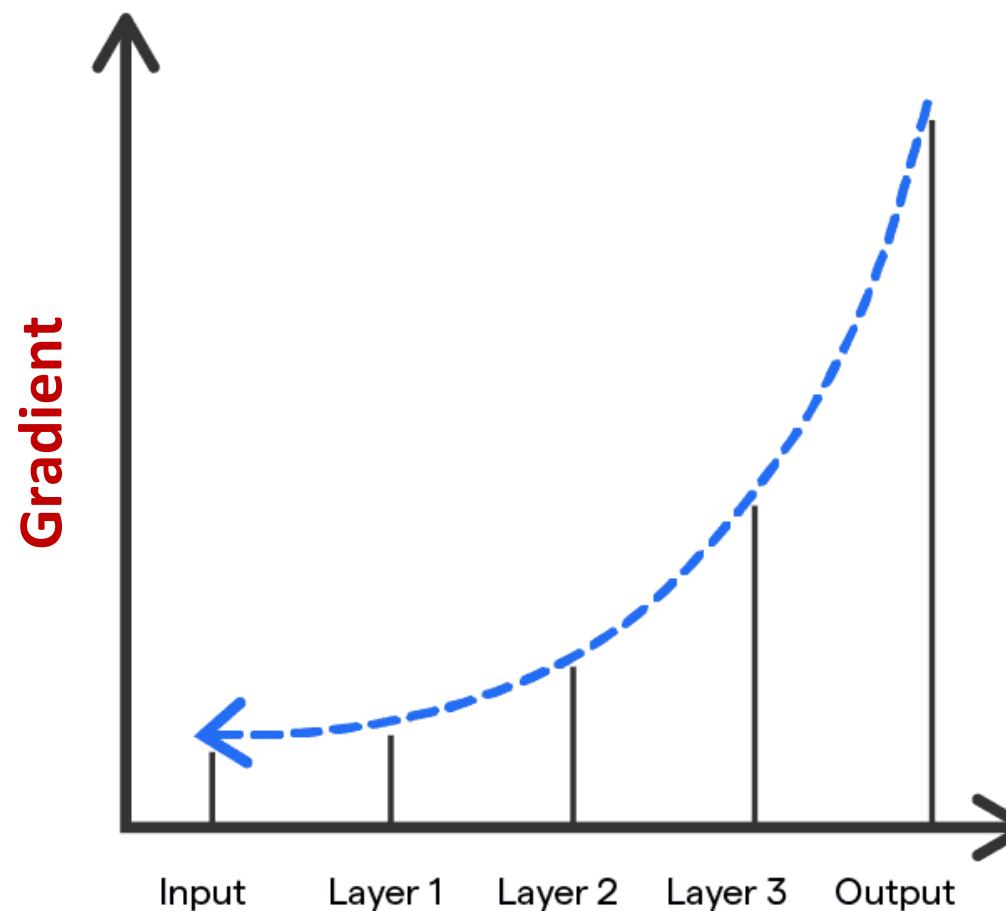
$$\frac{\partial \mathcal{L}}{\partial W_l} = \delta_l A_{l-1}^\top \quad \delta_l = \frac{\partial \mathcal{L}}{\partial Z_l} = (W_{l+1}^\top \delta_{l+1}) \odot \sigma'(Z_l)$$

Podemos desarrollar esto recursivamente:

$$\delta_l = \left(W_{l+1}^\top \left(W_{l+2}^\top \left(\dots \left(W_L^\top \delta_L \right) \right) \right) \right) \odot \sigma'(Z_l)$$



Vanishing gradient problem

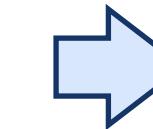


Supongamos una red con L capas, y cada capa l tiene un conjunto de pesos W_l y un conjunto de activaciones A_l . En un forward pass, para una capa l , tenemos:

$$Z_l = W_l A_{l-1} + b_l \quad A_l = \sigma(Z_l)$$

Consideremos la derivada de la función de pérdida \mathcal{L} con respecto a los pesos en la capa l :

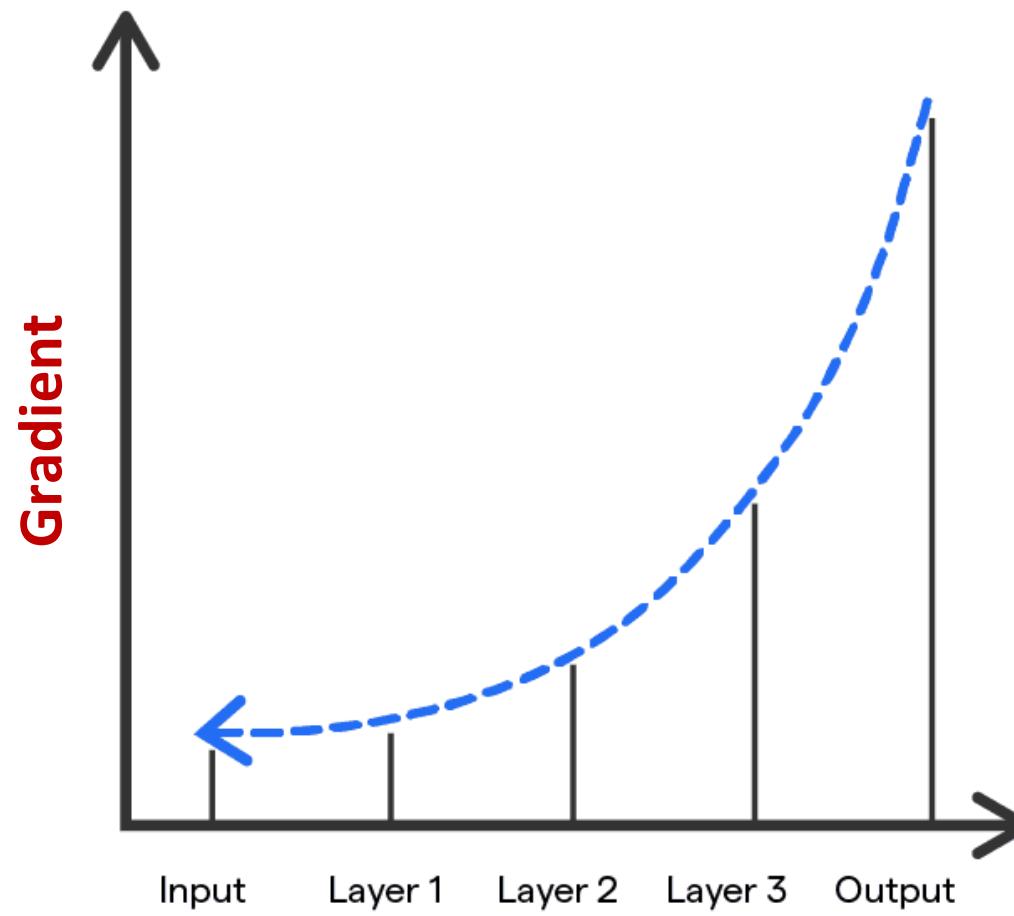
$$\frac{\partial \mathcal{L}}{\partial W_l} = \left(\prod_{k=l+1}^L W_k^\top \odot \sigma'(Z_{k-1}) \right) \delta_L A_0^\top$$

Para las funciones: $\sigma(z) = \tanh(z)$  $|\sigma(z)| < 1$
 $\sigma(z) = \text{Sigmoid}(z)$

Entonces: $\prod_{k=l+1}^L \sigma'(Z_{k-1}) \rightarrow 0$ cuando $L - l$ es muy grande



Vanishing gradient problem



Supongamos una red con L capas, y cada capa l tiene un conjunto de pesos W_l y un conjunto de activaciones A_l . En un forward pass, para una capa l , tenemos:

$$Z_l = W_l A_{l-1} + b_l \quad A_l = \sigma(Z_l)$$

Consideremos la derivada de la función de pérdida \mathcal{L} con respecto a los pesos en la capa l :

$$\frac{\partial \mathcal{L}}{\partial W_l} = \left(\prod_{k=l+1}^L W_k^\top \odot \sigma'(Z_{k-1}) \right) \delta_L A_0^\top$$

Para el caso de: $\sigma(z) = \text{ReLU}(z)$

- Cuando $Z_{k-1} > 0$ entonces $\sigma'(Z_{k-1}) = 1$.
- Cuando $Z_{k-1} \leq 0$ entonces $\sigma'(Z_{k-1}) = 0$.

Dead ReLU!

Algunas neuronas no son entrenadas



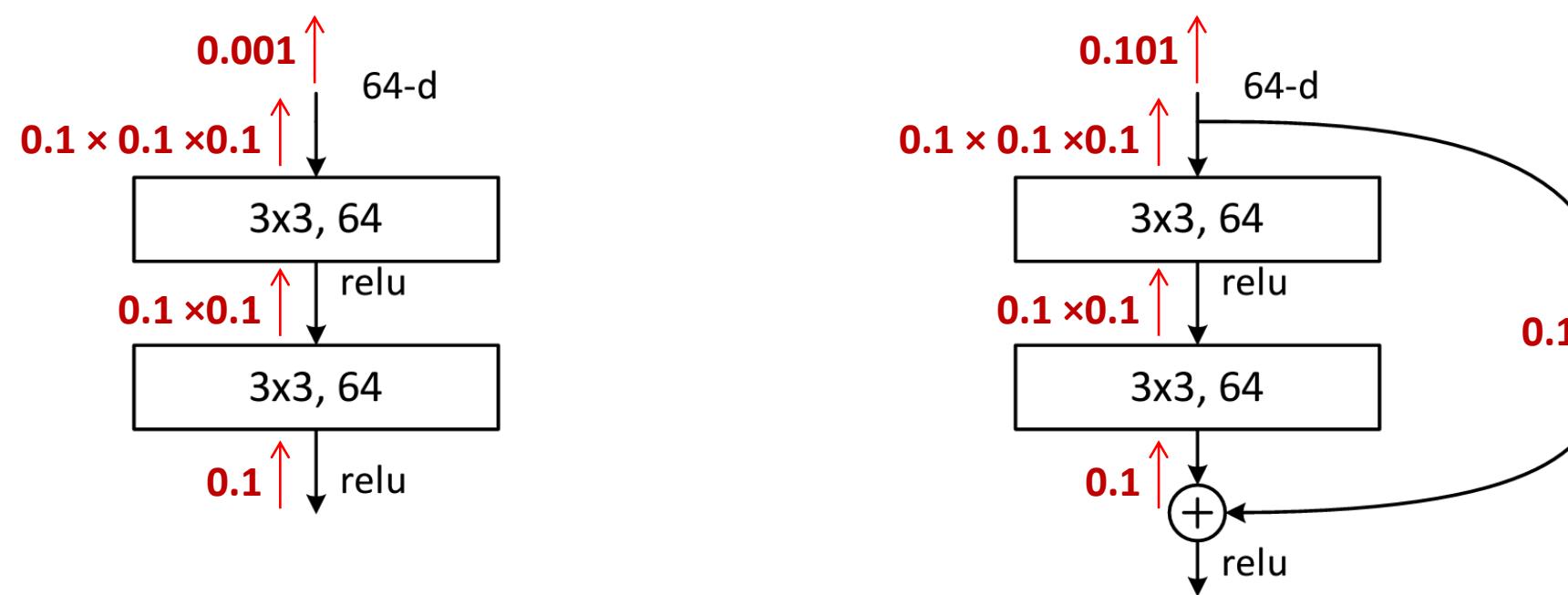
ResNet

Solución:

skip-connection

$$Z_l = \mathcal{F}(x) + Z_{l-1}$$

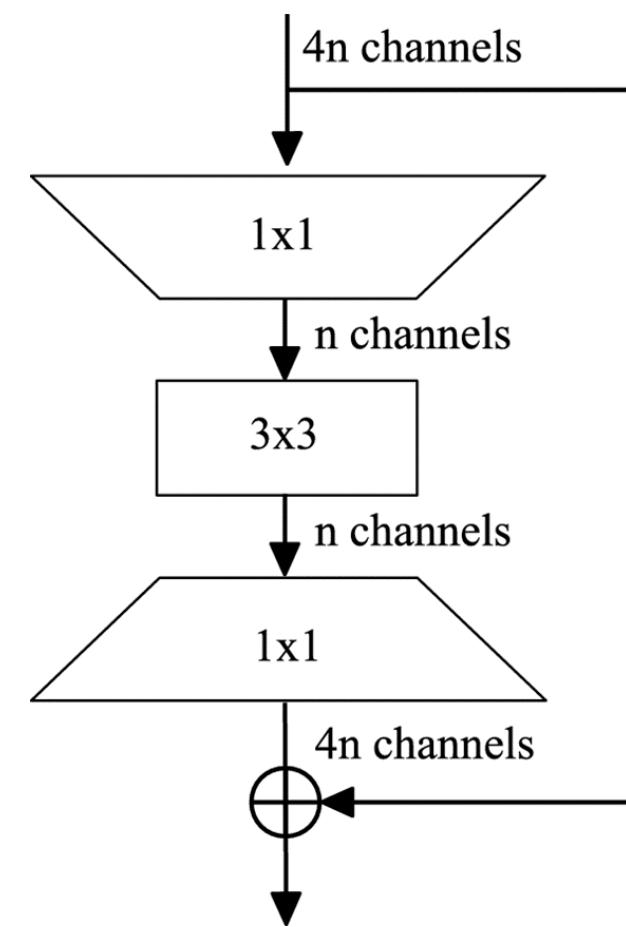
El gradiente del error en la capa $l - 1$: $\delta_{l-1} = \frac{\partial \mathcal{L}}{\partial Z_{l-1}} = \frac{\partial \mathcal{L}}{\partial Z_l} \left(\frac{\partial Z_l}{\partial Z_{l-1}} \right) = \frac{\partial \mathcal{L}}{\partial Z_l} \left(\frac{\partial \mathcal{F}}{\partial Z_{l-1}} + 1 \right)$



