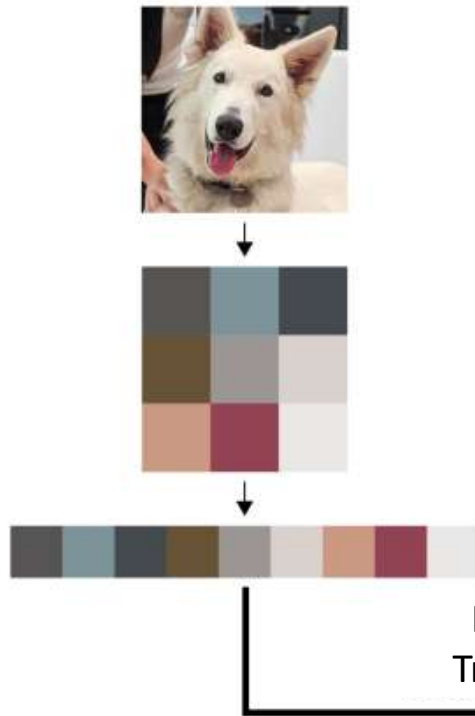


# Visión Computacional

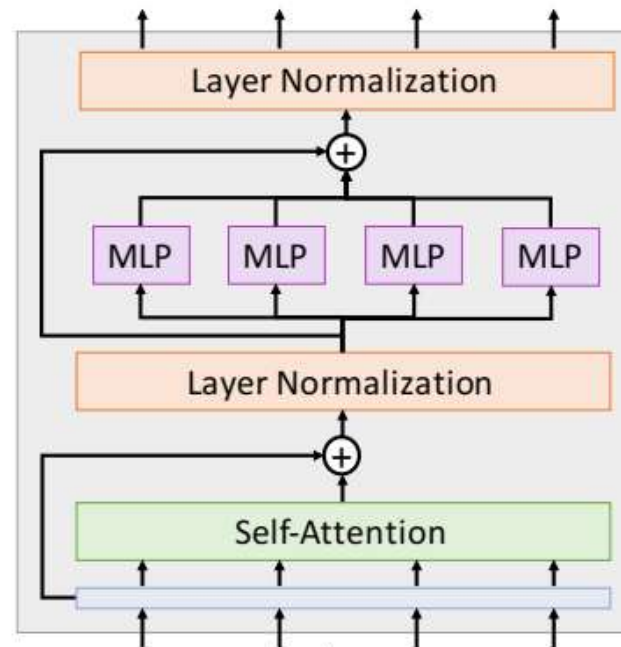
Ivan Sipiran

# Idea muy básica

Tratar una imagen como un conjunto de píxeles

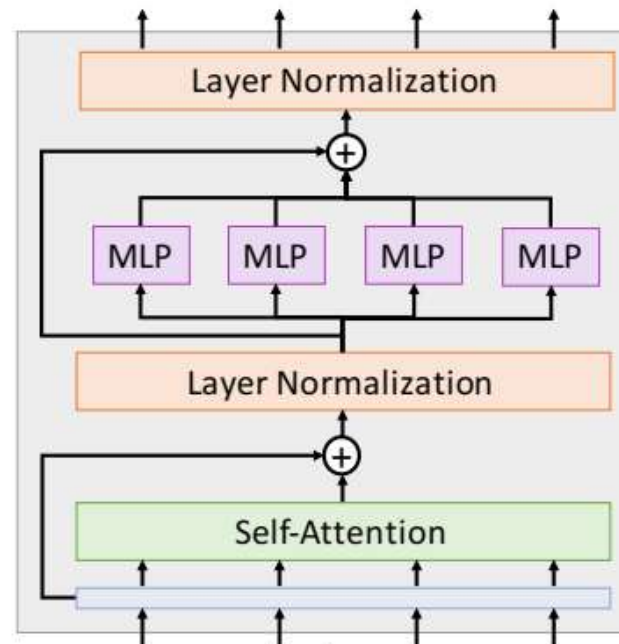
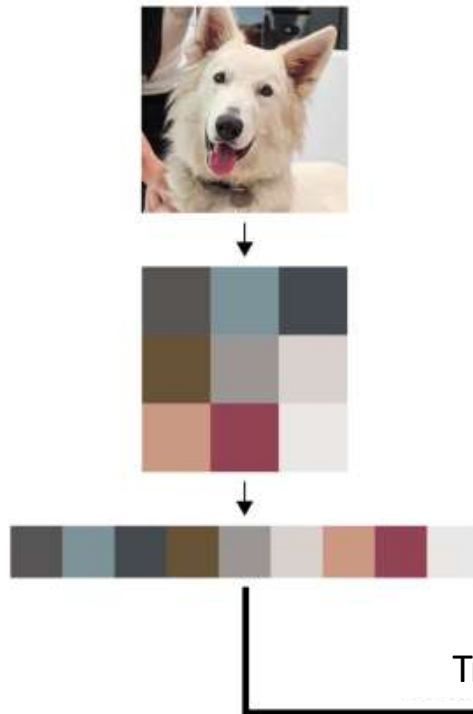


Pixels alimentan un Transformer estándar



# Idea muy básica

Tratar una imagen como un conjunto de píxeles

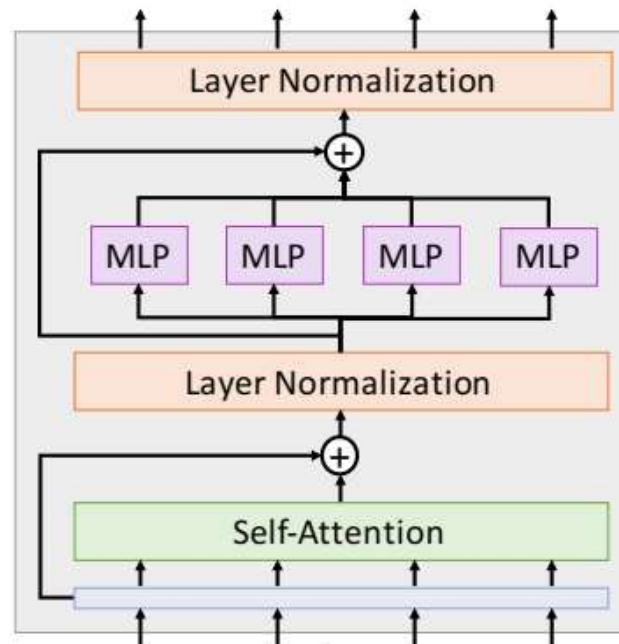
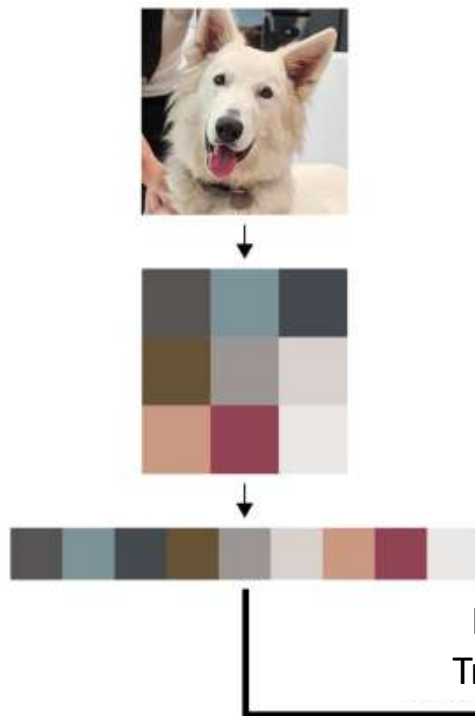


Problema: memoria!

Imagen de  $R \times R$   
necesita  $R^4$   
elementos en  
matriz de atención

# Idea muy básica

Tratar una imagen como un conjunto de píxeles

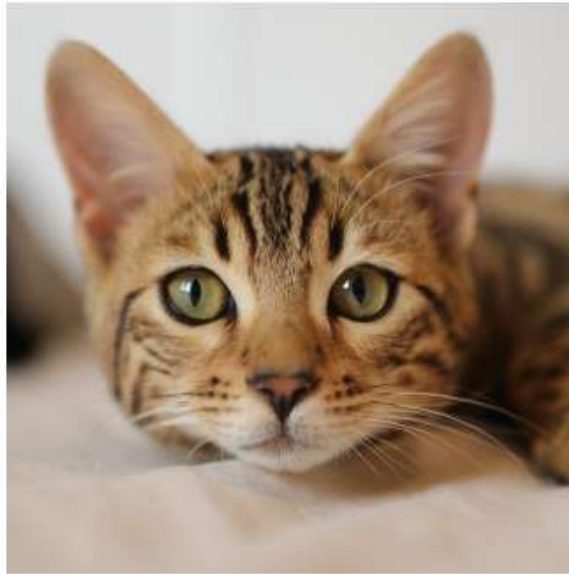


Problema: memoria!

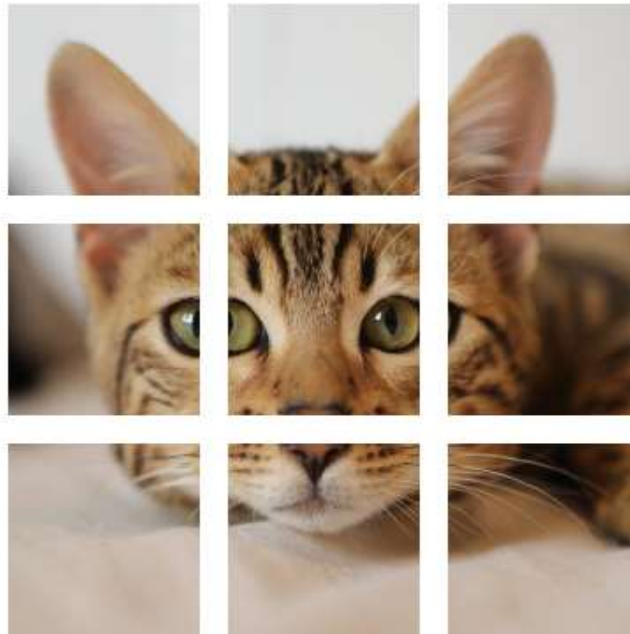
Imagen de  $R \times R$   
necesita  $R^4$   
elementos en  
matriz de atención

$R=128$ , 48 capas, 16  
cabezas de atención:  
768GB para matrices  
de atención de una  
sola imagen

Transformer en parches



# Transformer en parches



# Transformer en parches

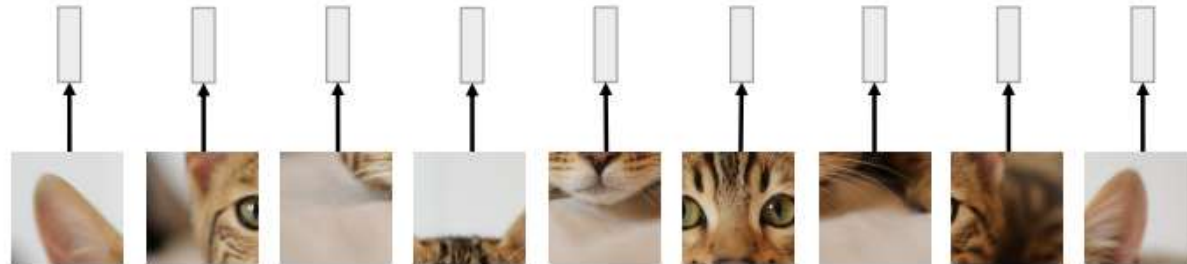
$N$  parches de entrada  
Tamaño: 3 x 16 x 16



# Transformer en parches

Proyección a vector de  
D dimensiones

$N$  parches de entrada  
Tamaño:  $3 \times 16 \times 16$



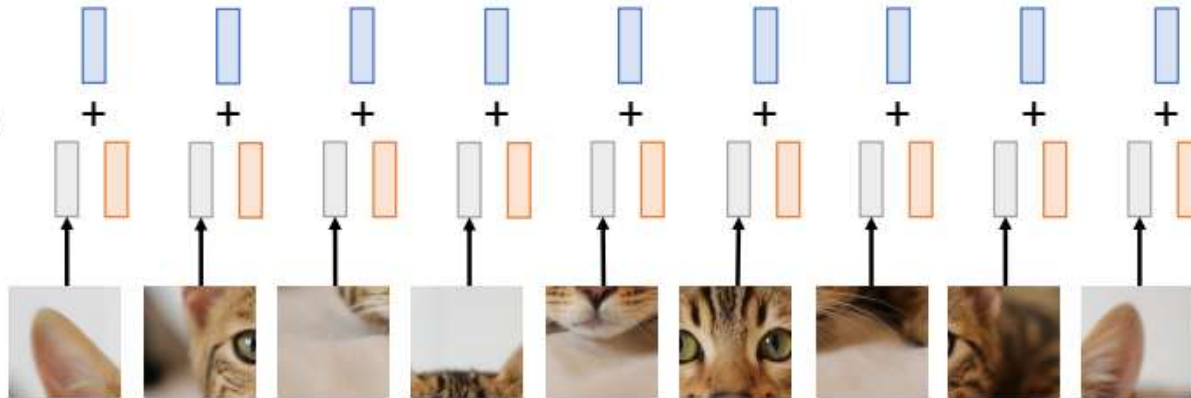


# Transformer en parches

Positional embedding

Proyección a vector de  $D$  dimensiones

$N$  parches de entrada  
Tamaño:  $3 \times 16 \times 16$

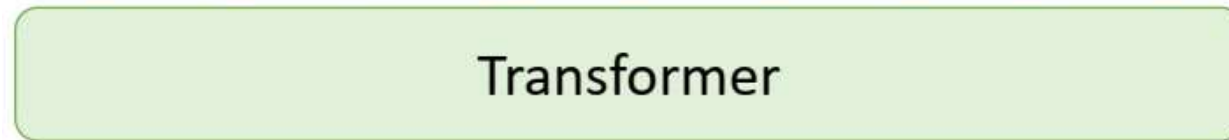


# Transformer en parches

Vector de salida



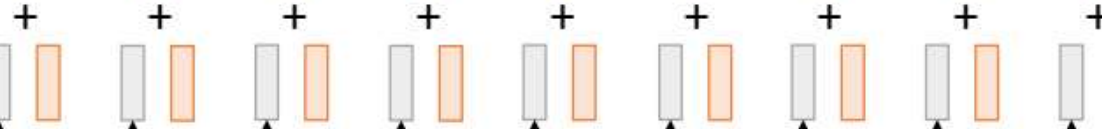
El mismo Transformer que NLP!