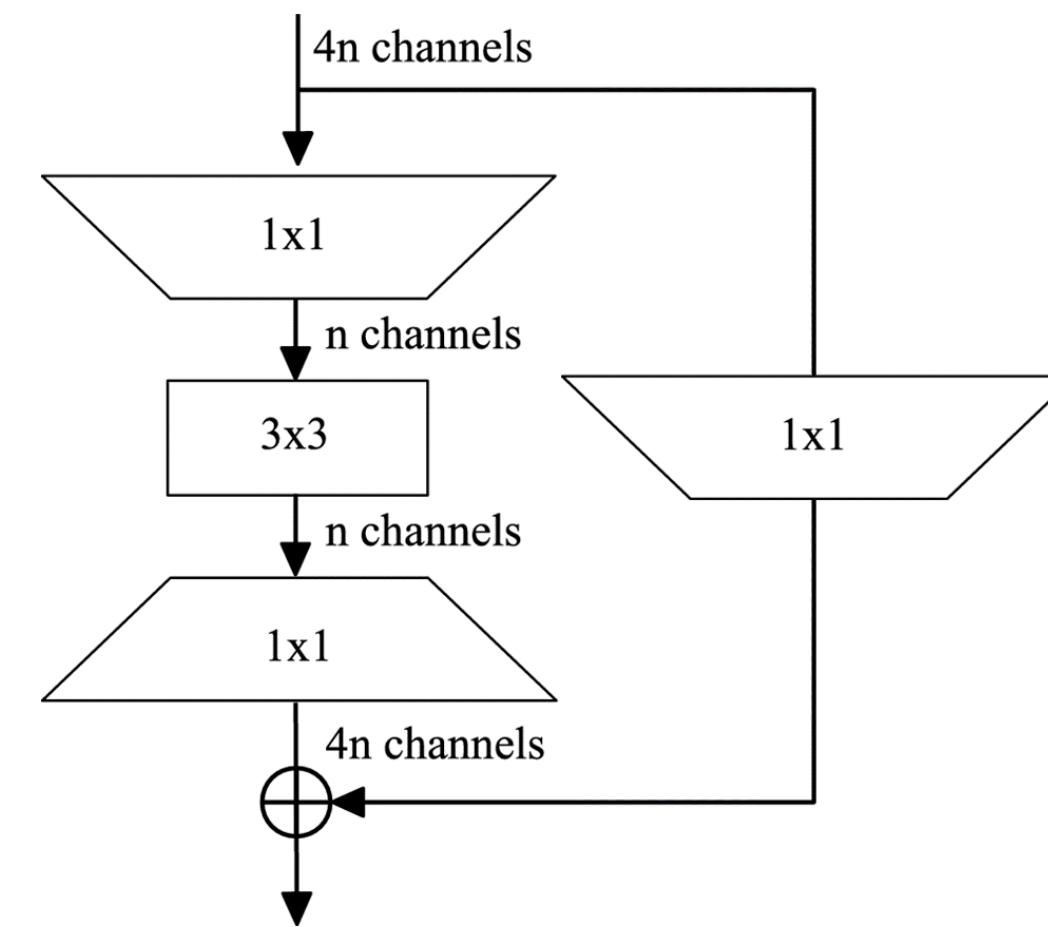
ResNet

Bottleneck architecture

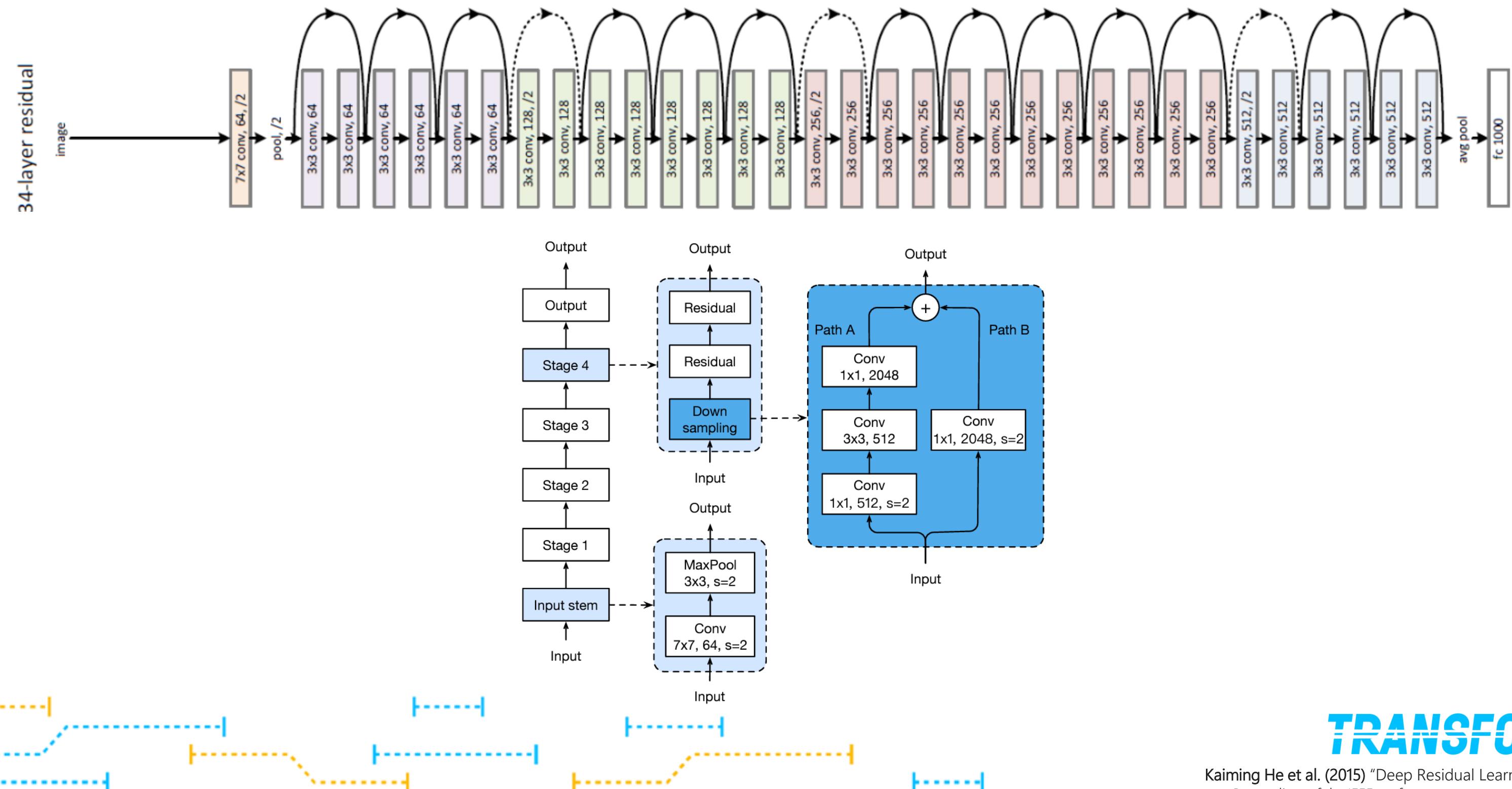
Identity Shortcuts



Projection Shortcuts



ResNet



ResNet

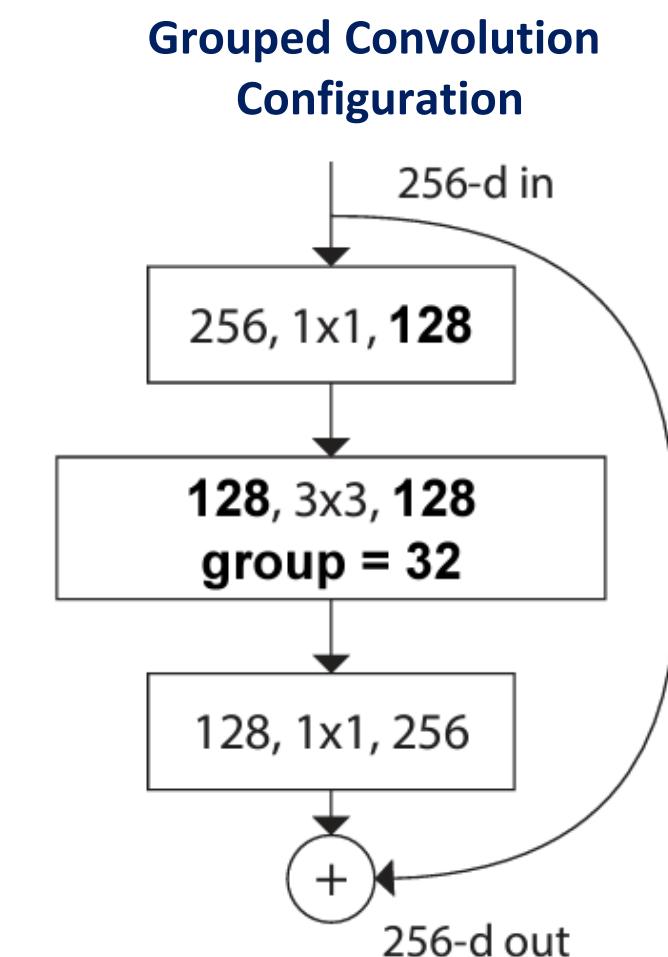
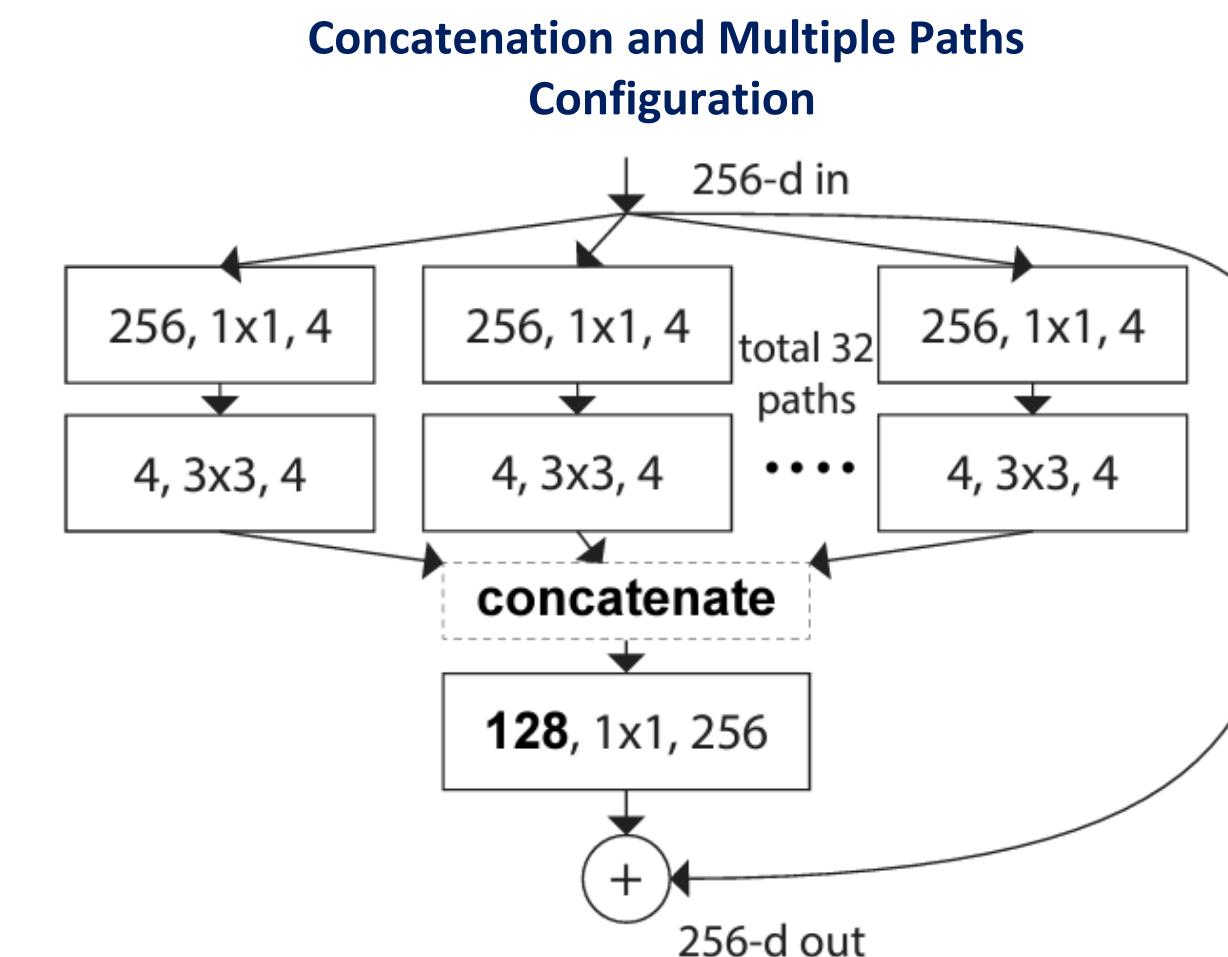
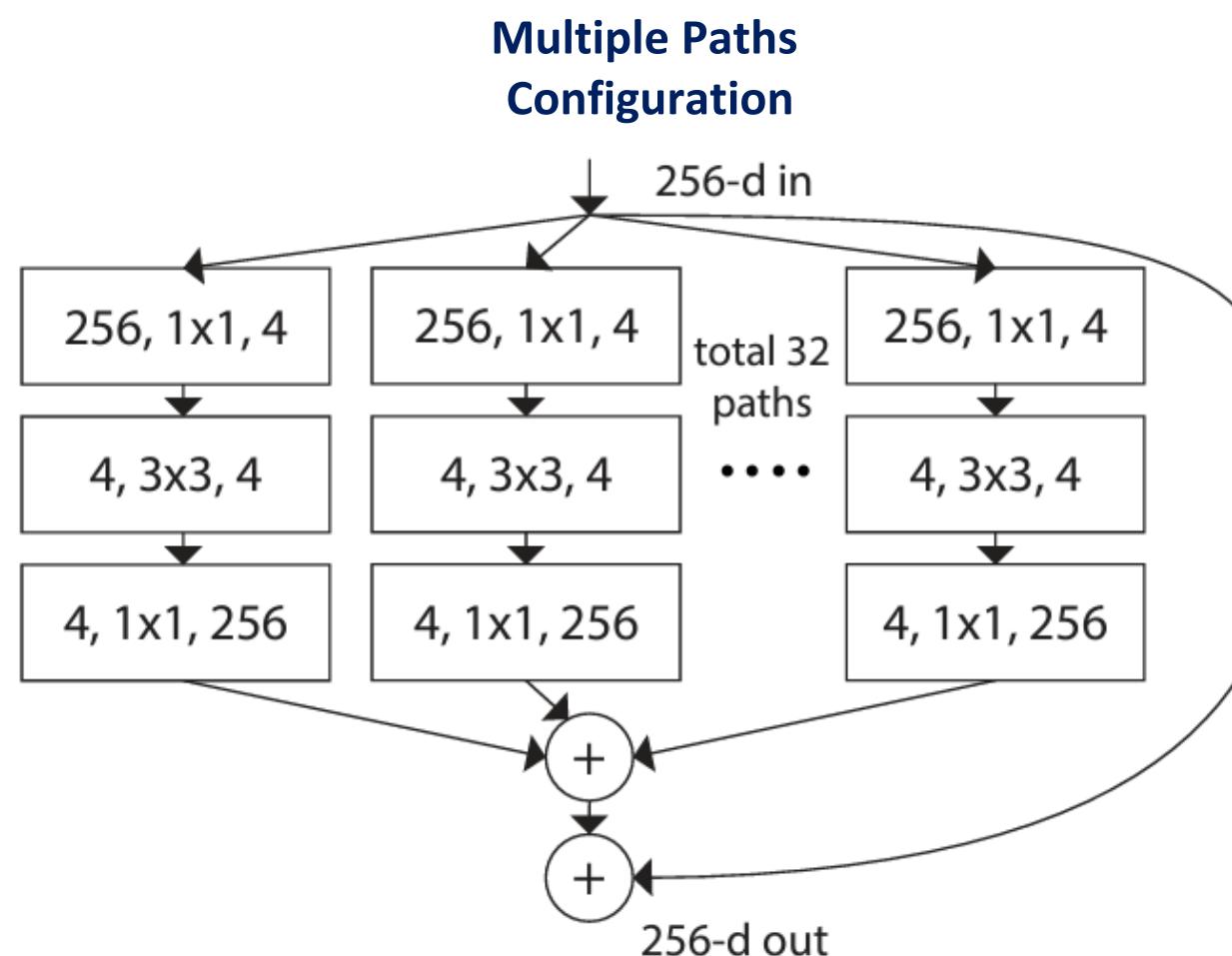
layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			$7 \times 7, 64, \text{stride } 2$		
				$3 \times 3 \text{ max pool, stride } 2$		
conv2_x	56×56	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	28×28	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	1×1			average pool, 1000-d fc, softmax		
FLOPs		1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9



TRANSFORMATEC

Kaiming He et al. (2015) "Deep Residual Learning for Image Recognition".
Proceedings of the IEEE conference on computer vision and pattern recognition

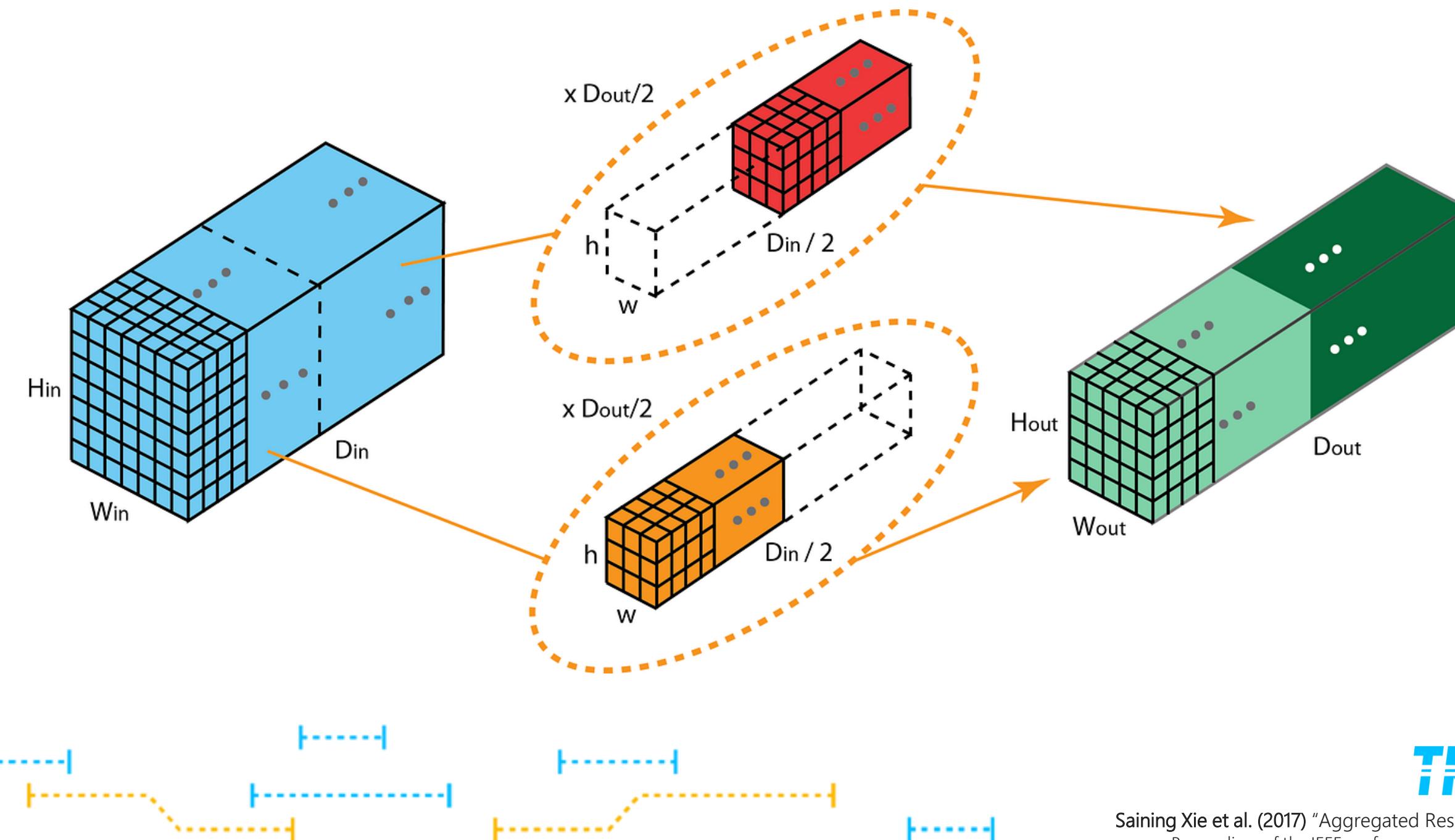
ResNeXt



TRANSFORMATEC

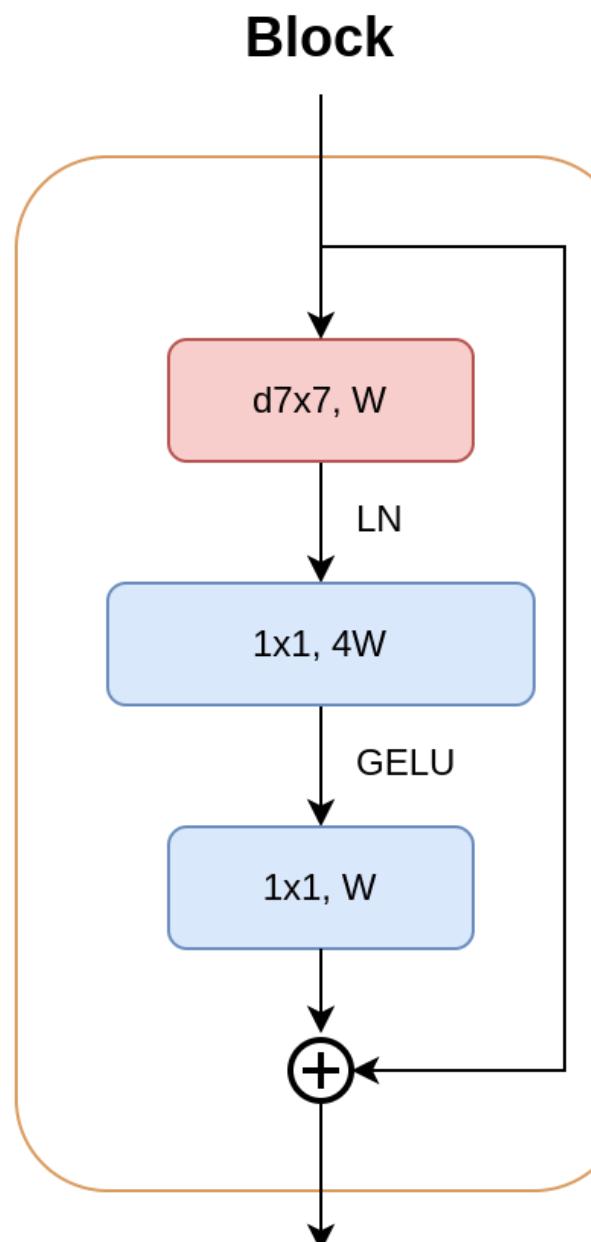
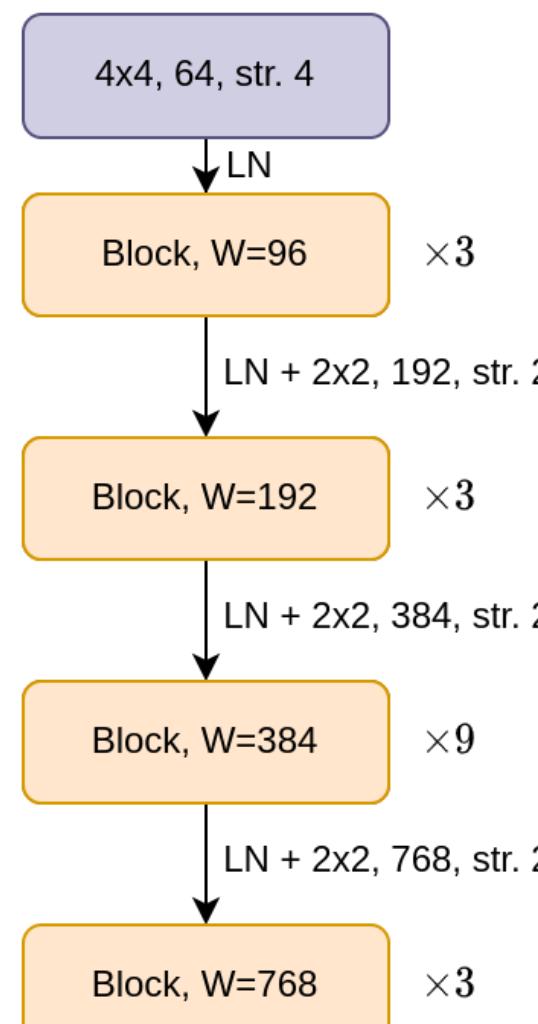
ResNeXt

Grouped Convolutions

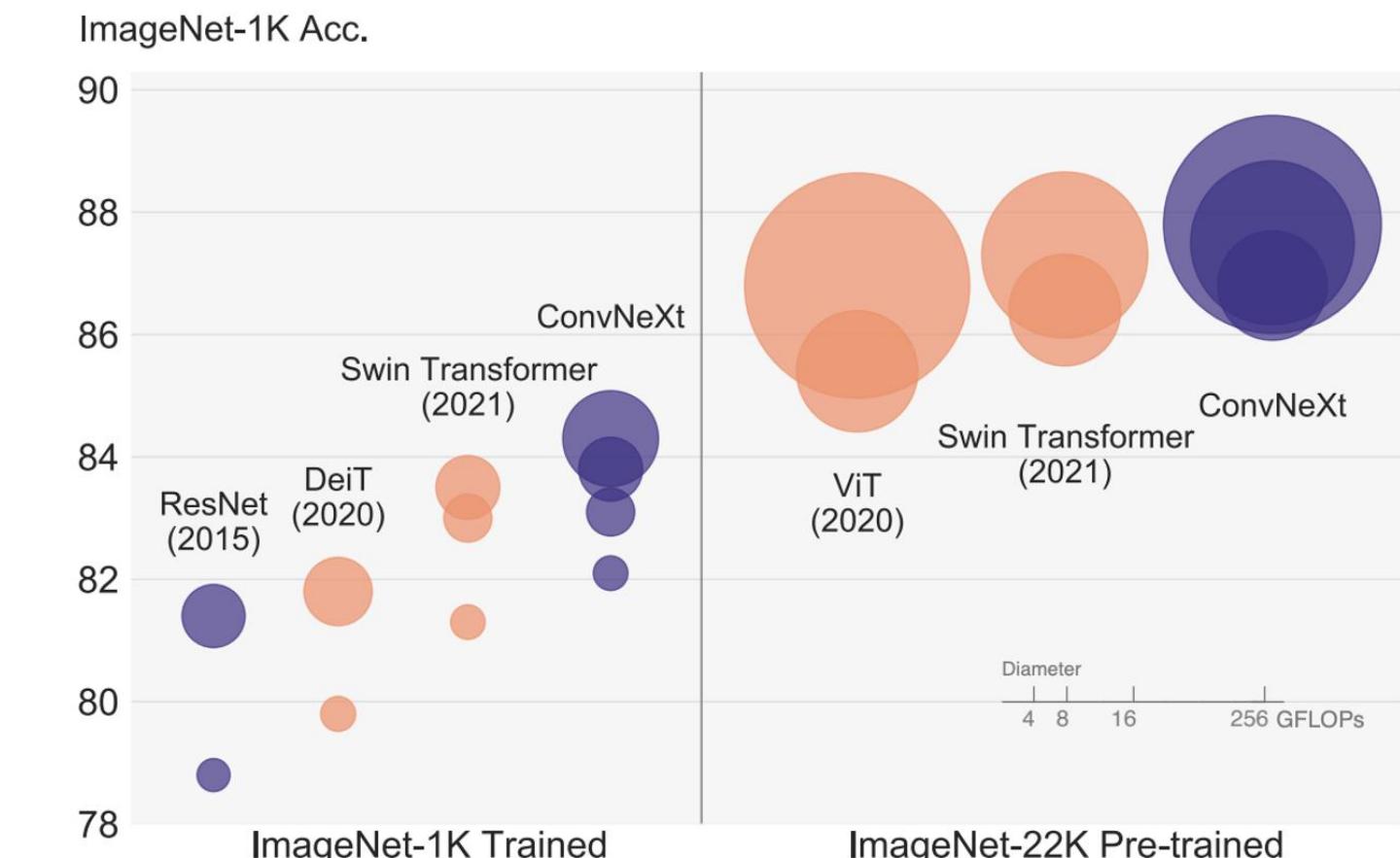


ConvNeXt

ConvNeXt Architecture



Params: 28.6M, FLOPs: 4.5G, IN-1k acc: 82.0%



ConvNeXt

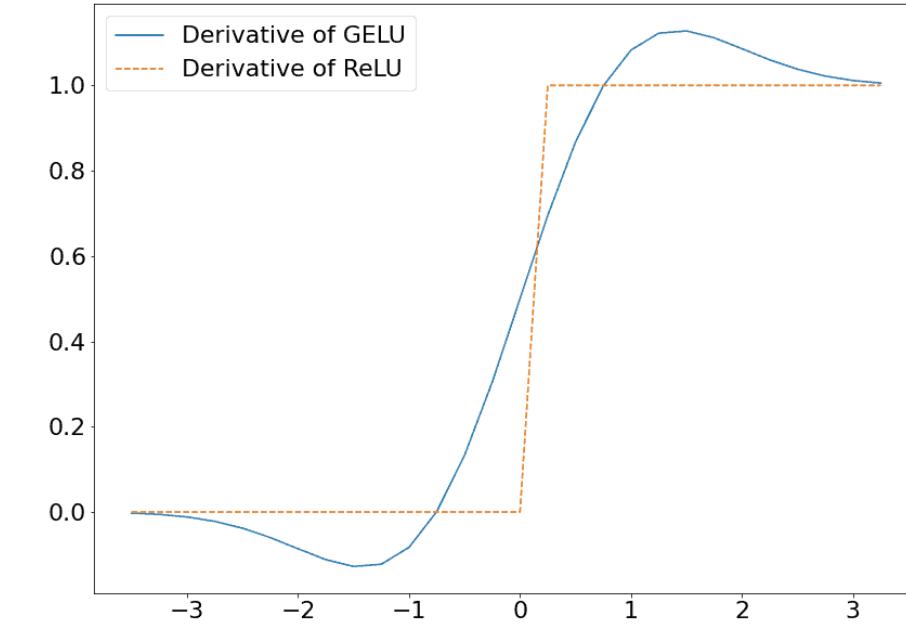
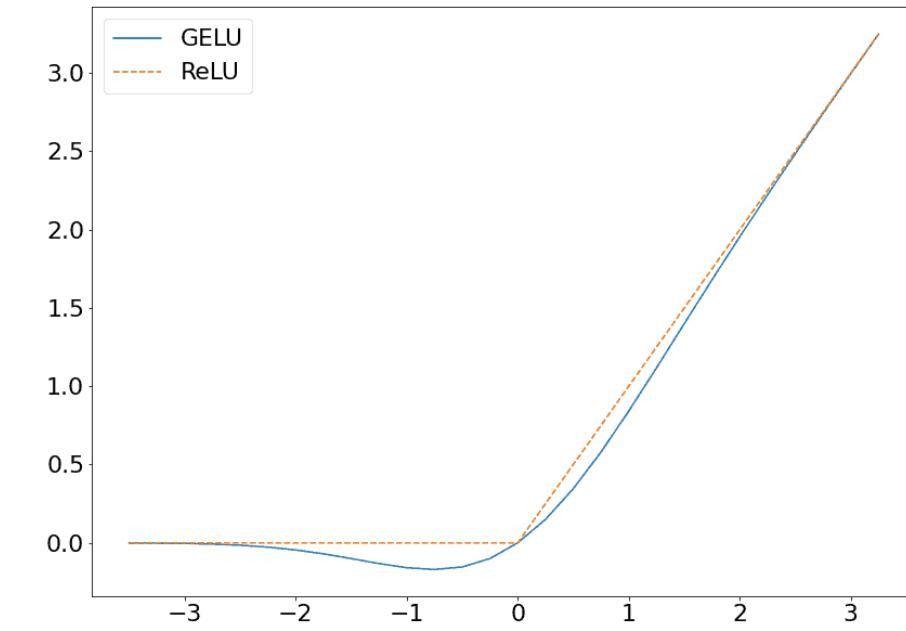
GELU activation

$$\text{GELU}(x) = xP(X \leq x) = x\Phi(x) = x \cdot \frac{1}{2} \left[1 + \text{erf}\left(\frac{x}{\sqrt{2}}\right) \right]$$

Podemos aproximar la GELU con:

$$\text{GELU}(x) \approx 0.5x \left(1 + \tanh \left[\sqrt{\frac{2}{\pi}} (x + 0.044715x^3) \right] \right)$$