Positional embedding



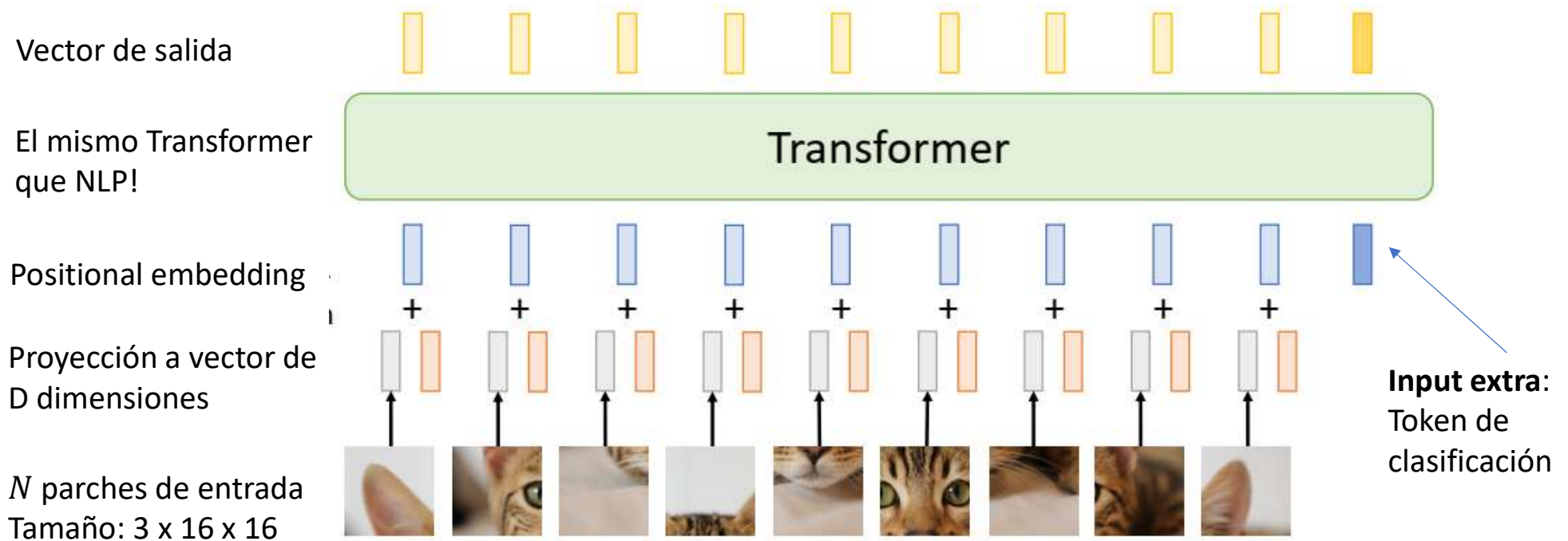
Proyección a vector de D dimensiones



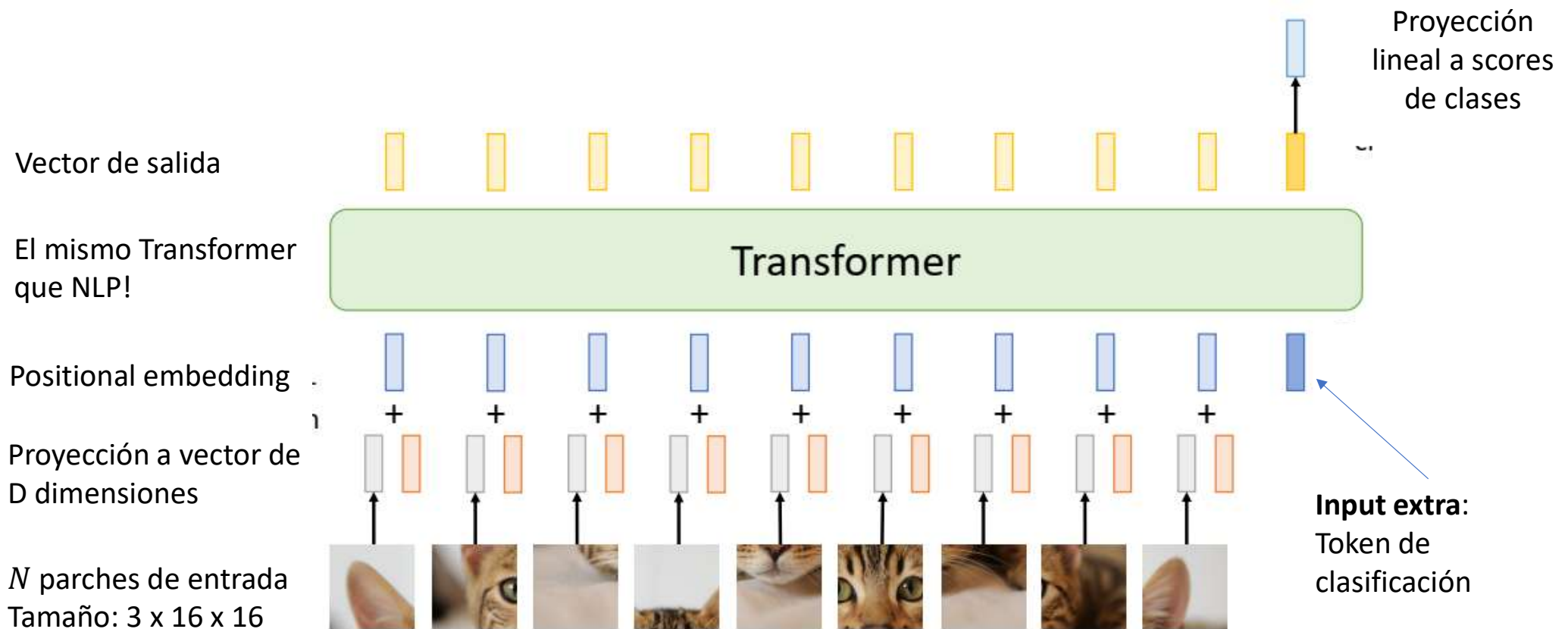
$N$  parches de entrada  
Tamaño: 3 x 16 x 16



# Transformer en parches

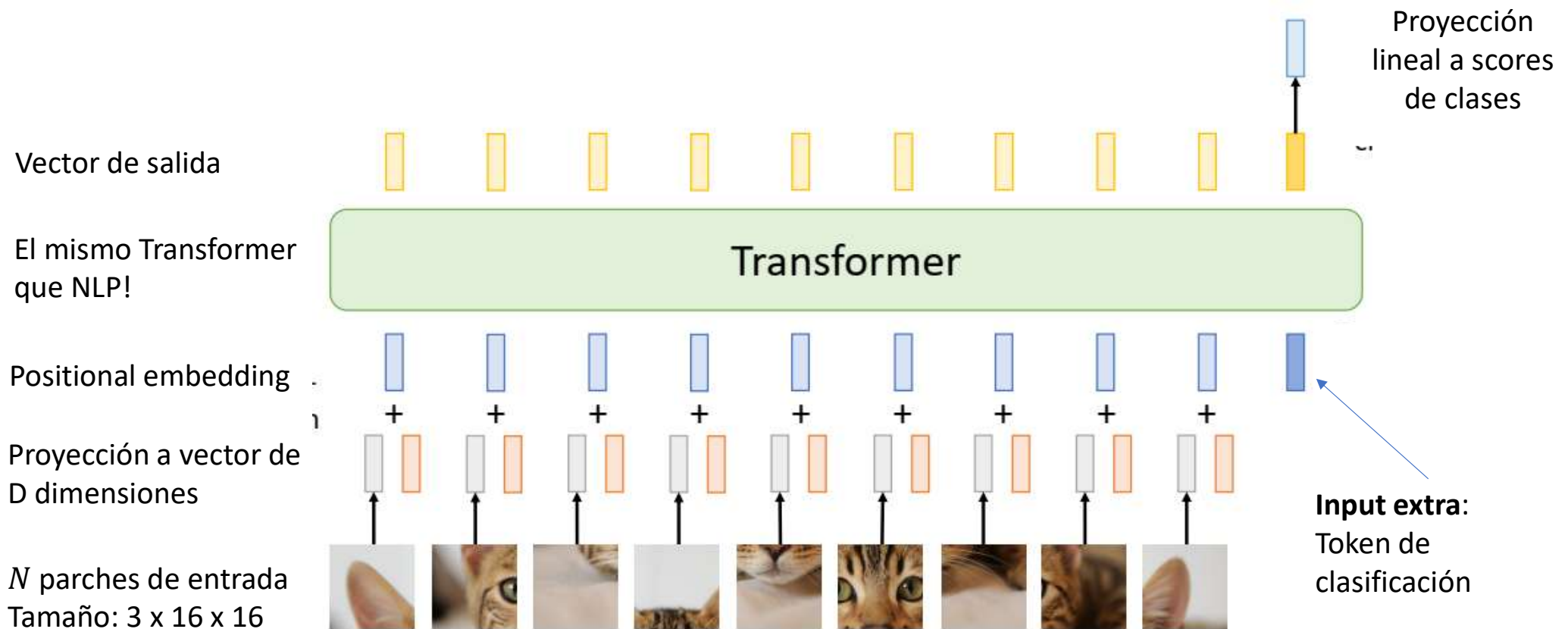


# Transformer en parches



# Transformer en parches

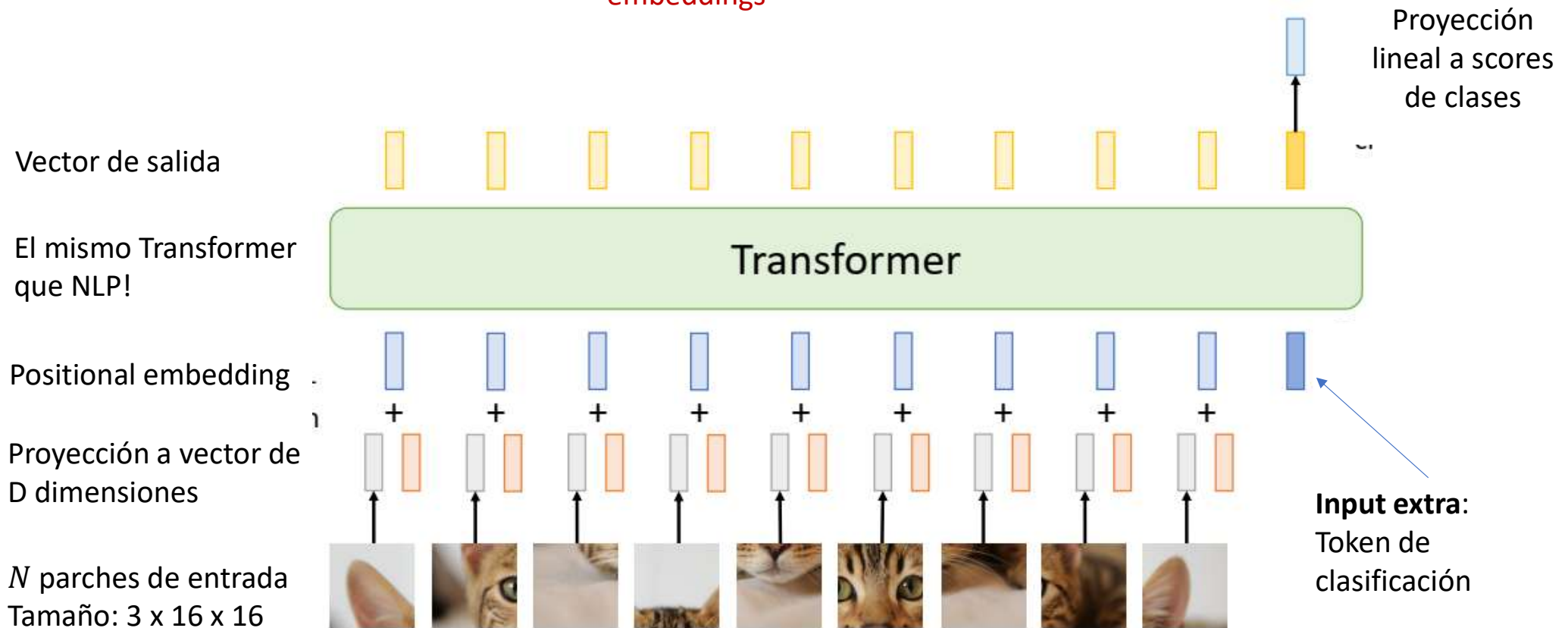
Modelo de visión sin convoluciones!



# Transformer en parches

Modelo de visión sin convoluciones!

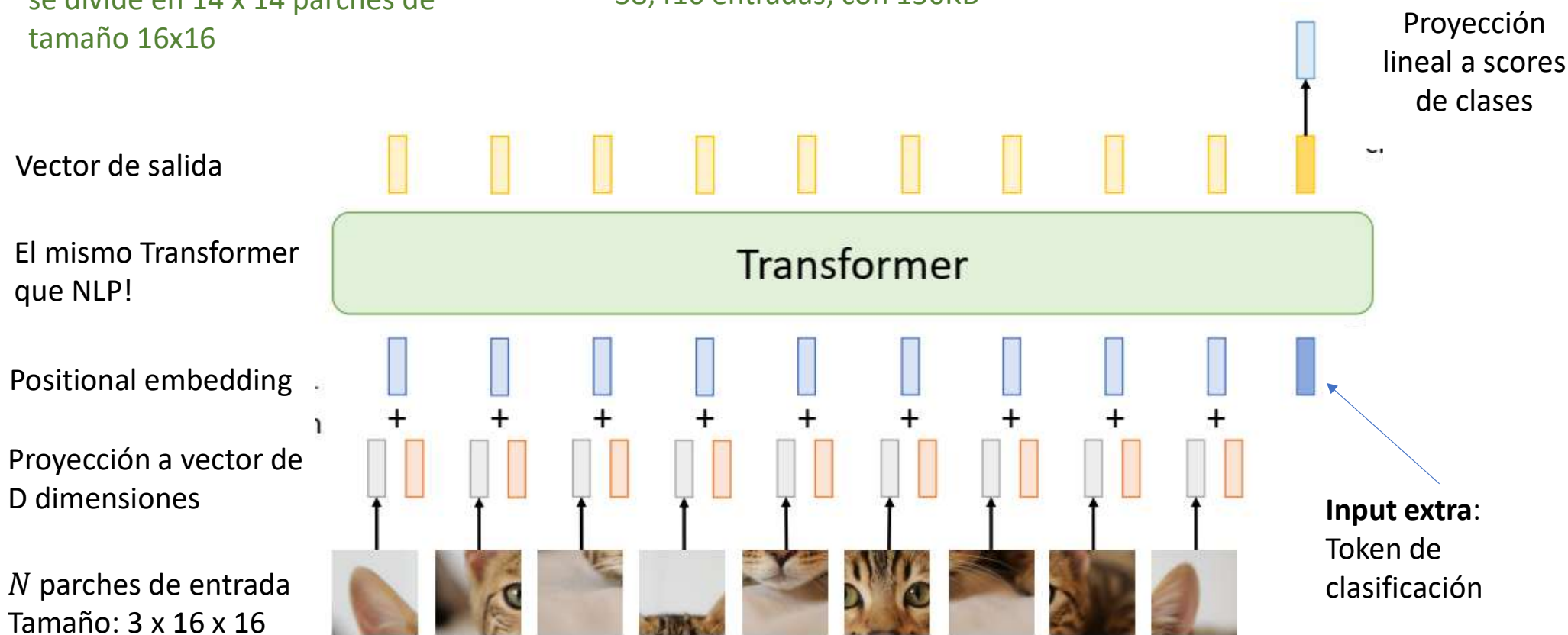
No del todo: la primera capa es Conv2D para embeddings