$$\text{GELU}(x) \approx x\sigma(1.702x)$$

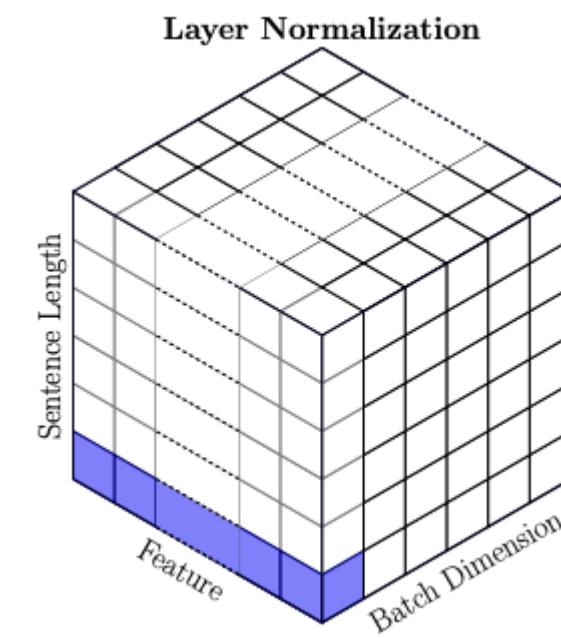
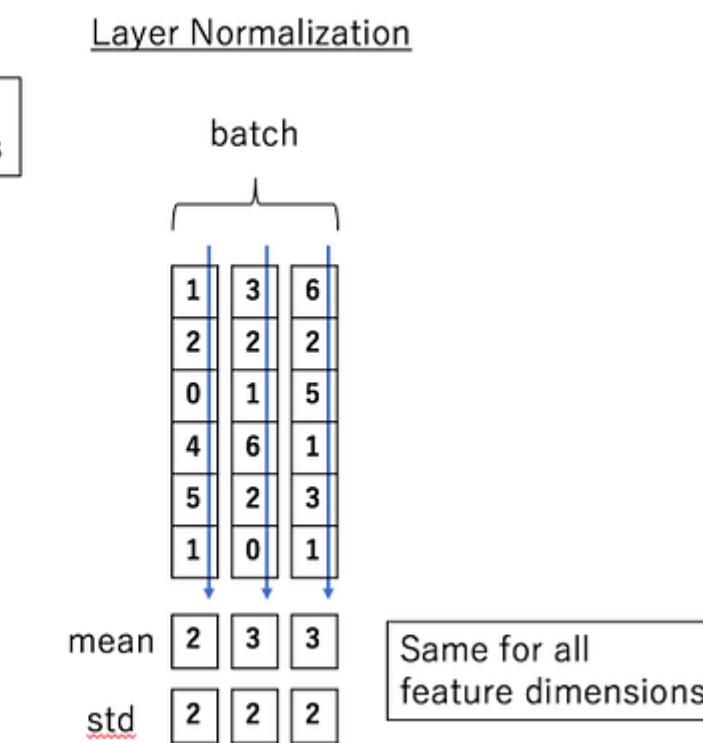
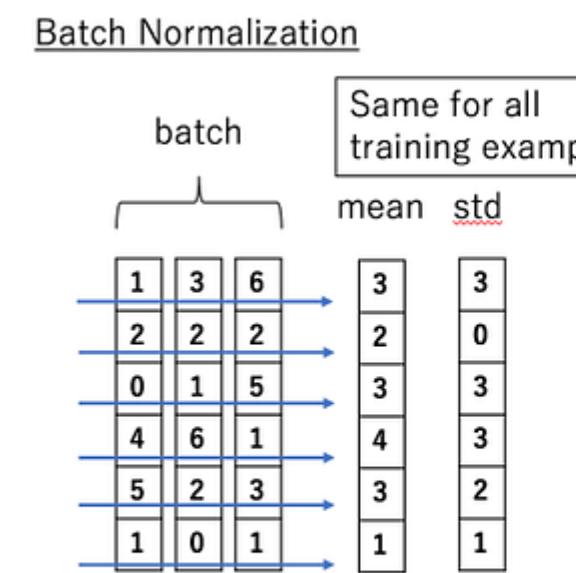
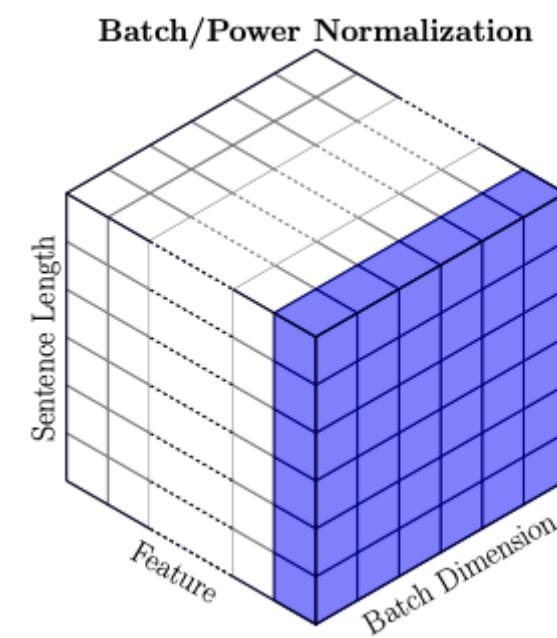


TRANSFORMATEC

Dan Hendrycks, and Kevin Gimpel (2016) "Gaussian Error Linear Units (GELUs)".
arXiv preprint arXiv:1606.08415.

ConvNeXt

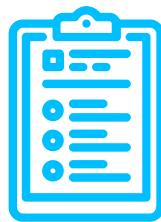
Layer Norm



$$y = \frac{x - \mathbb{E}[x]}{\sqrt{\text{Var}[x] + \varepsilon}} \cdot \gamma + \beta$$



7.



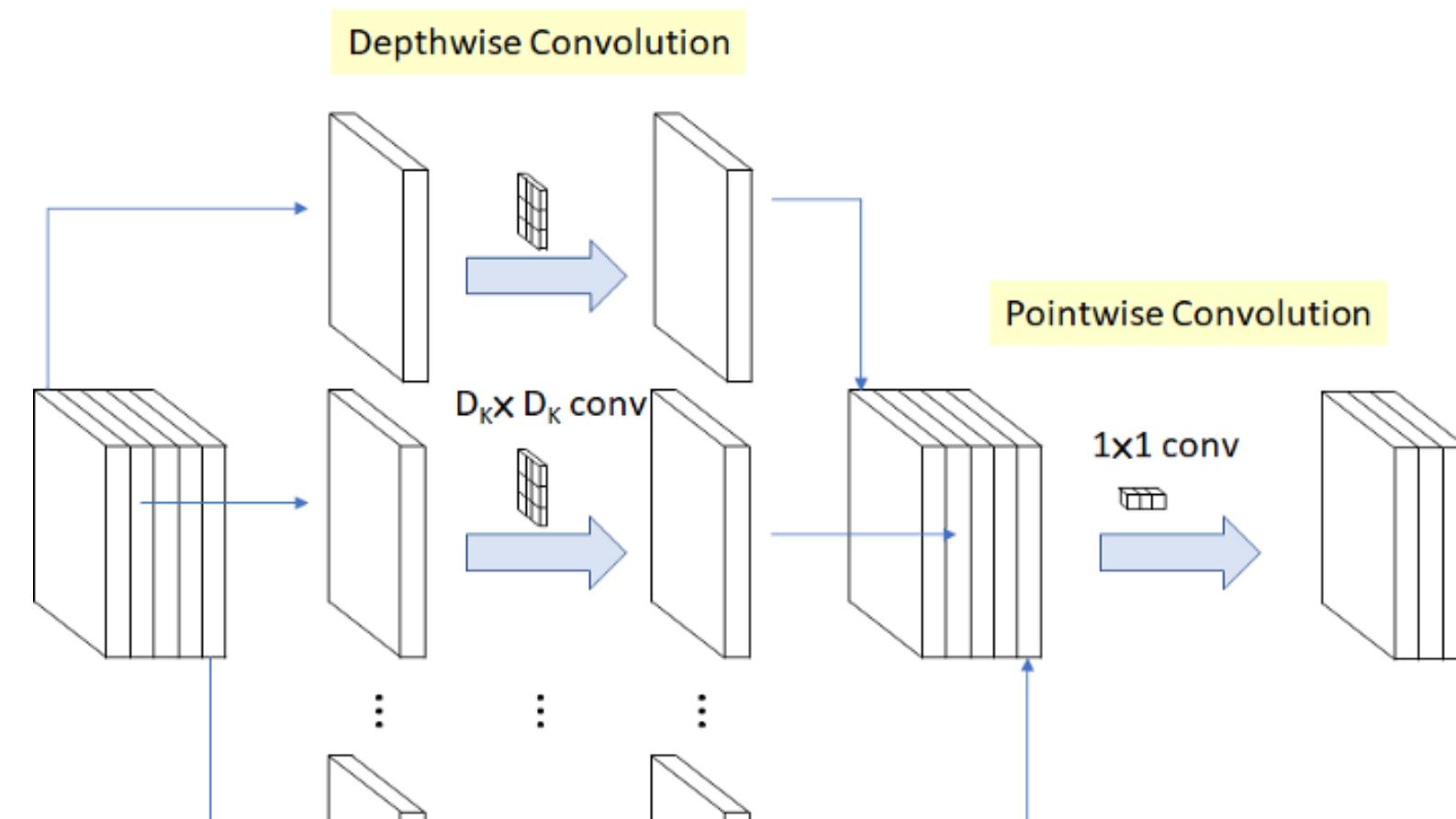
MobilNet

TRANSFORMATEC

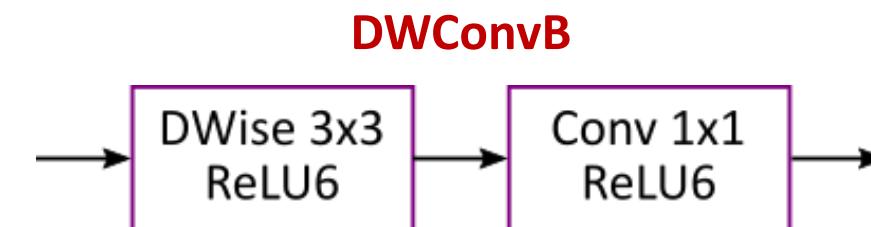
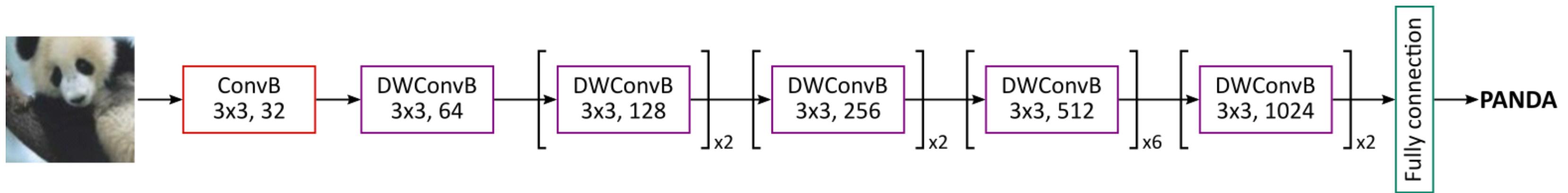
> Reinventa el mundo <



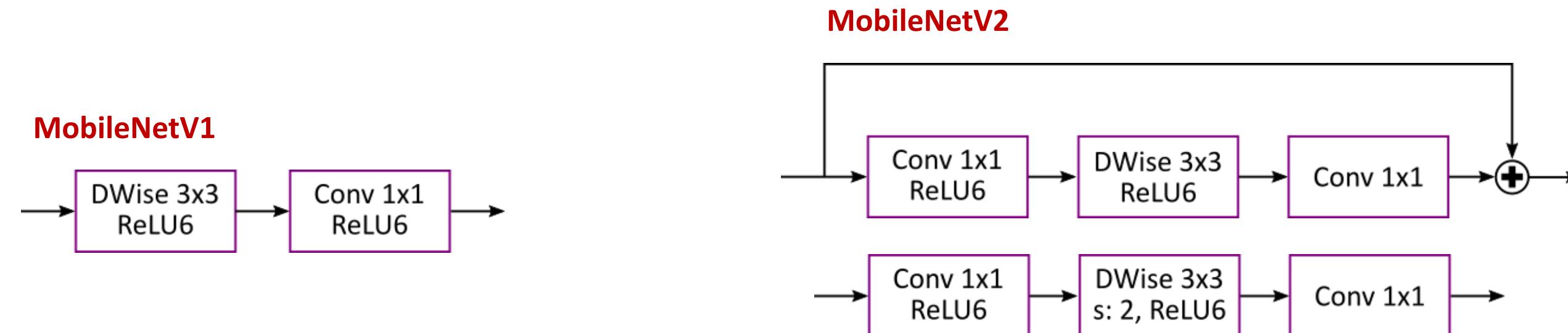
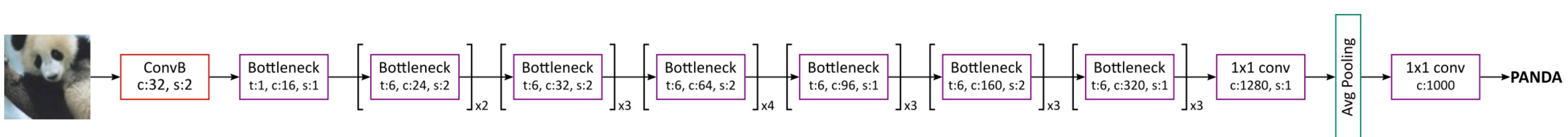
MobilNet



MobilNet



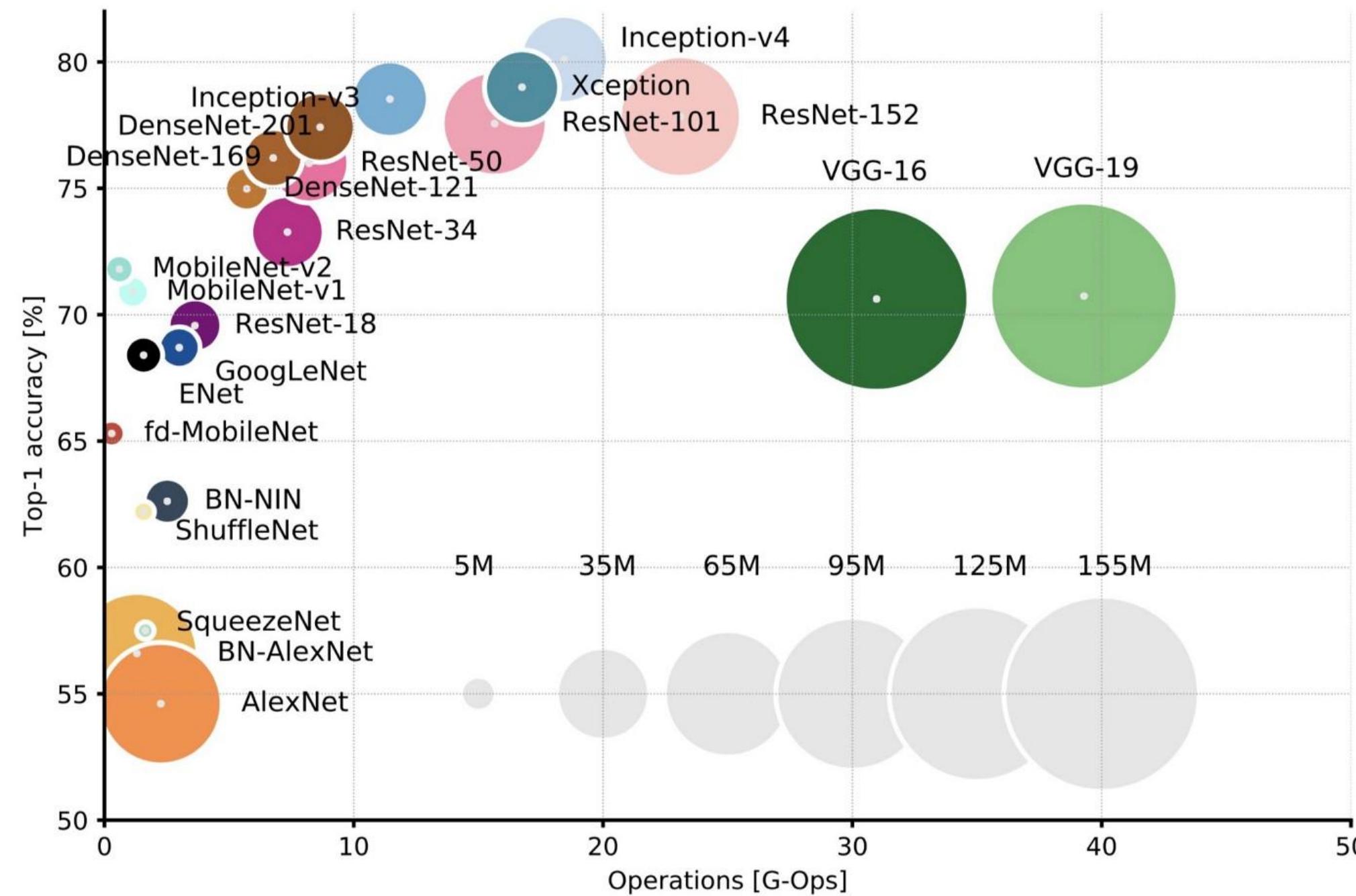
MobilNetV2



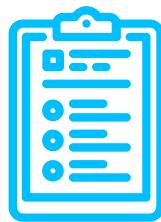
TRANSFORMATEC

Mark Sandler et al. (2018) "MobileNetV2: Inverted Residuals and Linear Bottlenecks".
Proceedings of the IEEE conference on computer vision and pattern recognition. 2018.

MobilNetV2



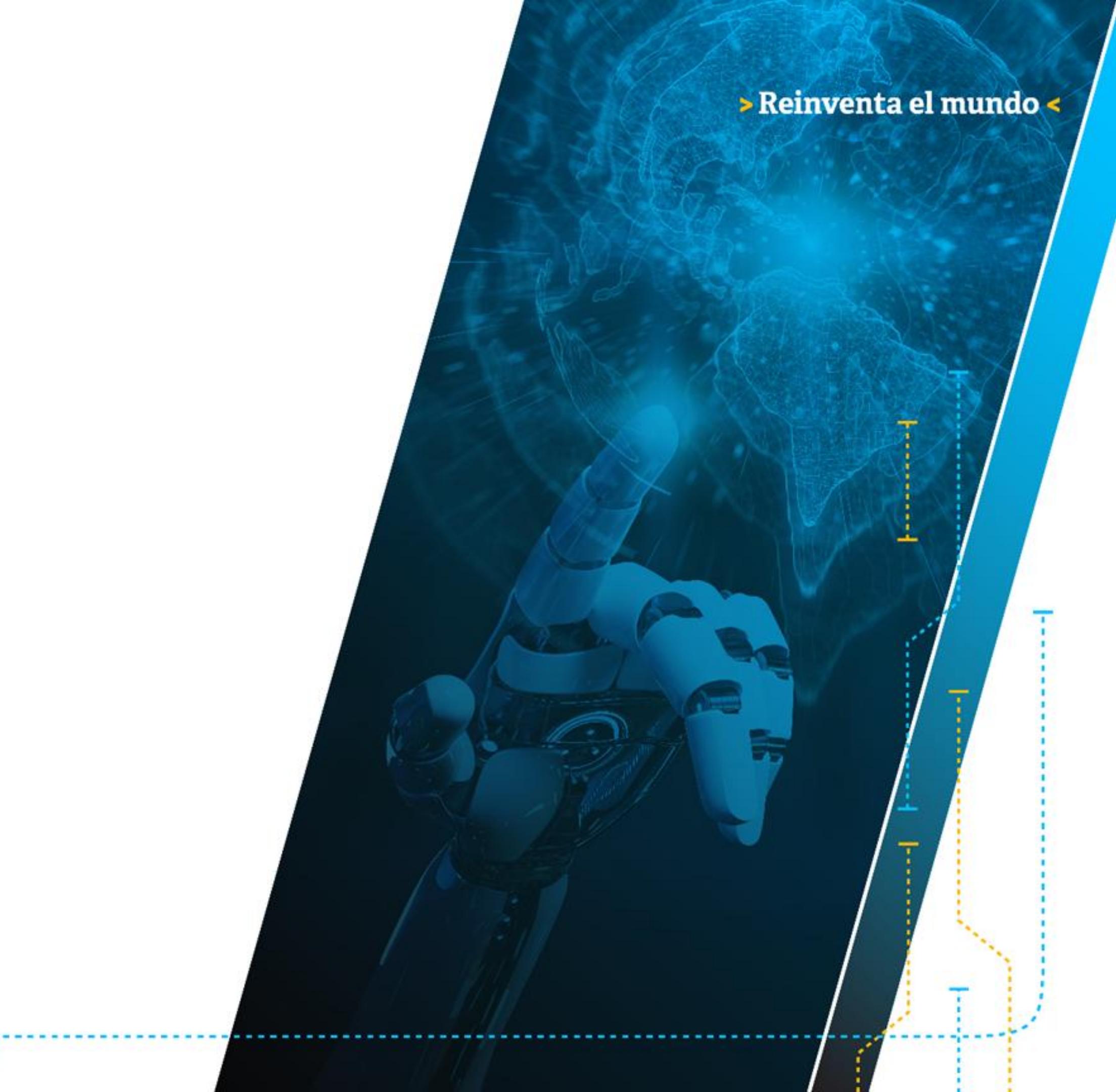
3.