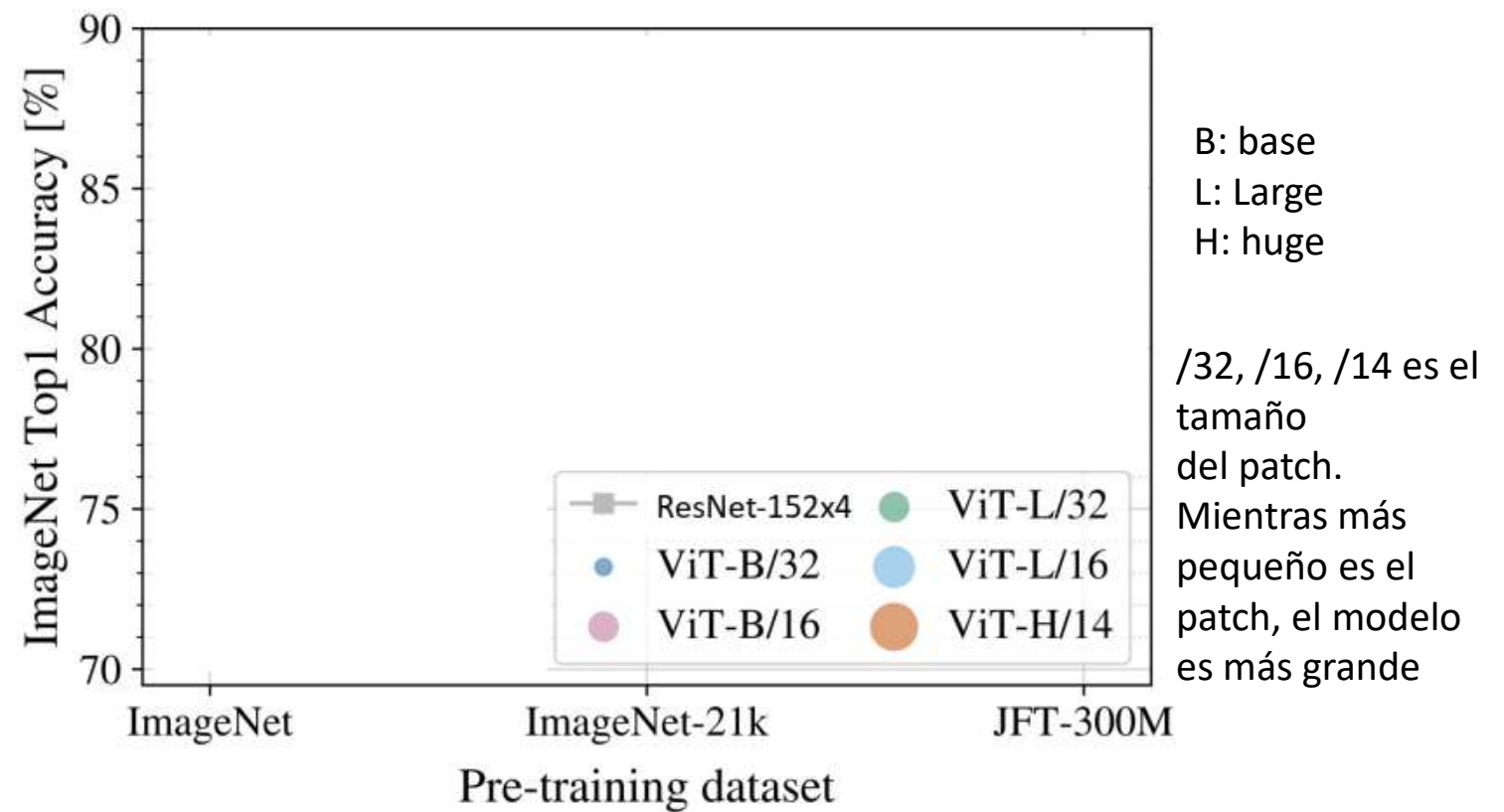
# Transformer en parches

En la práctica: Imagen de 224x224,  
se divide en 14 x 14 parches de  
tamaño 16x16

Cada matriz de atención tiene  
38,416 entradas, con 150KB



# ViT vs. ResNets

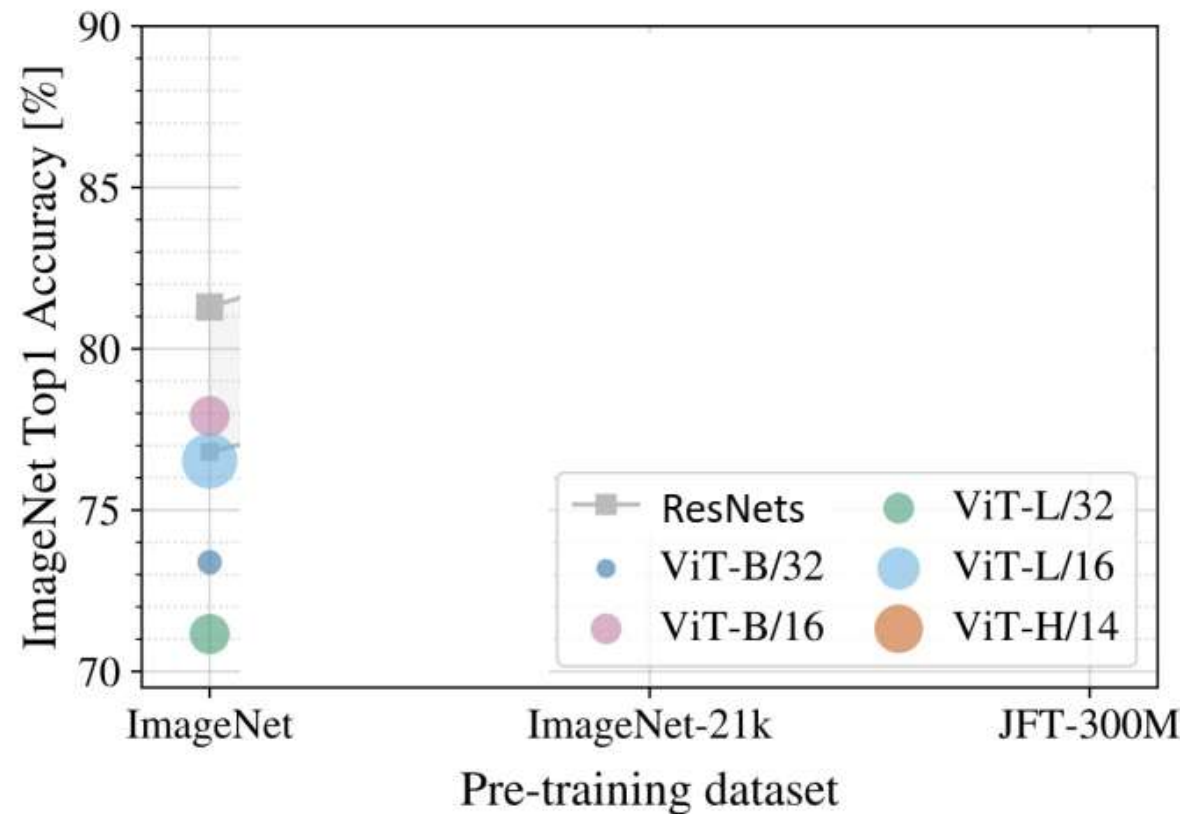




# ViT vs. ResNets

ImageNet tiene 1k clases  
y 1.2M imágenes

Entrenado en ImageNet  
ViT no es mejor que  
ResNets



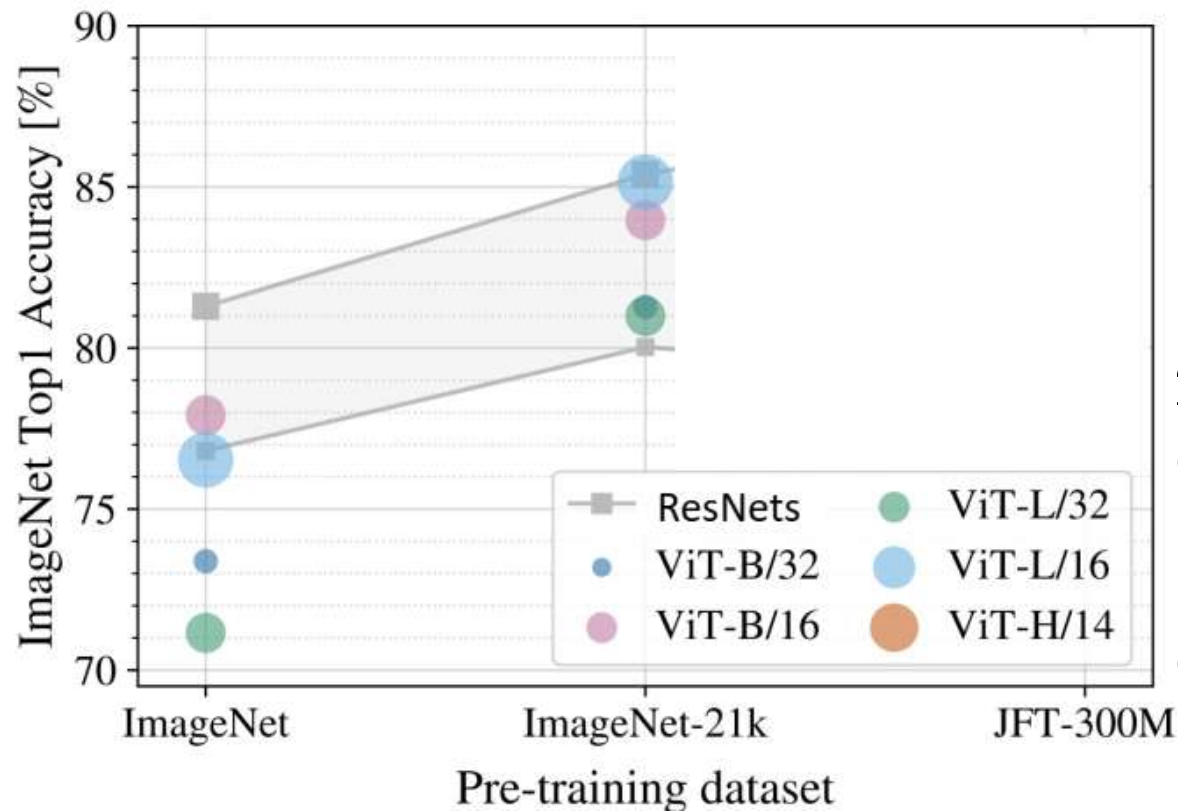
B: base  
L: Large  
H: huge

/32, /16, /14 es el  
tamaño  
del patch.  
Mientras más  
pequeño es el  
patch, el modelo  
es más grande

# ViT vs. ResNets

ImageNet-21k tiene 21k clases y 14M imágenes

Pre-entrenamiento con ImageNet-21k hace que ViT alcance el performance de ResNets grandes



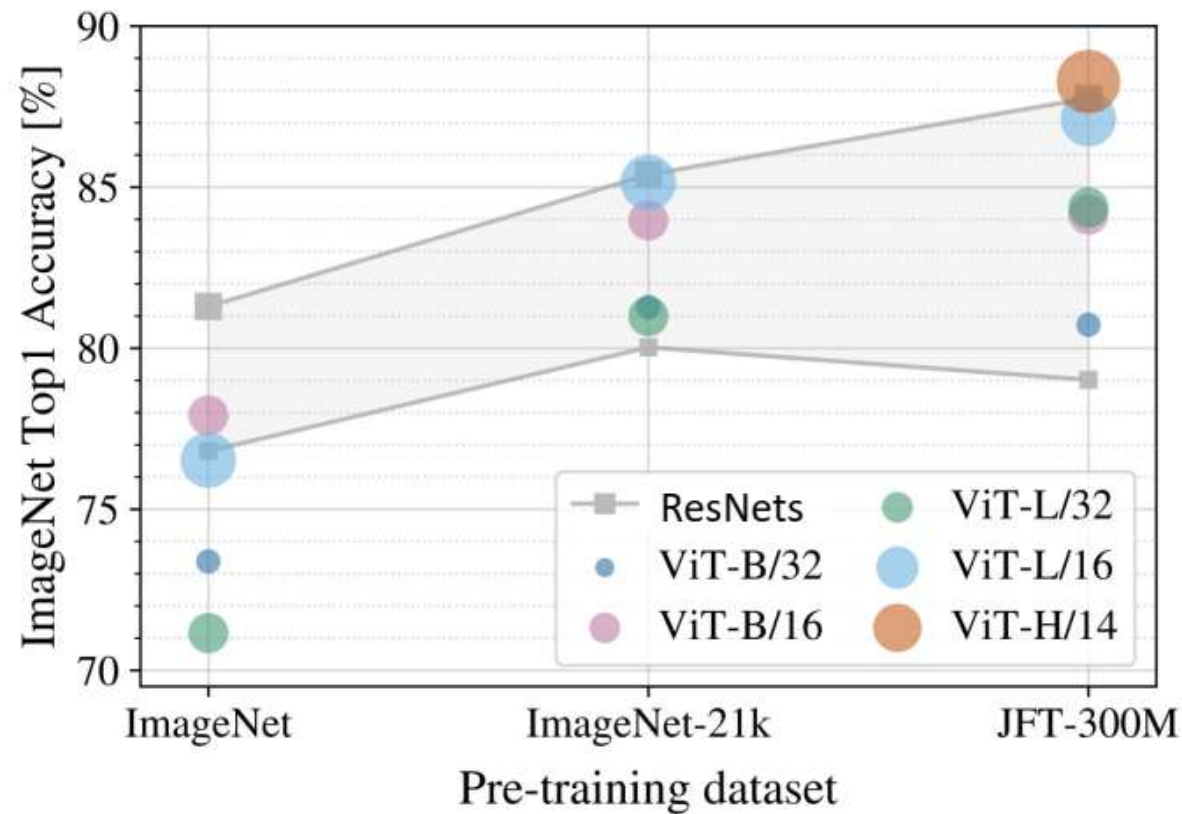
B: base  
L: Large  
H: huge

/32, /16, /14 es el tamaño del patch. Mientras más pequeño es el patch, el modelo es más grande

# ViT vs. ResNets

JFT-300M dataset interno de Google con 300M de imágenes etiquetadas

Pre-entrenamiento con JFT-300M hace que ViT supere a ResNets



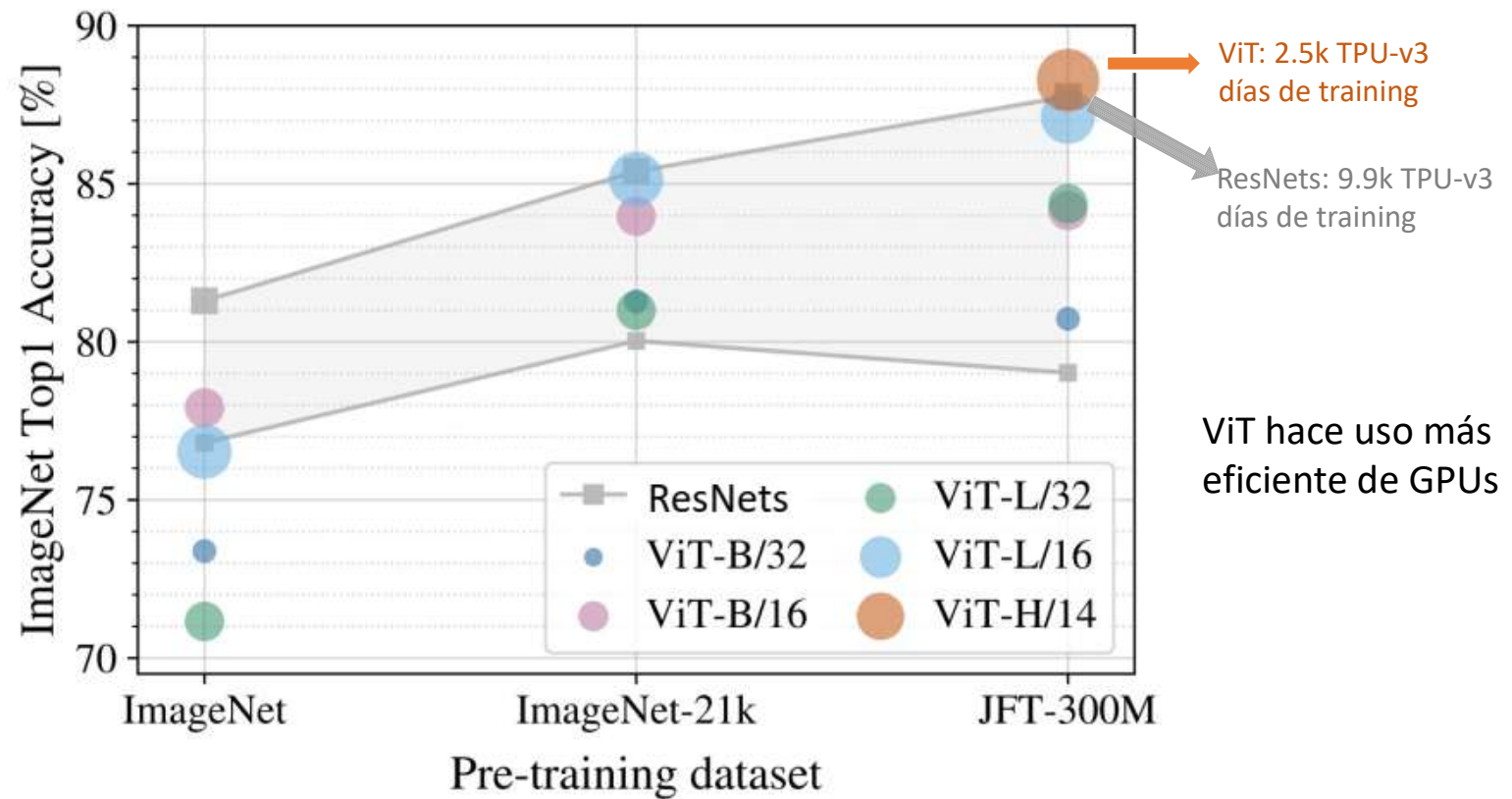
B: base  
L: Large  
H: huge

/32, /16, /14 es el tamaño del patch. Mientras más pequeño es el patch, el modelo es más grande