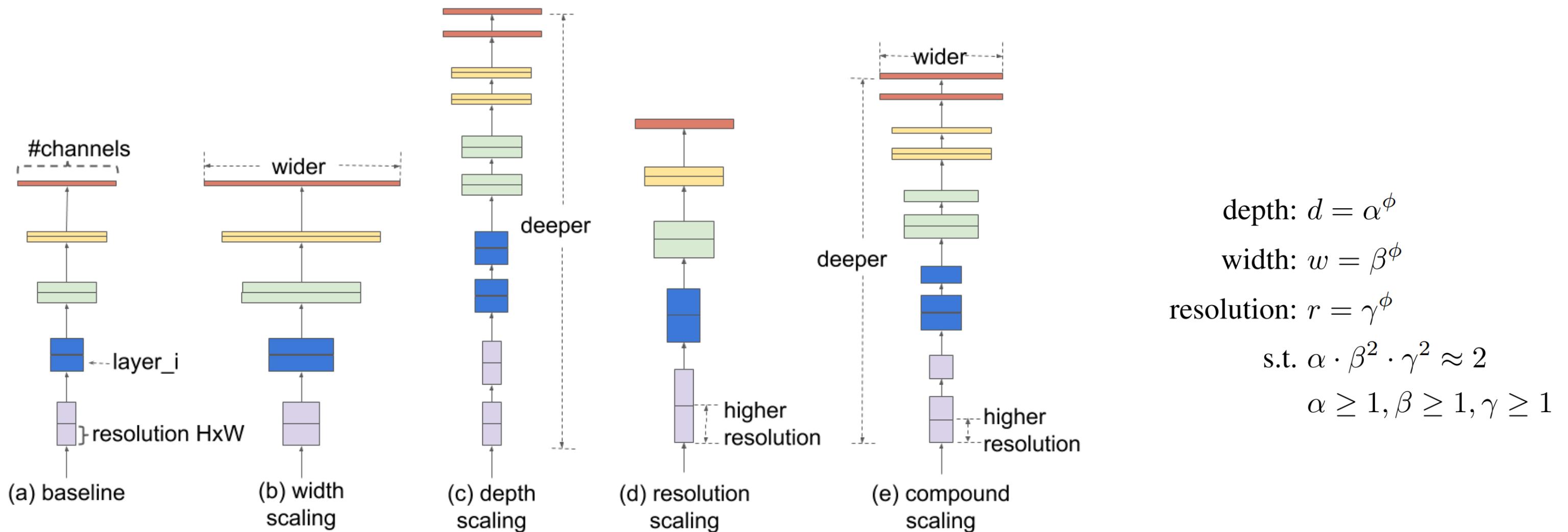
EfficientNet

TRANSFORMATEC

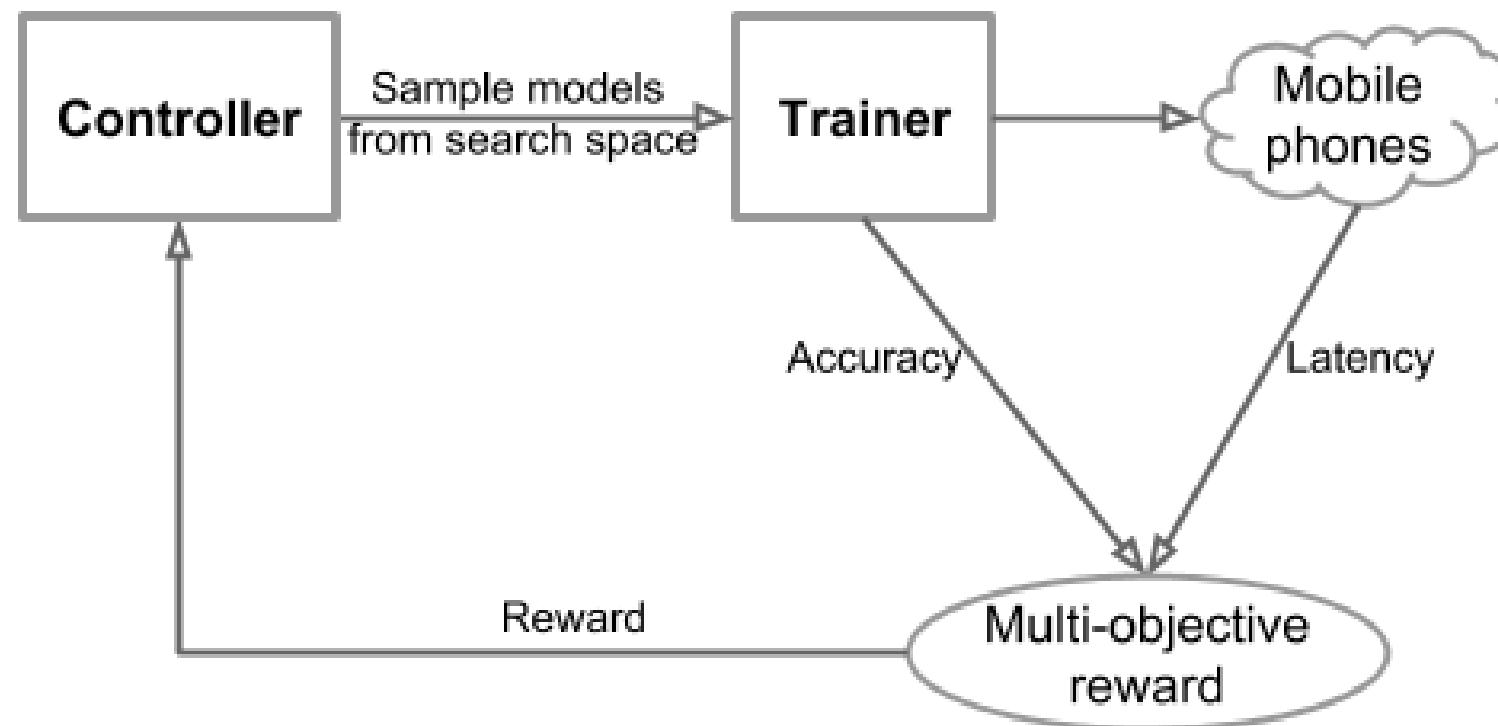
> Reinventa el mundo <



EfficientNet



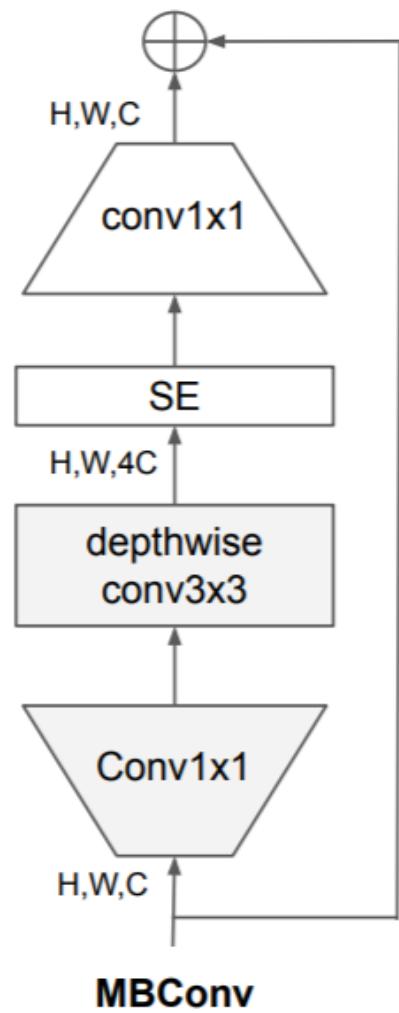
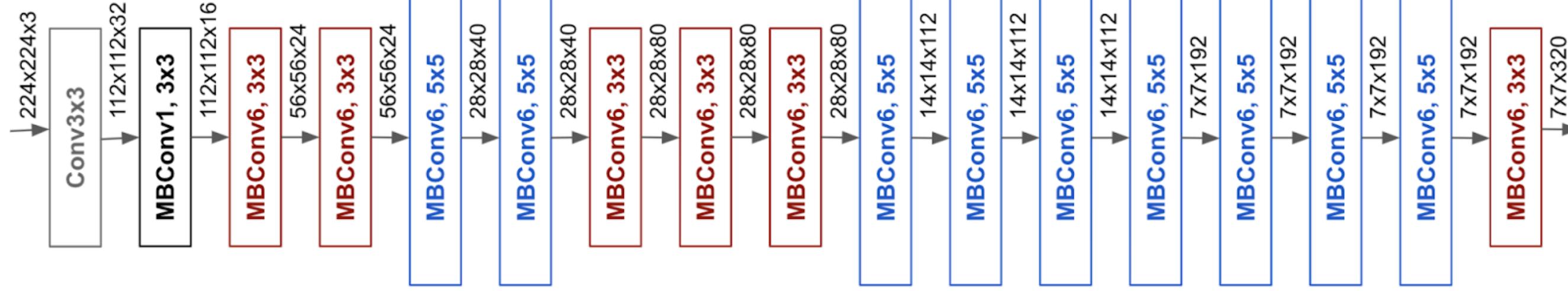
MnasNet



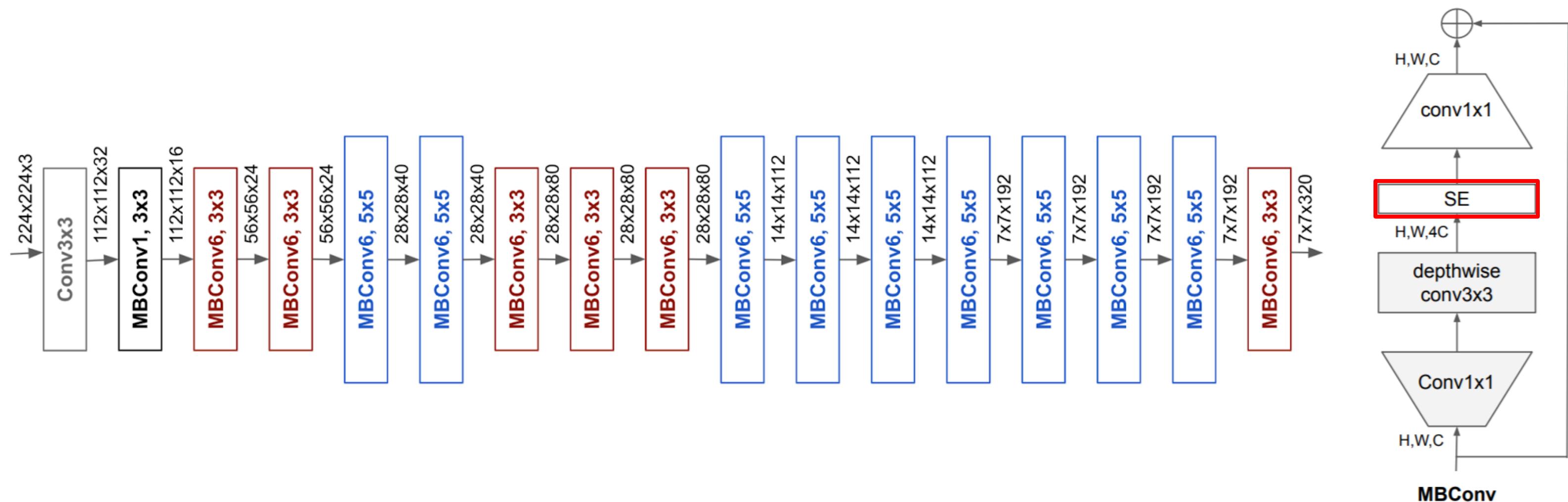
$$\begin{aligned}
 & \underset{m}{\text{maximize}} \quad ACC(m) \times \left[\frac{LAT(m)}{T} \right]^w \\
 w = & \begin{cases} \alpha, & \text{if } LAT(m) \leq T \\ \beta, & \text{otherwise} \end{cases}
 \end{aligned}$$



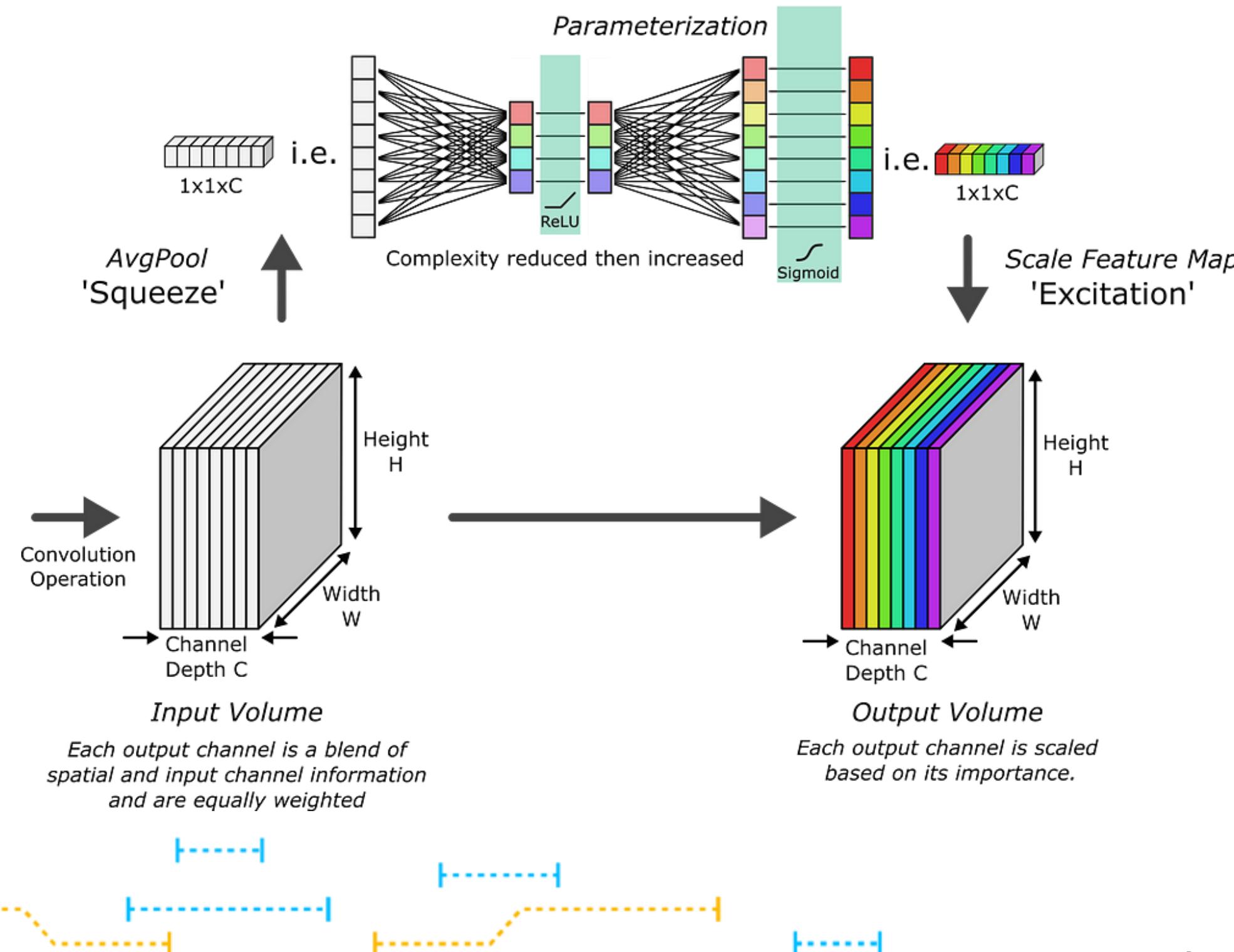
EfficientNet



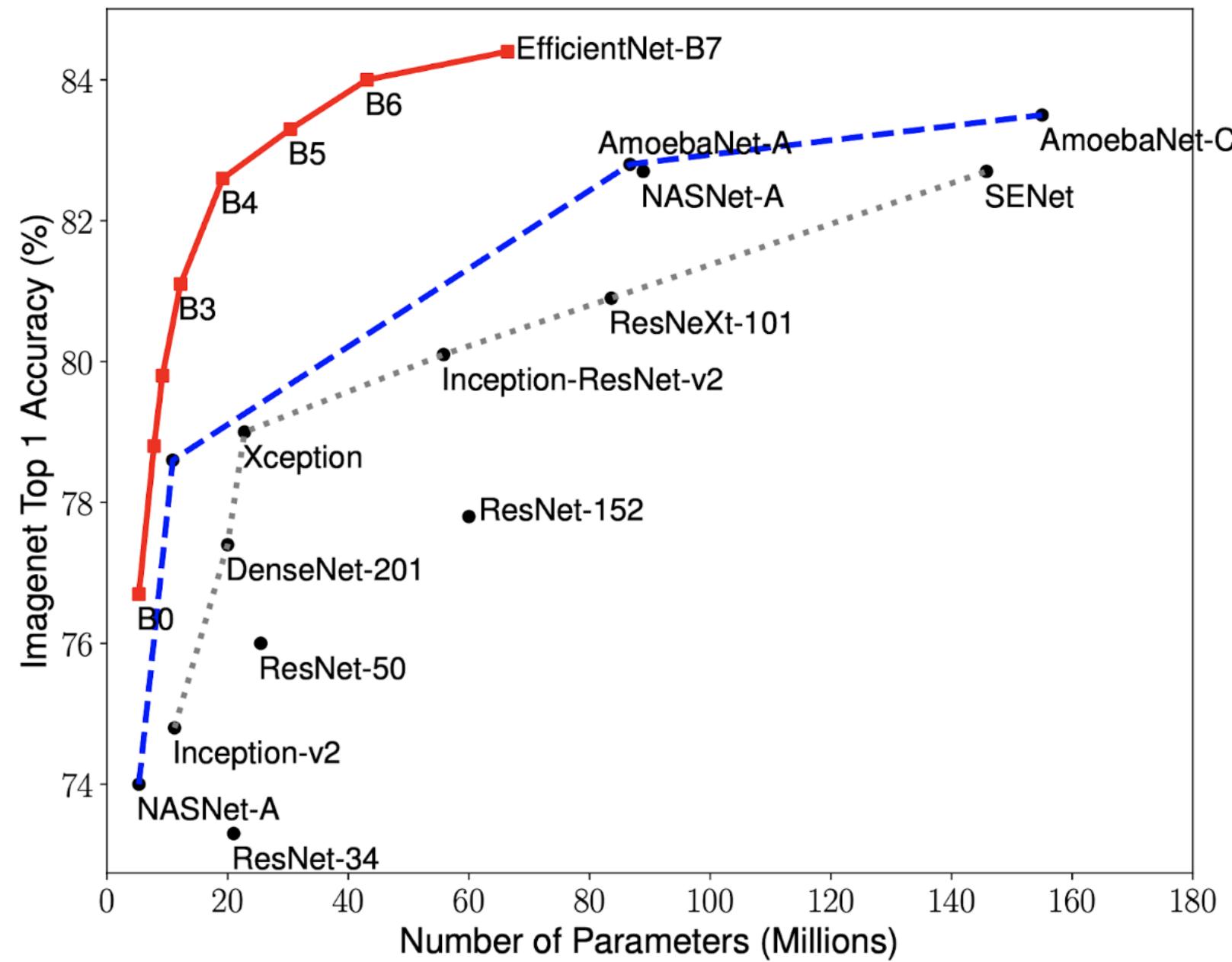
EfficientNet



Squeeze-and-Excitation

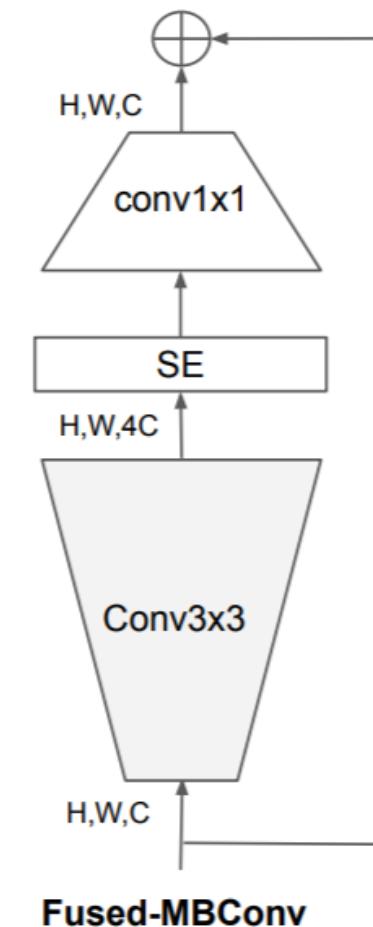
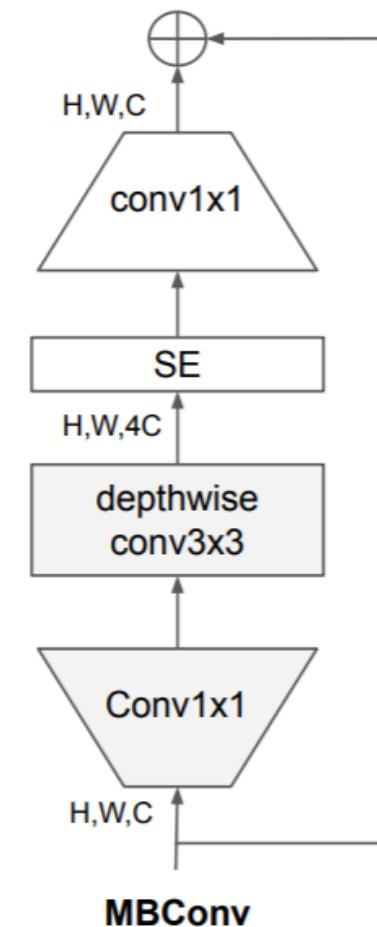


EfficientNet



EfficientNetV2

Stage	Operator	Stride	#Channels	#Layers
0	Conv3x3	2	24	1
1	Fused-MBConv1, k3x3	1	24	2
2	Fused-MBConv4, k3x3	2	48	4
3	Fused-MBConv4, k3x3	2	64	4
4	MBConv4, k3x3, SE0.25	2	128	6
5	MBConv6, k3x3, SE0.25	1	160	9
6	MBConv6, k3x3, SE0.25	2	272	15
7	Conv1x1 & Pooling & FC	-	1792	1

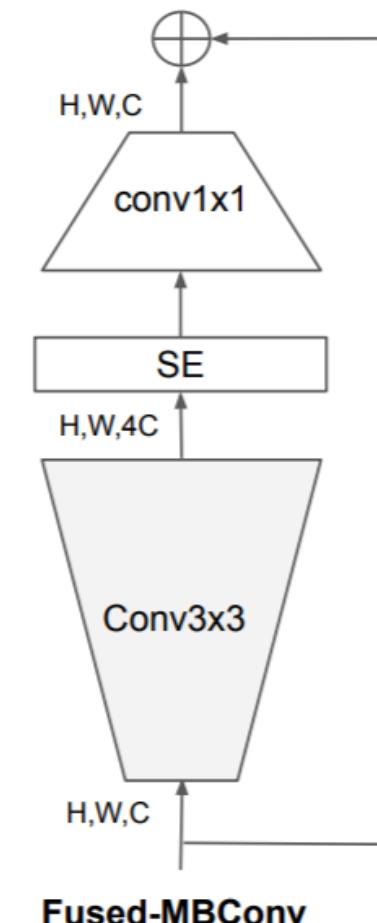
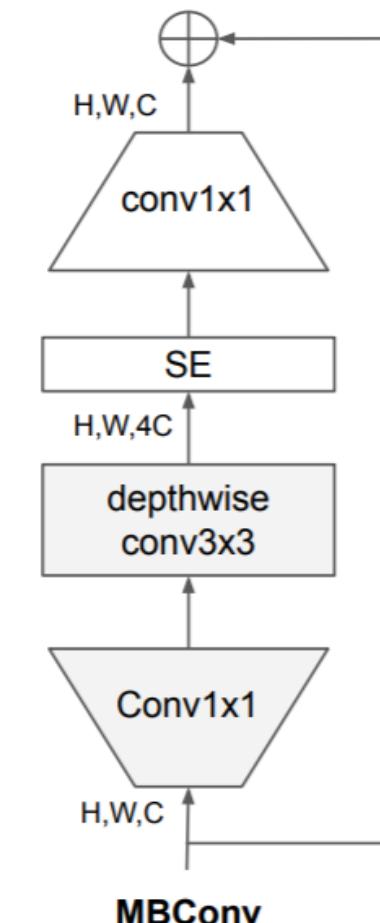


TRANSFORMATEC

Mingxing Tan et al. (2021) "EfficientNetV2: Smaller Models and Faster Training".