# ViT vs. ResNets

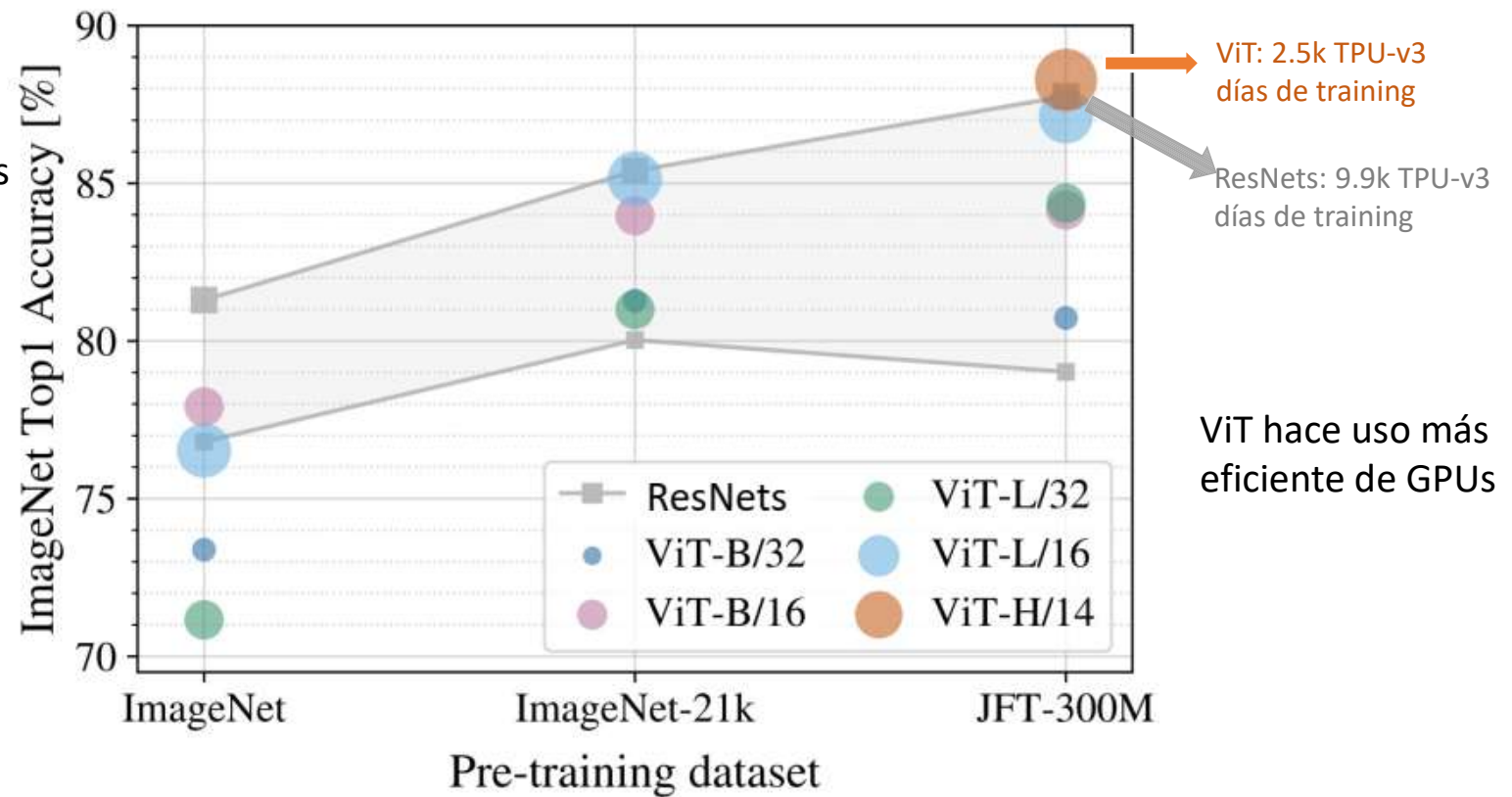
JFT-300M dataset interno de Google con 300M de imágenes etiquetadas

Pre-entrenamiento con JFT-300M hace que ViT supere a ResNets



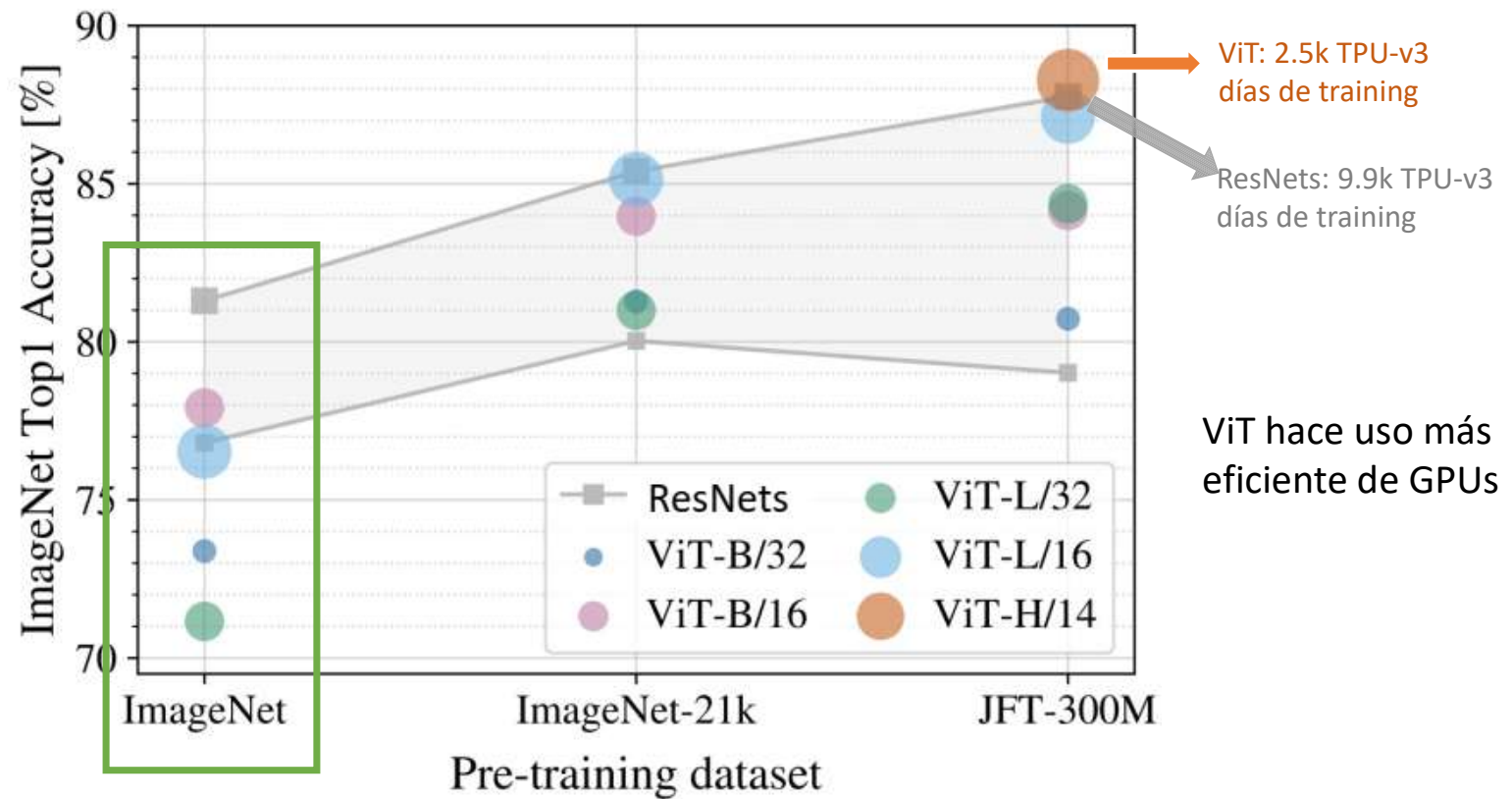
# ViT vs. ResNets

ViT tiene “menos vías inductivo” que ConvNets, por eso necesitamos más datos para pre-entrenar



# ViT vs. ResNets

Cómo podemos mejorar  
el performance con ImageNet?



# Mejorando ViT: Aumentación y Regularización

Regularización para ViT:

- Weight decay
- Dropout (en capas MLP del Transformer)

Data augmentation para ViT:

- MixUp
- RandAugment

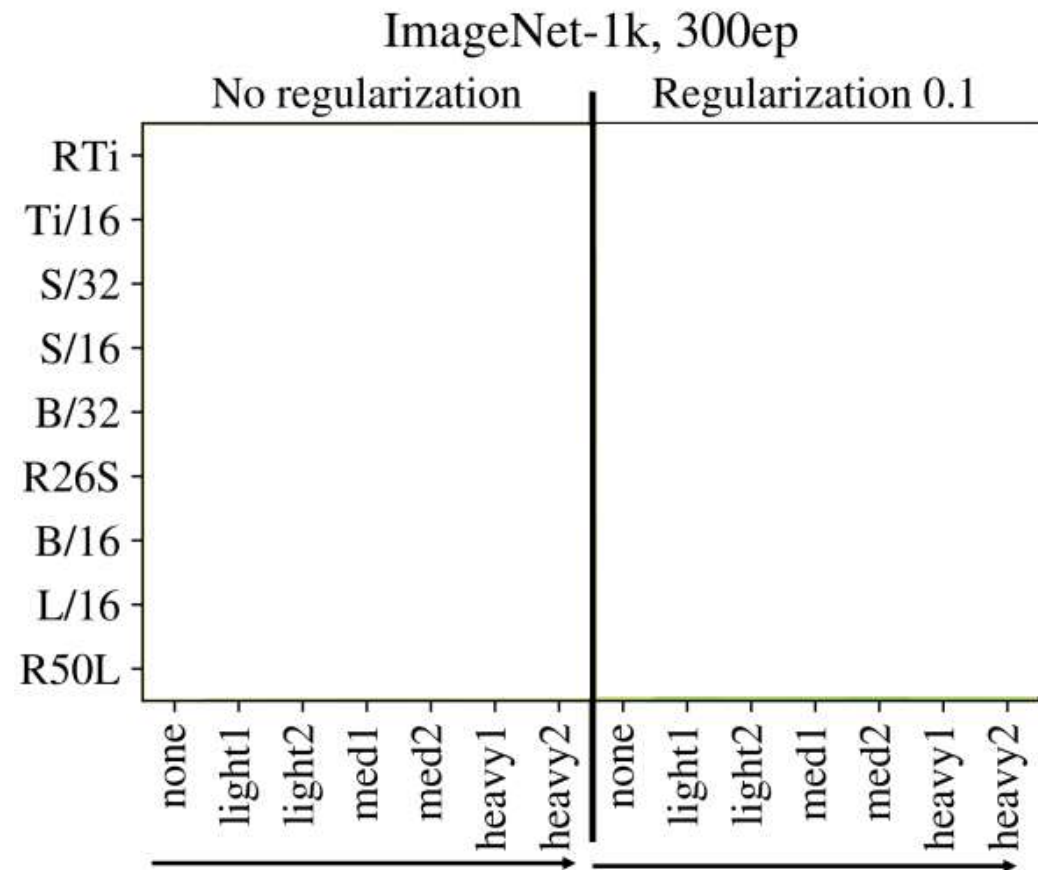
# Mejorando ViT: Aumentación y Regularización

Regularización para ViT:

- Weight decay
- Dropout (en capas MLP del Transformer)

Data augmentation para ViT:

- MixUp
- RandAugment





# Mejorando ViT: Aumentación y Regularización

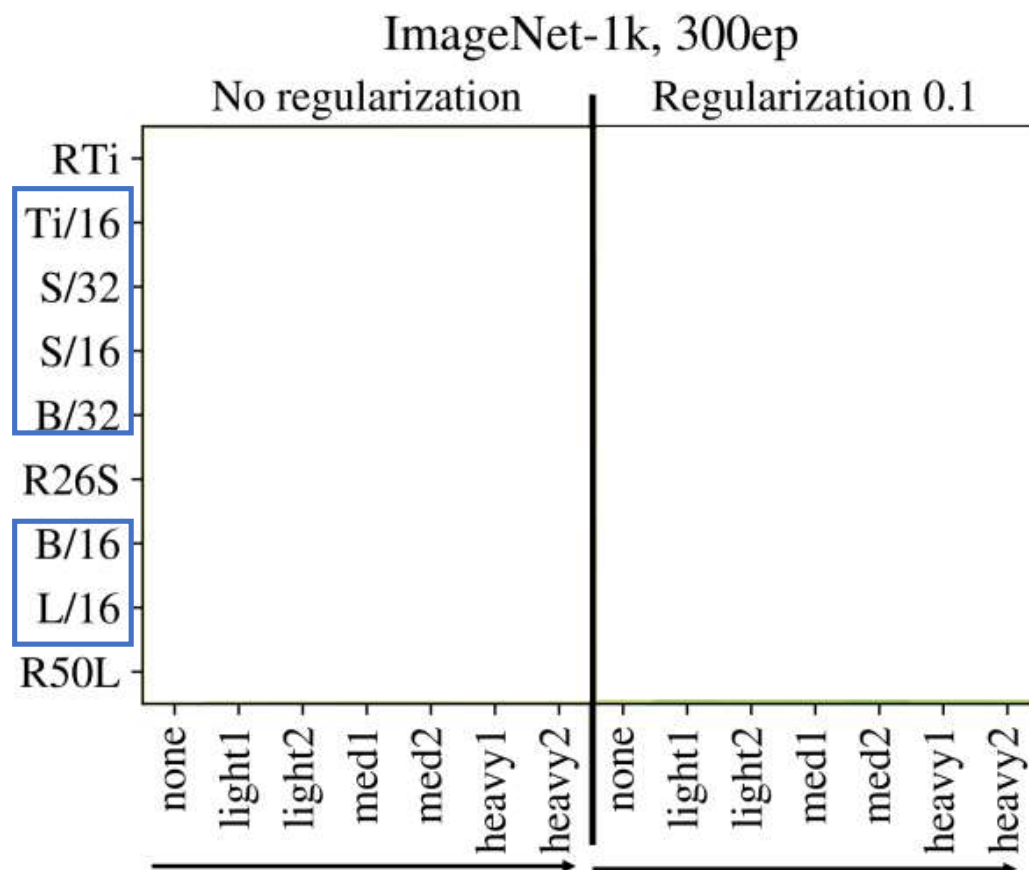
Regularización para ViT:

- Weight decay
- Dropout (en capas MLP del Transformer)

Data augmentation para ViT:

- MixUp
- RandAugment

ViT:  
Ti: tiny  
S: Small  
B: Base  
L: Large



# Mejorando ViT: Aumentación y Regularización

Regularización para ViT:

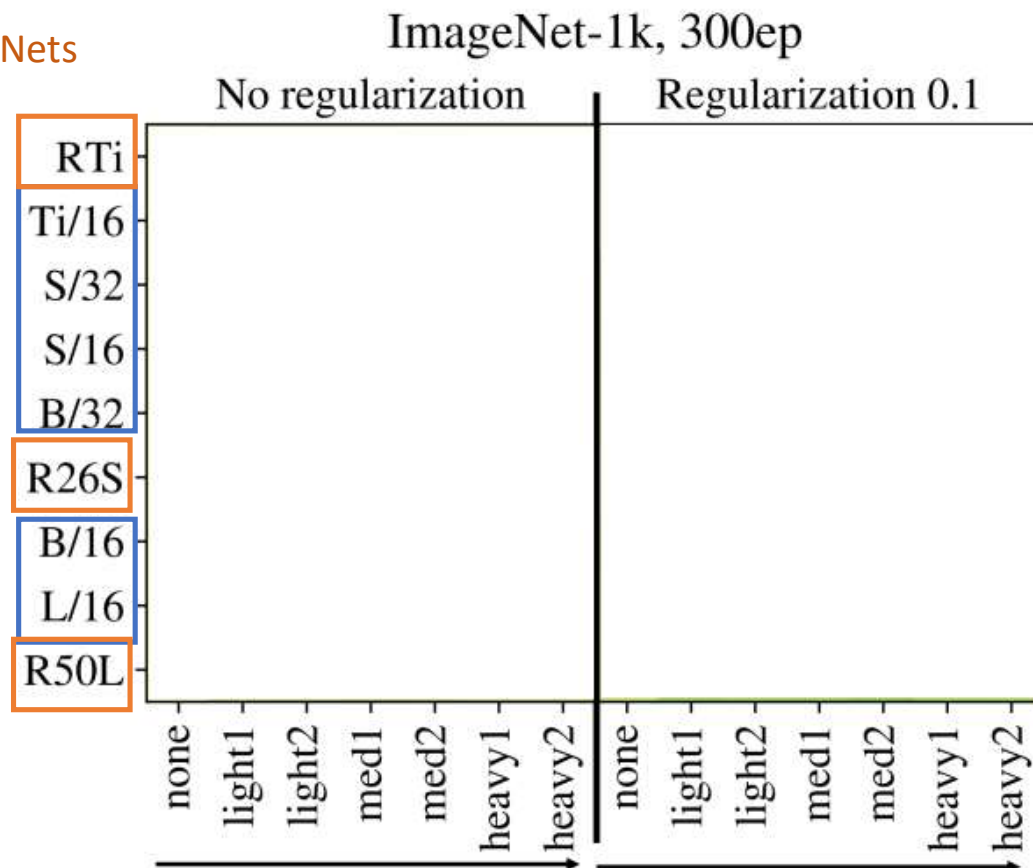
- Weight decay
- Dropout (en capas MLP del Transformer)

Data augmentation para ViT:

- MixUp
- RandAugment

Híbridos: bloques ResNets  
Con bloques ViT

ViT:  
Ti: tiny  
S: Small  
B: Base  
L: Large



# Mejorando ViT: Aumentación y Regularización

Regularización para ViT:

- Weight decay
- Dropout (en capas MLP del Transformer)

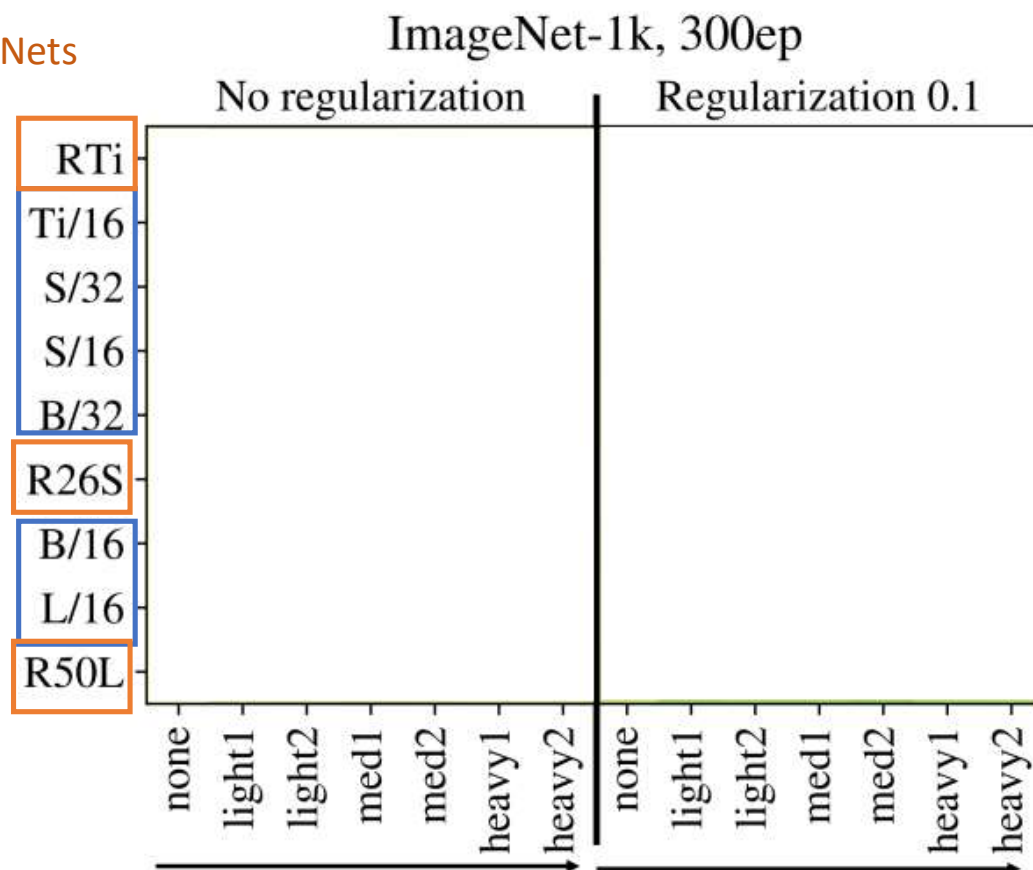
Data augmentation para ViT:

- MixUp
- RandAugment

Híbridos: bloques ResNets  
Con bloques ViT

ViT:  
Ti: tiny  
S: Small  
B: Base  
L: Large

Original:  
77.9  
76.53



# Mejorando ViT: Aumentación y Regularización

Regularización para ViT:

- Weight decay
- Dropout (en capas MLP del Transformer)

Data augmentation para ViT:

- MixUp
- RandAugment

Regularización siempre ayuda!

Híbridos: bloques ResNets  
Con bloques ViT

ViT:  
Ti: tiny  
S: Small  
B: Base  
L: Large

Original  
77.9  
76.53

		ImageNet-1k, 300ep													
		No regularization							Regularization 0.1						
		none	light1	light2	med1	med2	heavy1	heavy2	none	light1	light2	med1	med2	heavy1	heavy2
Híbridos: bloques ResNets Con bloques ViT	RTi	69							71						
	Ti/16	72							71						
	S/32	64							70						
	S/16	71							76						
	B/32	63							69						
	R26S	72							75						
Original	B/16	70							76						
	L/16	69							74						
	R50L	70							75						

# Mejorando ViT: Aumentación y Regularización

Regularización para ViT:

- Weight decay
- Dropout (en capas MLP del Transformer)

Data augmentation para ViT:

- MixUp
- RandAugment

Regularización + aumentación ayuda más!

Híbridos: bloques ResNets  
Con bloques ViT

ViT:  
Ti: tiny  
S: Small  
B: Base  
L: Large

Original  
77.9  
76.53

		No regularization							Regularization 0.1						
		none	light1	light2	med1	med2	heavy1	heavy2	none	light1	light2	med1	med2	heavy1	heavy2
Híbridos: bloques ResNets Con bloques ViT	RTi	69							71						
	Ti/16	72							71						
	S/32	64							70						
	S/16	71							76						
	B/32	63							69						
	R26S	72							75						
Original	B/16	70	76	79	79	81	80	80	76	79	81	82	83	82	82
	L/16	69	76	77	78	78	76	76	74	78	78	78	79	77	77
	R50L	70							75						

# Mejorando ViT: Aumentación y Regularización

Regularización para ViT:

- Weight decay
- Dropout (en capas MLP del Transformer)

Data augmentation para ViT:

- MixUp
- RandAugment

Todos los resultados

Híbridos: bloques ResNets  
Con bloques ViT

ViT:  
Ti: tiny  
S: Small  
B: Base  
L: Large

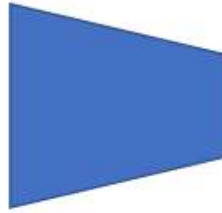
Original:  
77.9  
76.53

ImageNet-1k, 300ep

	No regularization								Regularization 0.1							
	none	light1	light2	med1	med2	heavy1	heavy2		none	light1	light2	med1	med2	heavy1	heavy2	
RTi	69	73	73	72	70	69	68		71	70	67	65	63	62	61	
Ti/16	72	76	75	75	74	72	71		71	72	68	65	63	63	62	
S/32	64	71	76	76	76	74	74		70	72	72	71	71	69	68	
S/16	71	77	79	81	82	80	80		76	79	80	79	79	77	77	
B/32	63	70	73	75	76	75	76		69	74	77	77	78	77	77	
R26S	72	76	78	79	80	80	80		75	78	81	82	82	81	81	
B/16	70	76	79	79	81	80	80		76	79	81	82	83	82	82	
L/16	69	76	77	78	78	76	76		74	78	78	78	79	77	77	
R50L	70	75	76	77	77	76	76		75	78	78	78	79	77	77	

# Mejorando ViT: Destilación

Paso 1: Entrenar un **modelo teacher** en imágenes con etiquetas

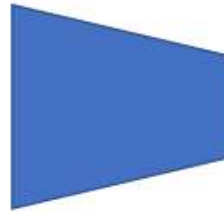


$P(\text{cat}) = 0.9$   
 $P(\text{dog}) = 0.1$

→ Cross Entropy Loss ← GT label: Cat

# Mejorando ViT: Destilación

Paso 1: Entrenar un **modelo teacher** en imágenes con etiquetas

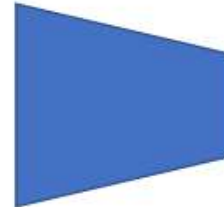


$P(\text{cat}) = 0.9$   
 $P(\text{dog}) = 0.1$

Cross  
Entropy  
Loss

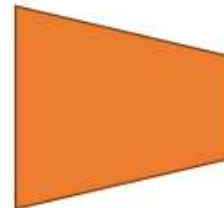
GT label:  
Cat

Paso 2: Entrenar un **modelo estudiante** para matchear predicciones del **teacher**

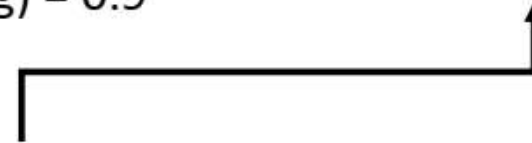


$P(\text{cat}) = 0.1$   
 $P(\text{dog}) = 0.9$

KL Divergence Loss



$P(\text{cat}) = 0.2$   
 $P(\text{dog}) = 0.8$





# Mejorando ViT: Destilación

Paso 1: Entrenar un **modelo teacher** en imágenes con etiquetas

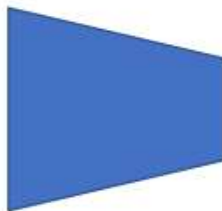


$P(\text{cat}) = 0.9$   
 $P(\text{dog}) = 0.1$

Cross  
Entropy  
Loss

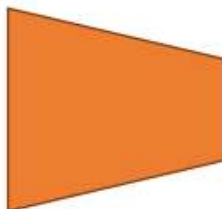
GT label:  
Cat

Paso 2: Entrenar un **modelo estudiante** para matchear predicciones del **teacher** (algunas veces también las etiquetas)



$P(\text{cat}) = 0.1$   
 $P(\text{dog}) = 0.9$

KL Divergence Loss



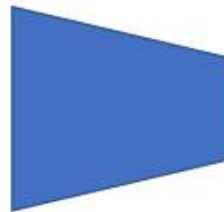
$P(\text{cat}) = 0.2$   
 $P(\text{dog}) = 0.8$

Cross  
Entropy  
Loss

GT label:  
Dog

# Mejorando ViT: Destilación

Paso 1: Entrenar un **modelo Teacher CNN** en imágenes con etiquetas

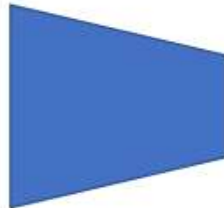


$P(\text{cat}) = 0.9$   
 $P(\text{dog}) = 0.1$

Cross  
Entropy  
Loss

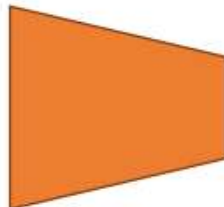
GT label:  
Cat

Paso 2: Entrenar un **modelo Estudiante ViT** para matchear predicciones del **teacher CNN** (algunas veces también las etiquetas)



$P(\text{cat}) = 0.1$   
 $P(\text{dog}) = 0.9$

KL Divergence Loss

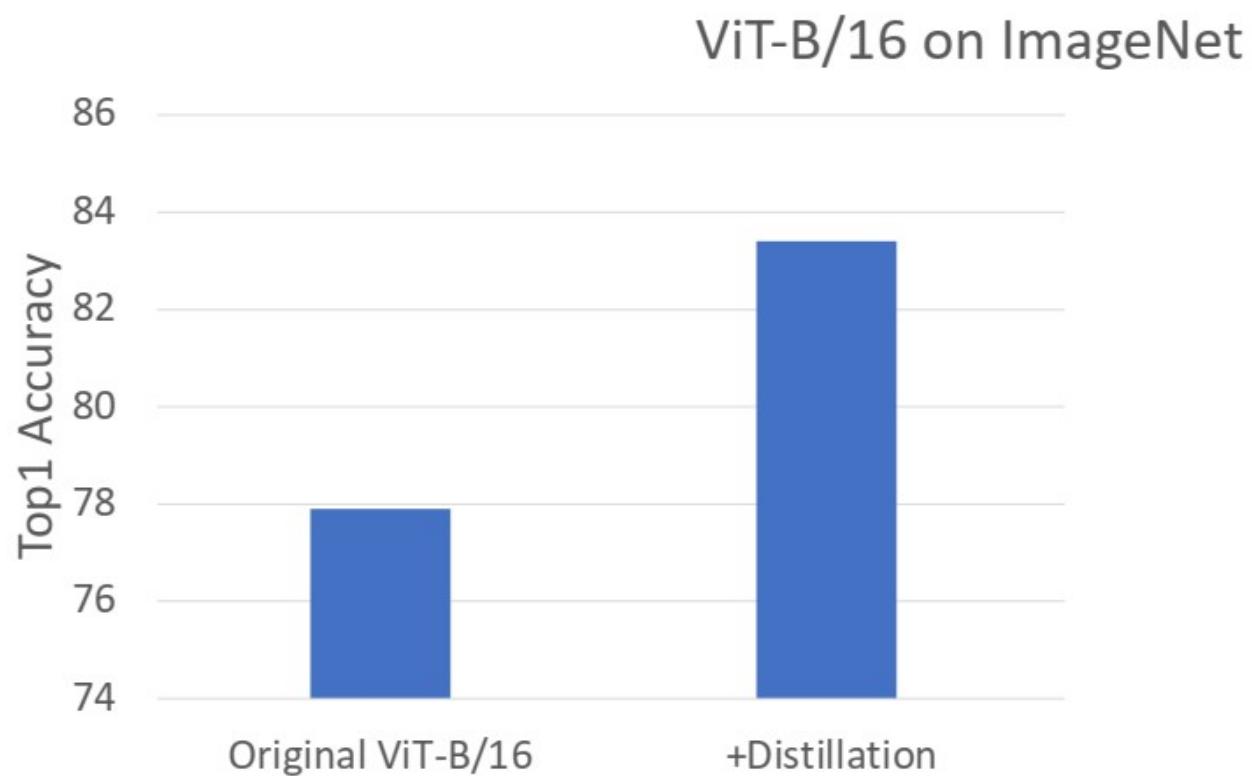


$P(\text{cat}) = 0.2$   
 $P(\text{dog}) = 0.8$

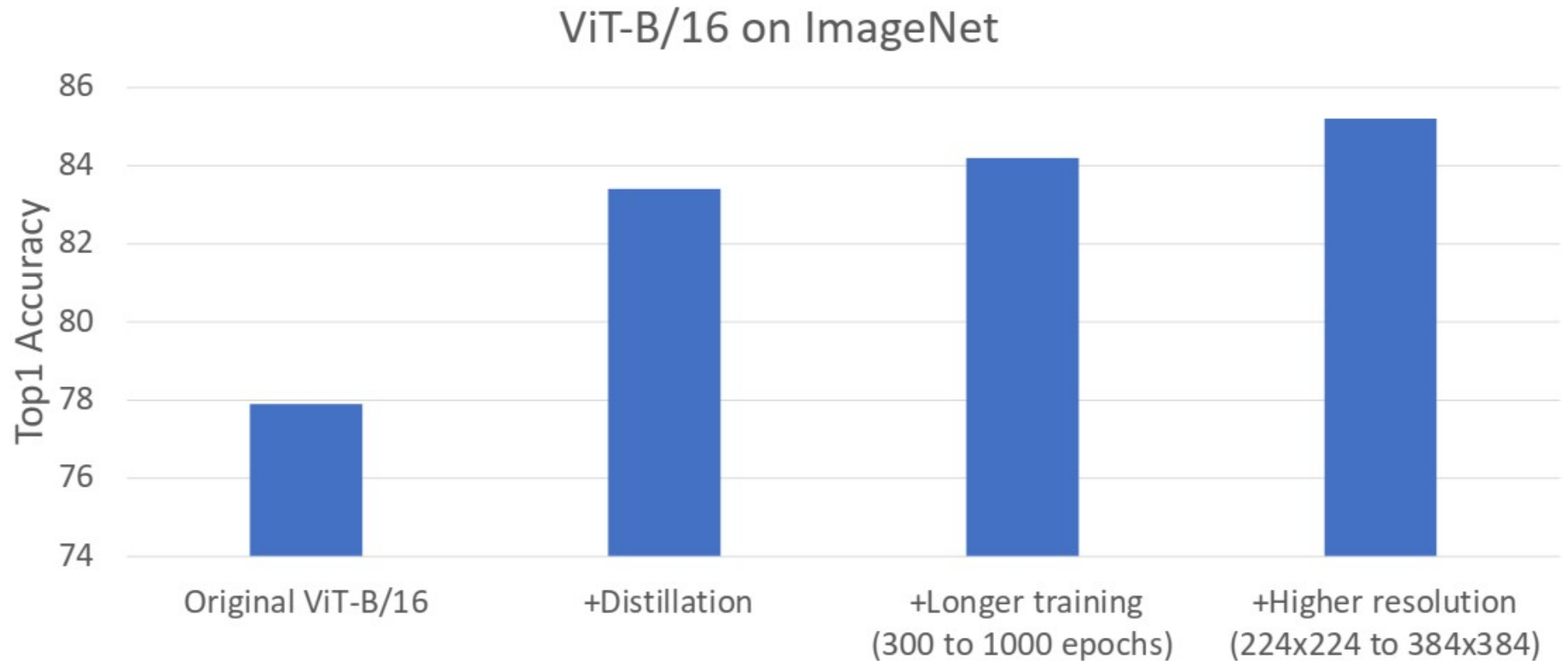
Cross  
Entropy  
Loss

GT label:  
Dog

# Mejorando ViT: Destilación

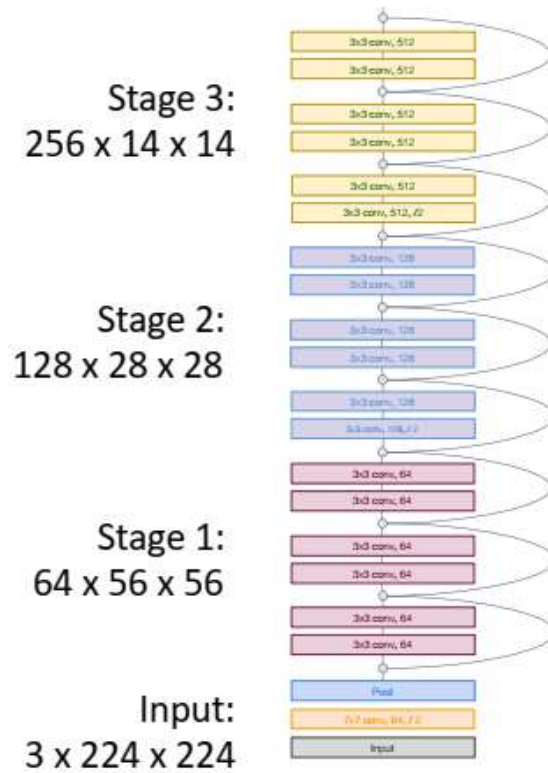


# Mejorando ViT: Destilación



# ViT vs. CNN

CNN se va reduciendo la resolución y ampliando la profundidad

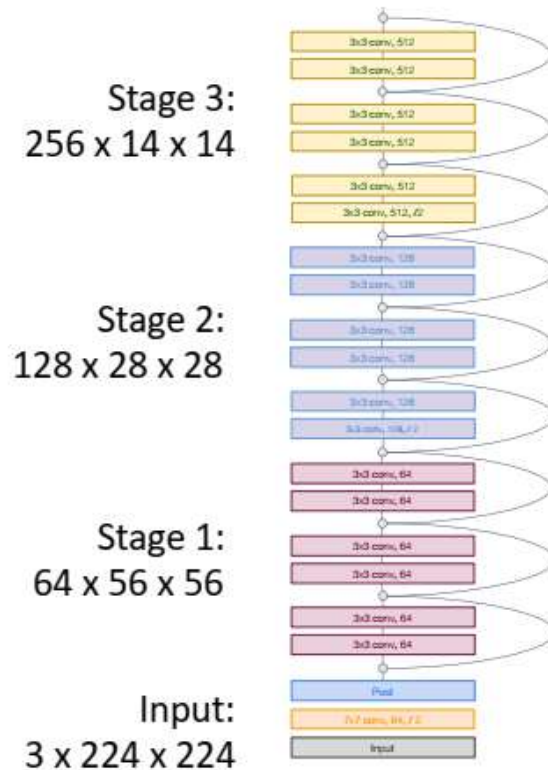


Útil para ser robustos a  
La escala

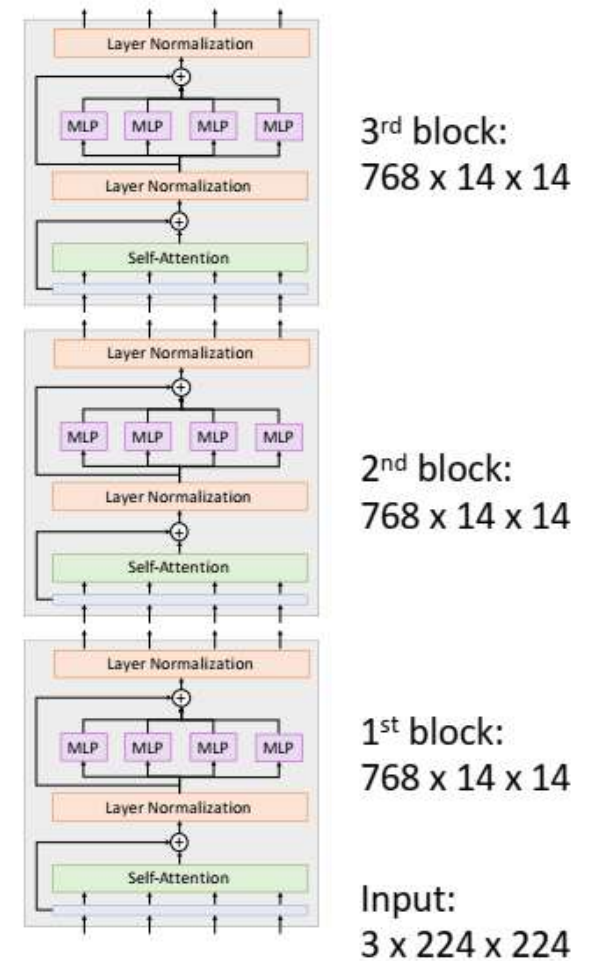
# ViT vs. CNN

CNN se va reduciendo la resolución y ampliando la profundidad

Útil para ser robustos a la escala



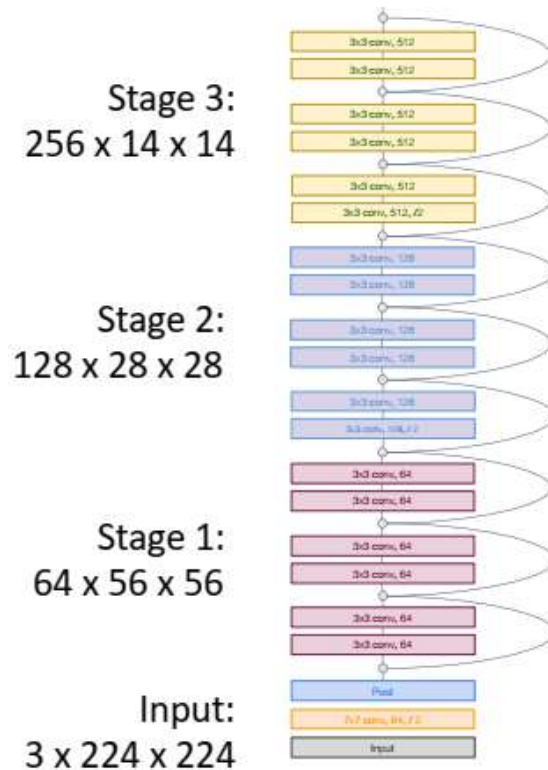
En ViT, todos los bloques tienen la misma resolución



# ViT vs. CNN

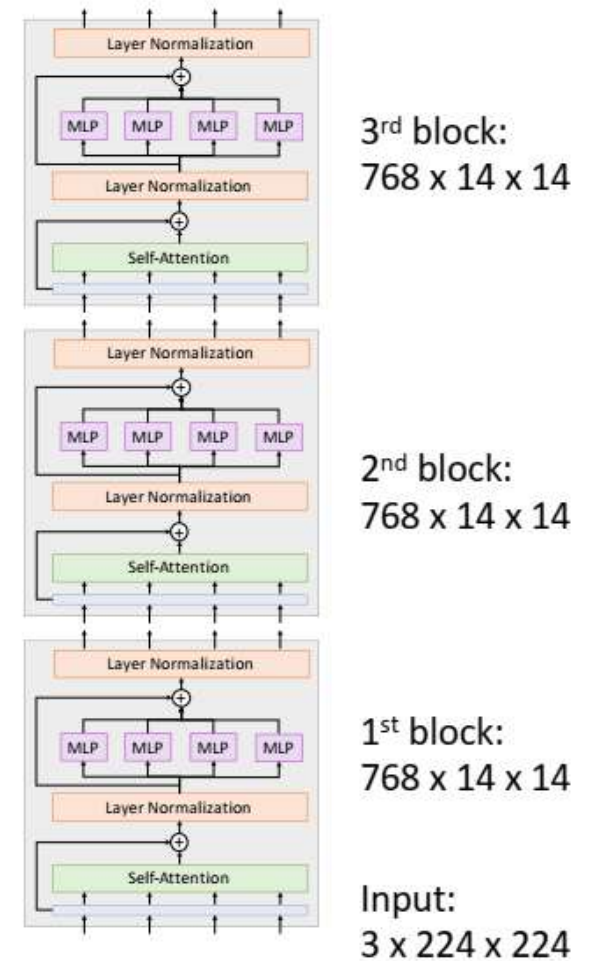
CNN se va reduciendo la resolución y ampliando la profundidad