arXiv preprint arXiv:2104.00298.

EfficientNetV2

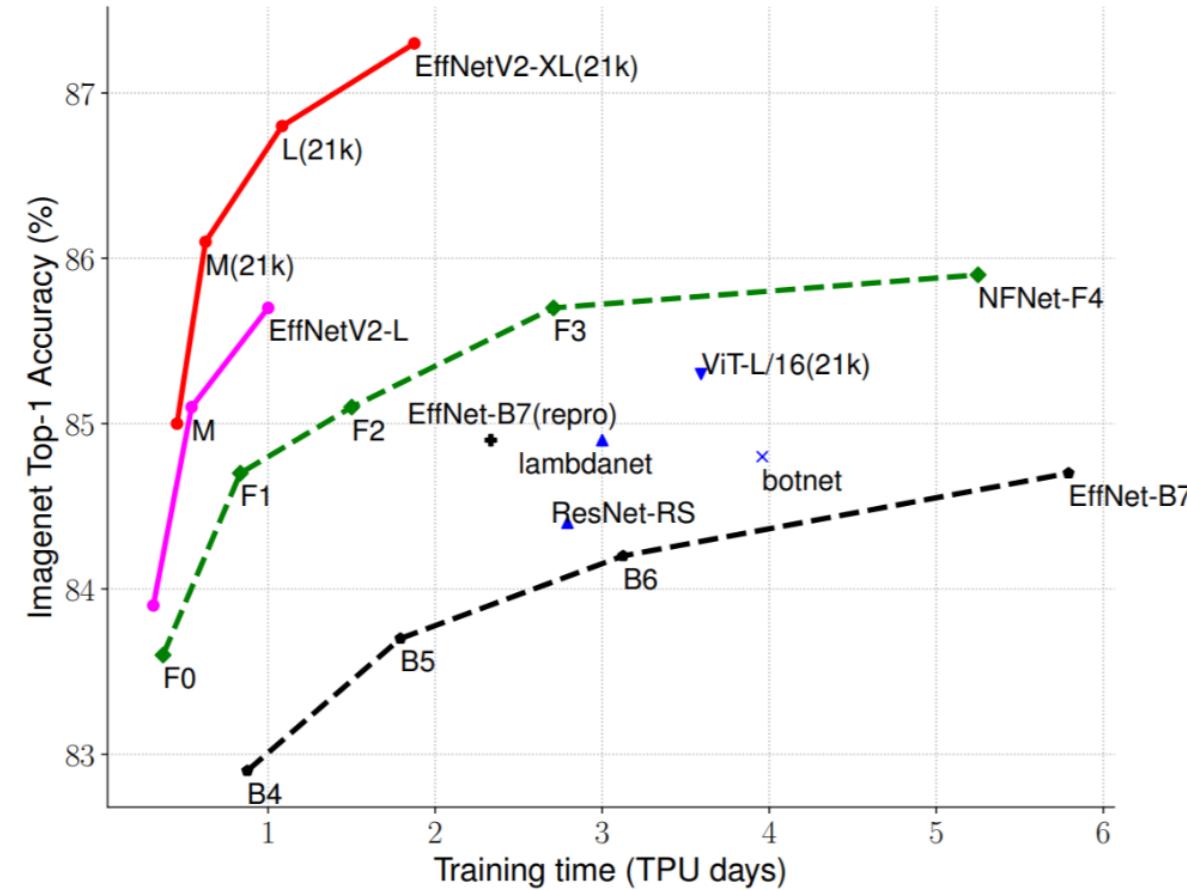
	Params (M)	FLOPs (B)	Top-1 Acc.	TPU imgs/sec/core	V100 imgs/sec/gpu
No fused	19.3	4.5	82.8%	262	155
Fused stage1-3	20.0	7.5	83.1%	362	216
Fused stage1-5	43.4	21.3	83.1%	327	223
Fused stage1-7	132.0	34.4	81.7%	254	206



TRANSFORMATEC

Mingxing Tan et al. (2021) "EfficientNetV2: Smaller Models and Faster Training".
arXiv preprint arXiv:2104.00298.

EfficientNetV2



(a) Training efficiency.

	EfficientNet (2019)	ResNet-RS (2021)	DeiT/ViT (2021)	EfficientNetV2 (ours)
Top-1 Acc.	84.3%	84.0%	83.1%	83.9%
Parameters	43M	164M	86M	24M

(b) Parameter efficiency.



GRACIAS

Victor Flores Benites