## Útil para ser robustos a La escala

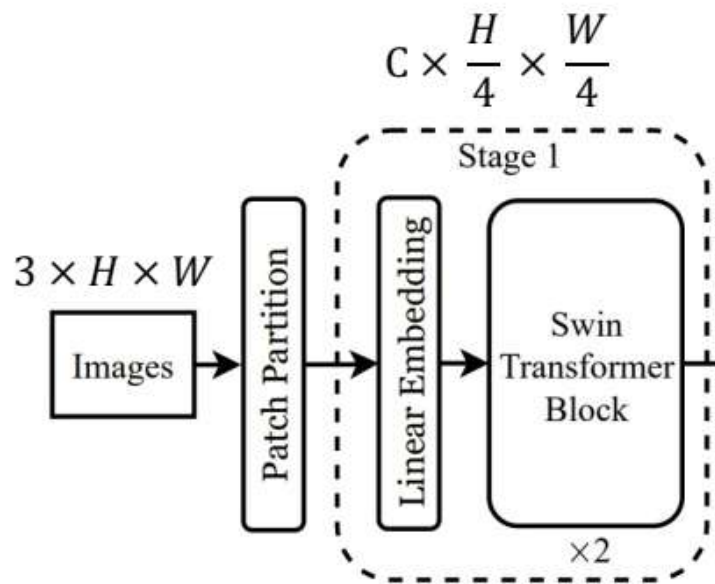


En ViT, todos los bloques tienen la misma resolución

## Se puede construir un ViT jerárquico?



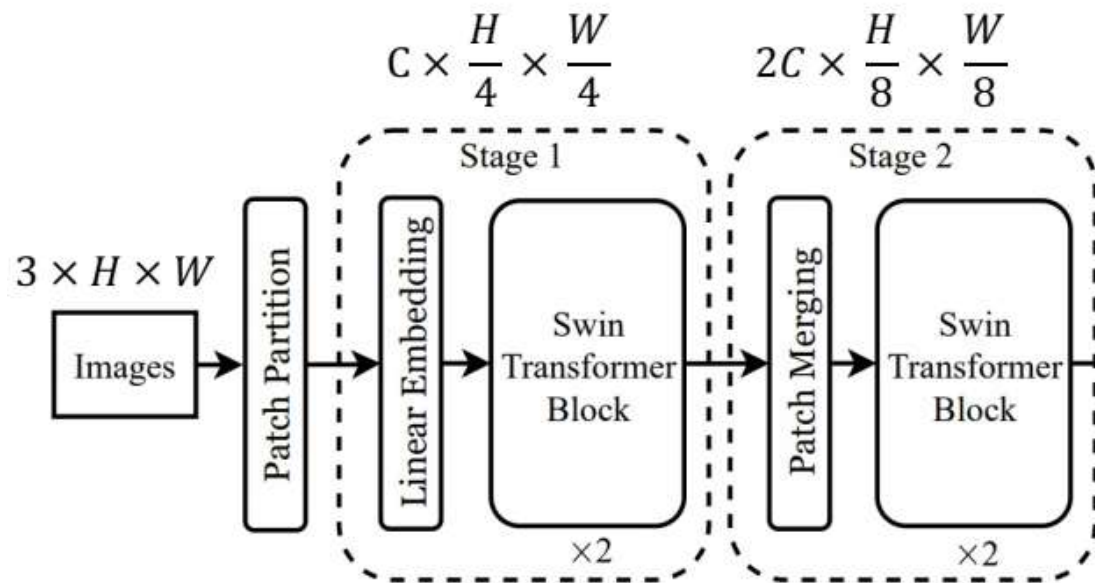
# ViT Jerárquico: Swin Transformer



Dividir imagen en 4x4  
Parches y proyectar a C  
dimensiones

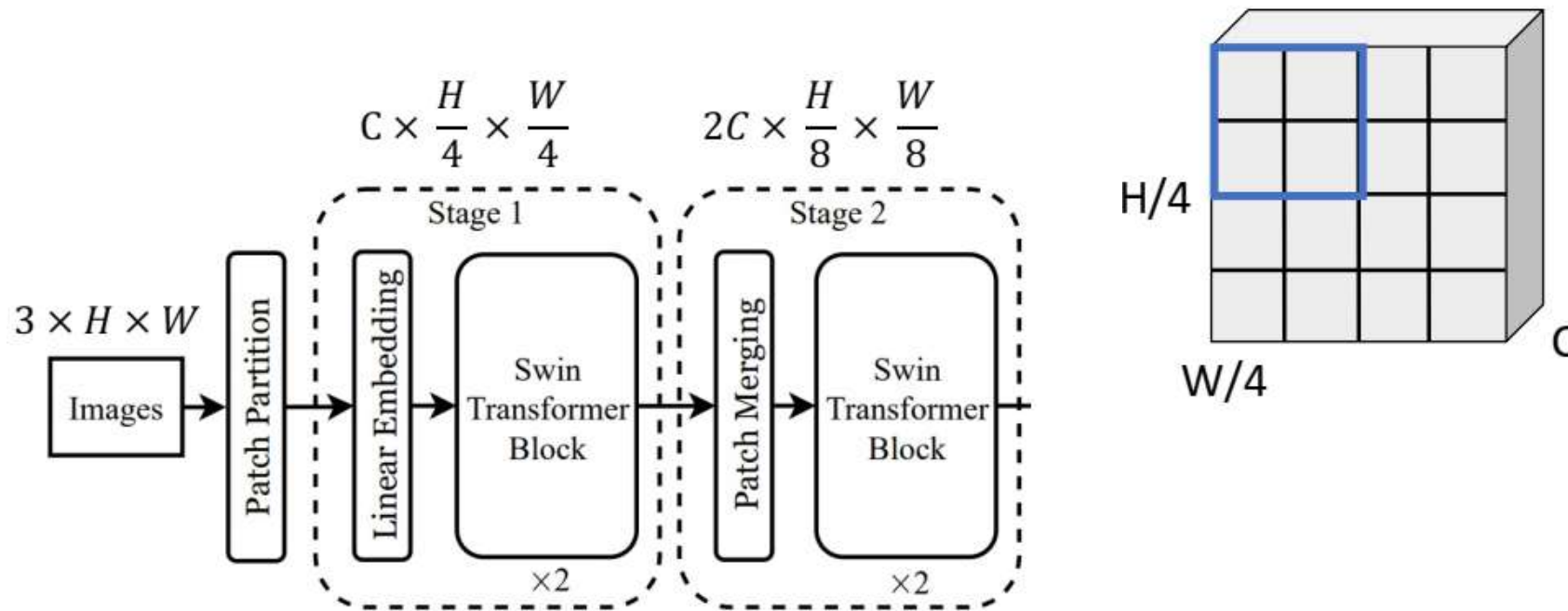


# ViT Jerárquico: Swin Transformer



Dividir imagen en  $4 \times 4$  Parches y proyectar a  $C$  dimensiones  
Unir  $2 \times 2$  vecinos: parches de  $8 \times 8$

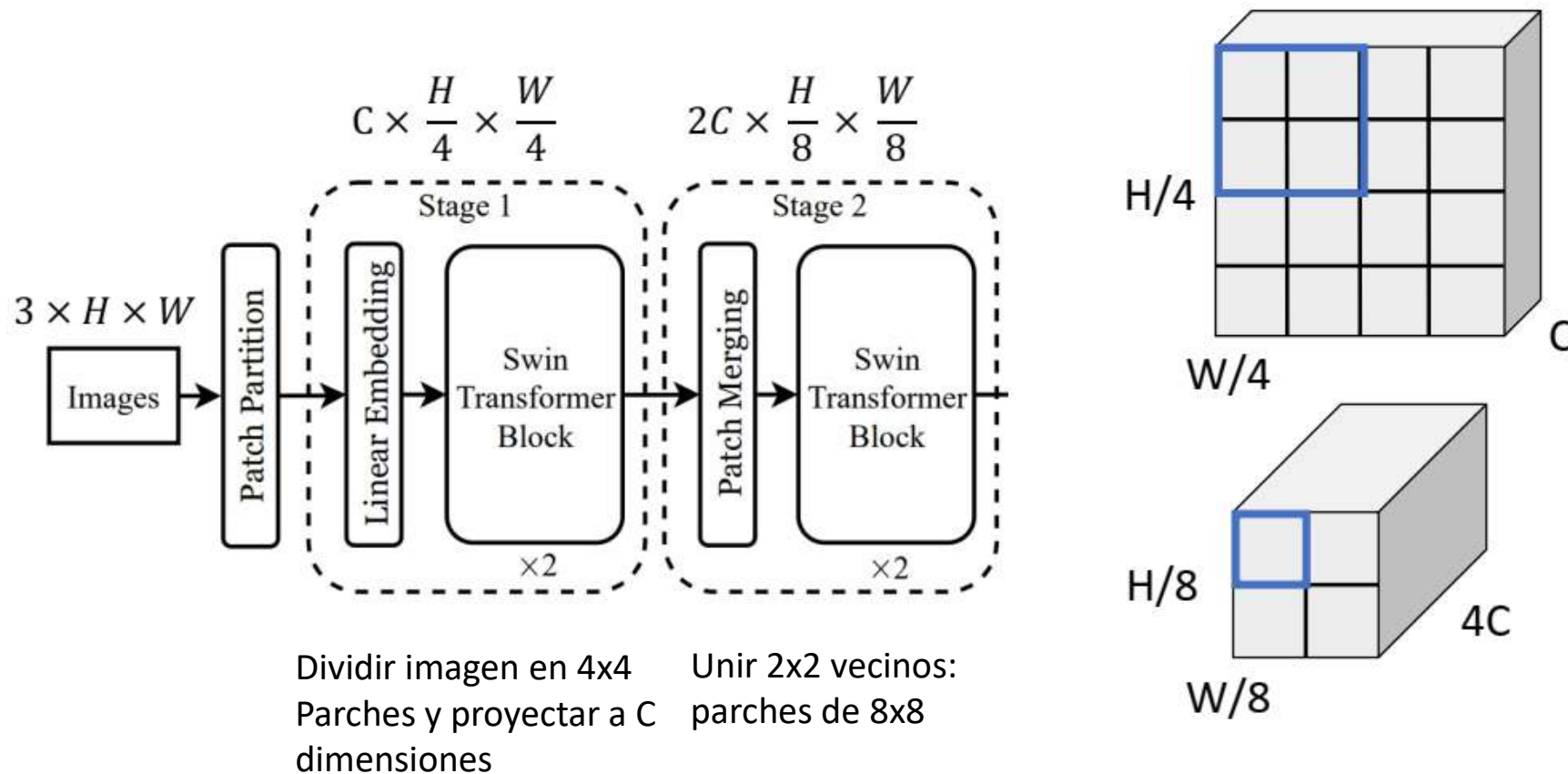
# ViT Jerárquico: Swin Transformer



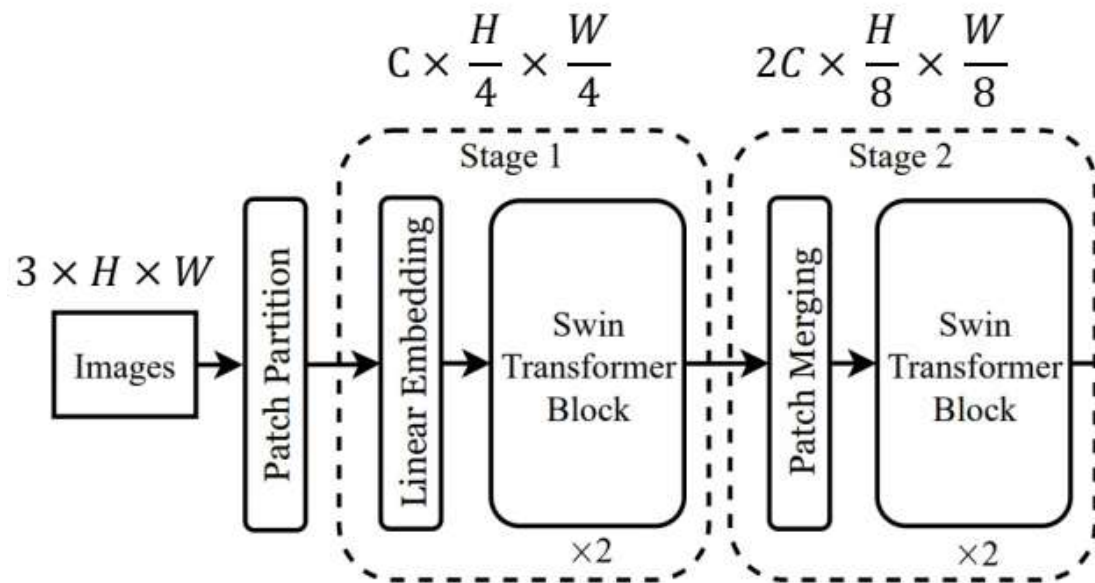
Dividir imagen en 4x4  
Parches y proyectar a C  
dimensiones

Unir 2x2 vecinos:  
parches de 8x8

# ViT Jerárquico: Swin Transformer

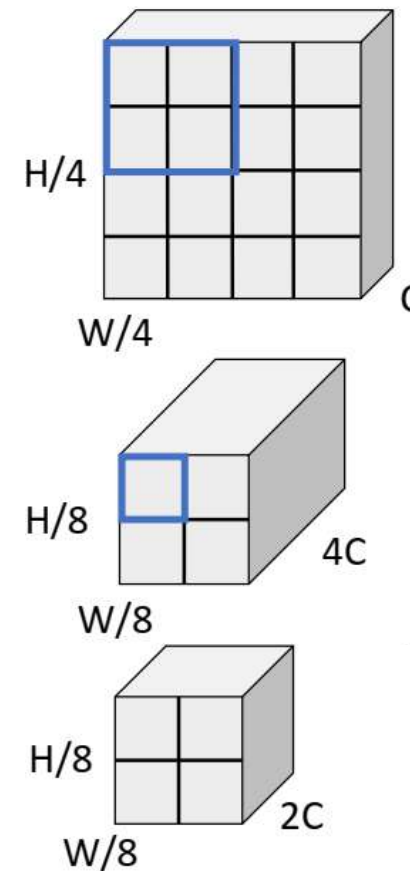


# ViT Jerárquico: Swin Transformer



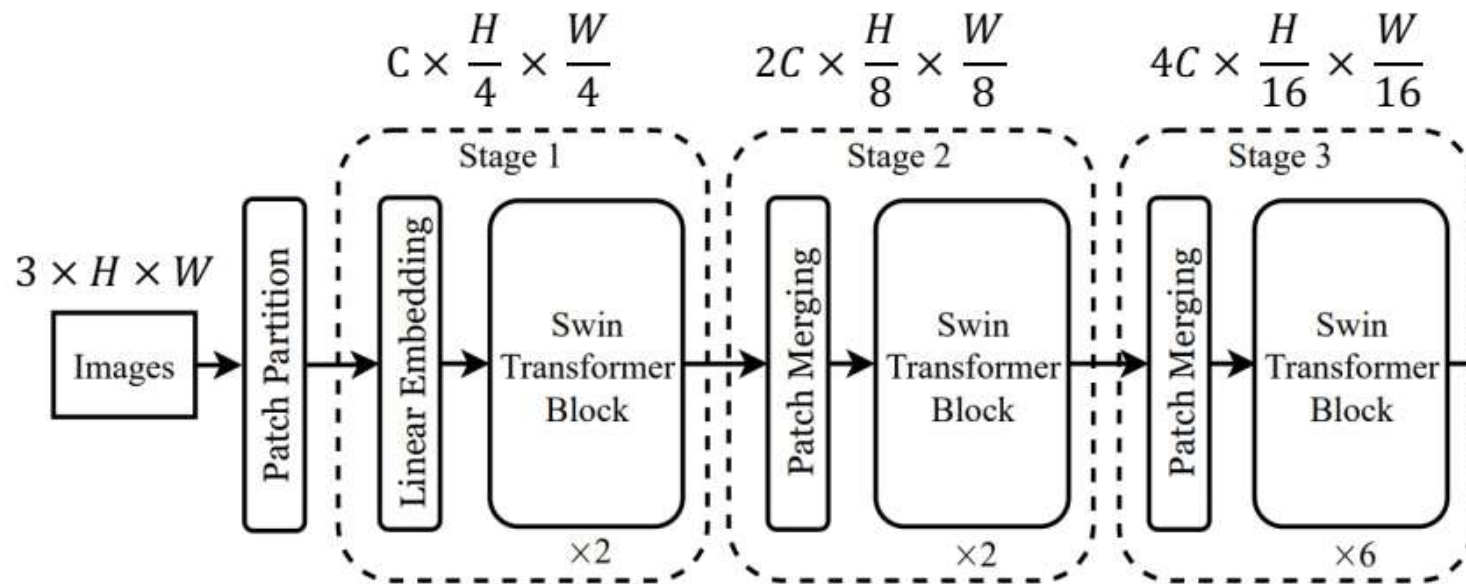
Dividir imagen en 4x4  
Parches y proyectar a C  
dimensiones

Unir 2x2 vecinos:  
parches de 8x8



Proyección lineal  
De  $4C$  a  $2C$  canales

# ViT Jerárquico: Swin Transformer

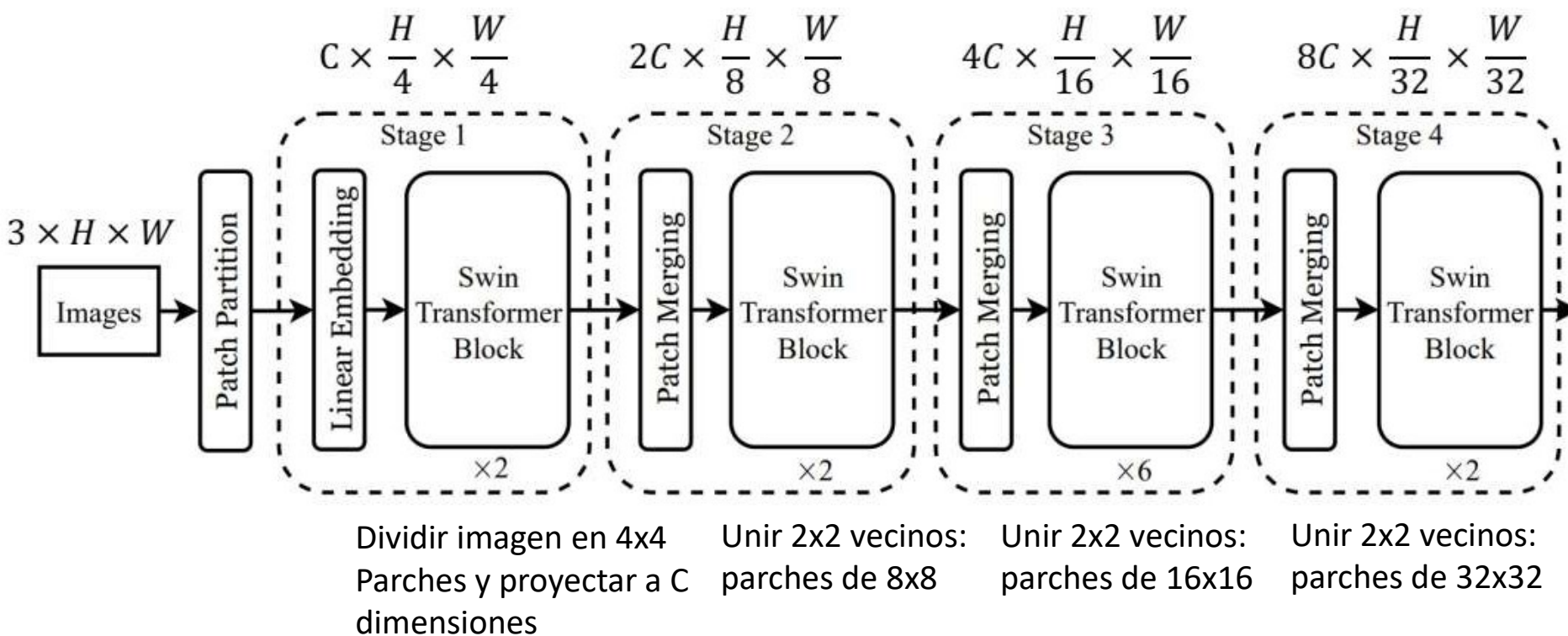


Dividir imagen en 4x4  
Parches y proyectar a C  
dimensiones

Unir 2x2 vecinos:  
parches de 8x8

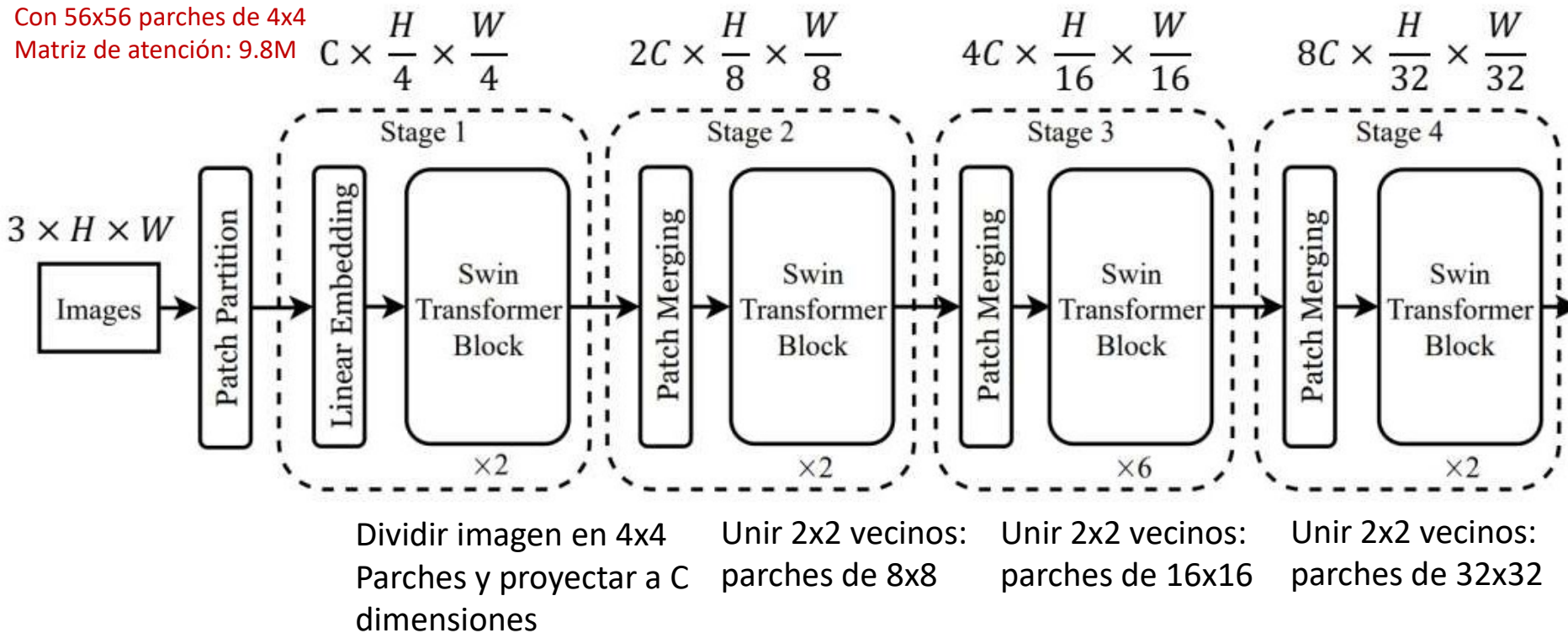
Unir 2x2 vecinos:  
parches de 16x16

# ViT Jerárquico: Swin Transformer



# ViT Jerárquico: Swin Transformer

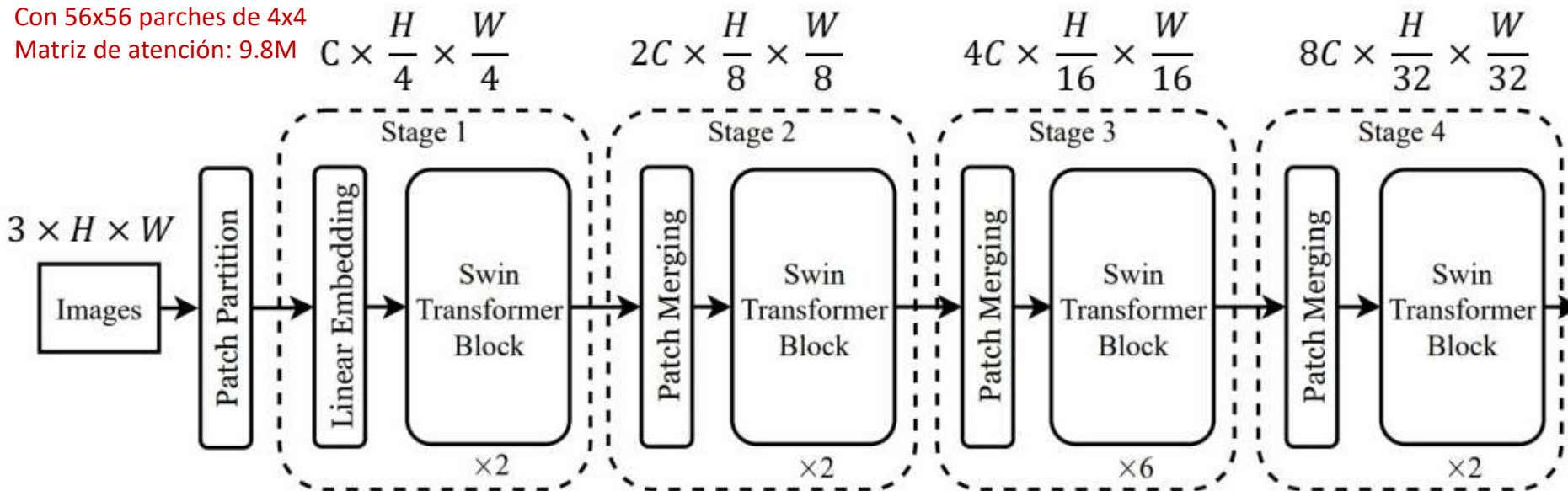
Problema: 224x224 imagen  
Con 56x56 parches de 4x4  
Matriz de atención: 9.8M





# ViT Jerárquico: Swin Transformer

Problema: 224x224 imagen  
Con 56x56 parches de 4x4  
Matriz de atención: 9.8M



Solución: no usar full attention  
Atención sobre parches

Dividir imagen en 4x4  
Parches y proyectar a C  
dimensiones

Unir 2x2 vecinos:  
parches de 8x8

Unir 2x2 vecinos:  
parches de 16x16

Unir 2x2 vecinos:  
parches de 32x32

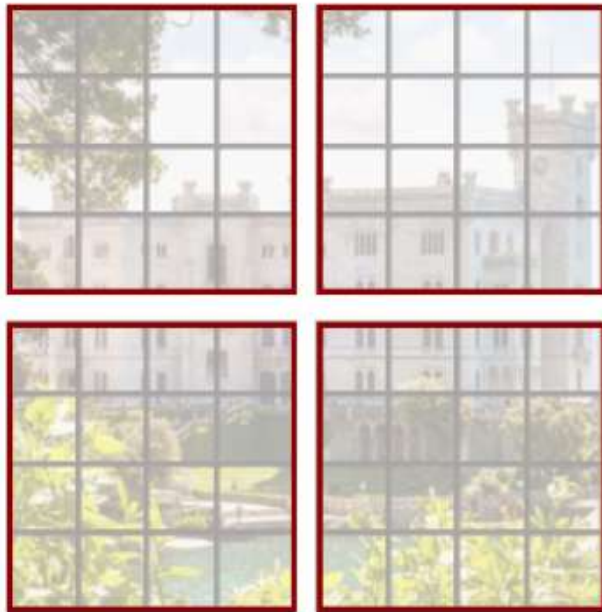


# Swin Transformer: Window Attention

Con una grilla de  $H \times W$  de tokens, cada matriz de atención necesita  $H^2W^2$  - **cuadrático**

# Swin Transformer: Window Attention

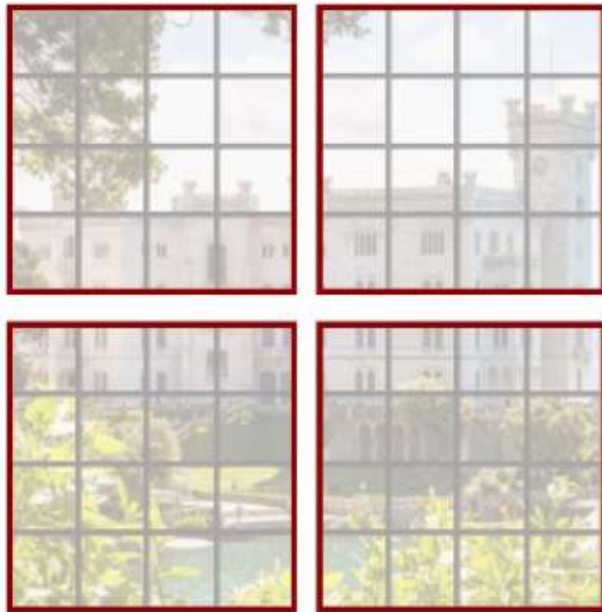
Con una grilla de  $H \times W$  de tokens, cada matriz de atención necesita  $H^2W^2$  - **cuadrático**



En vez que un token atienda a todos los otros tokens,  
Dividir en ventanas de  $M \times M$  tokens (aquí  $M = 4$ );  
Solo computar atención dentro de las ventanas

# Swin Transformer: Window Attention

Con una grilla de  $H \times W$  de tokens, cada matriz de atención necesita  $H^2W^2$  - **cuadrático**



En vez que un token atiende a todos los otros tokens,  
Dividir en ventanas de  $M \times M$  tokens (aquí  $M = 4$ );  
Solo computar atención dentro de las ventanas

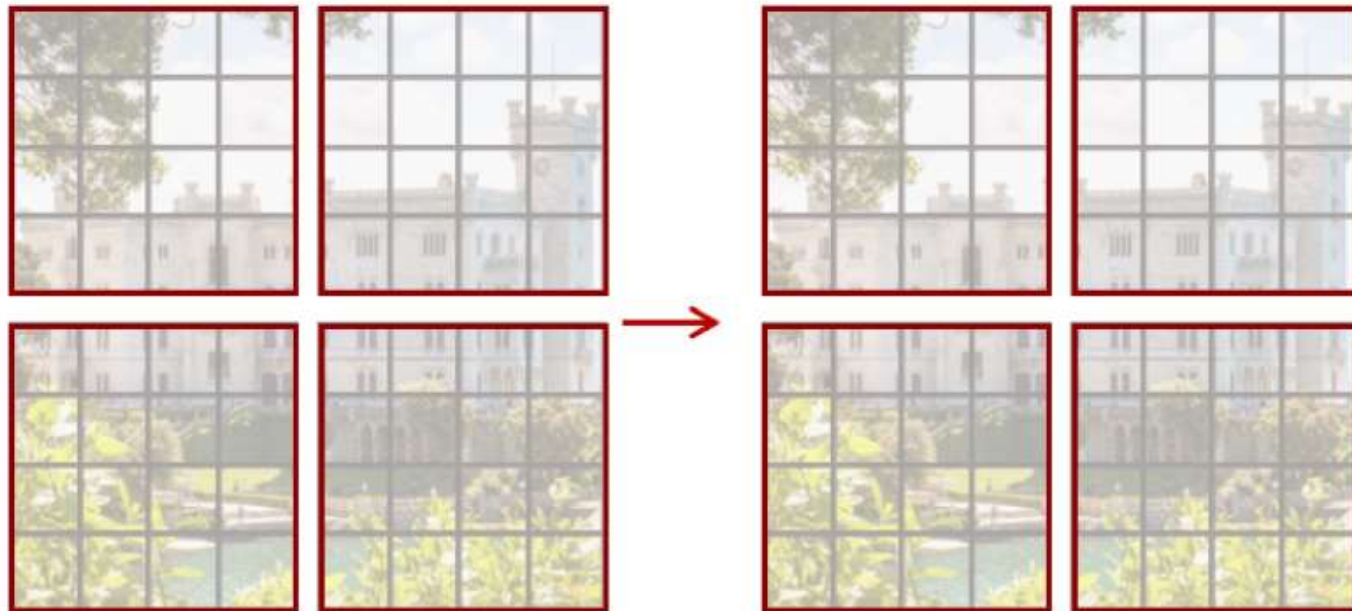
Tamaño total de matriz de atención:

$$M^4 \left( \frac{H}{M} \right) \left( \frac{W}{M} \right) = M^2 HW$$

Lineal cuando  $M$  es fijo. Swin usa  $M=7$

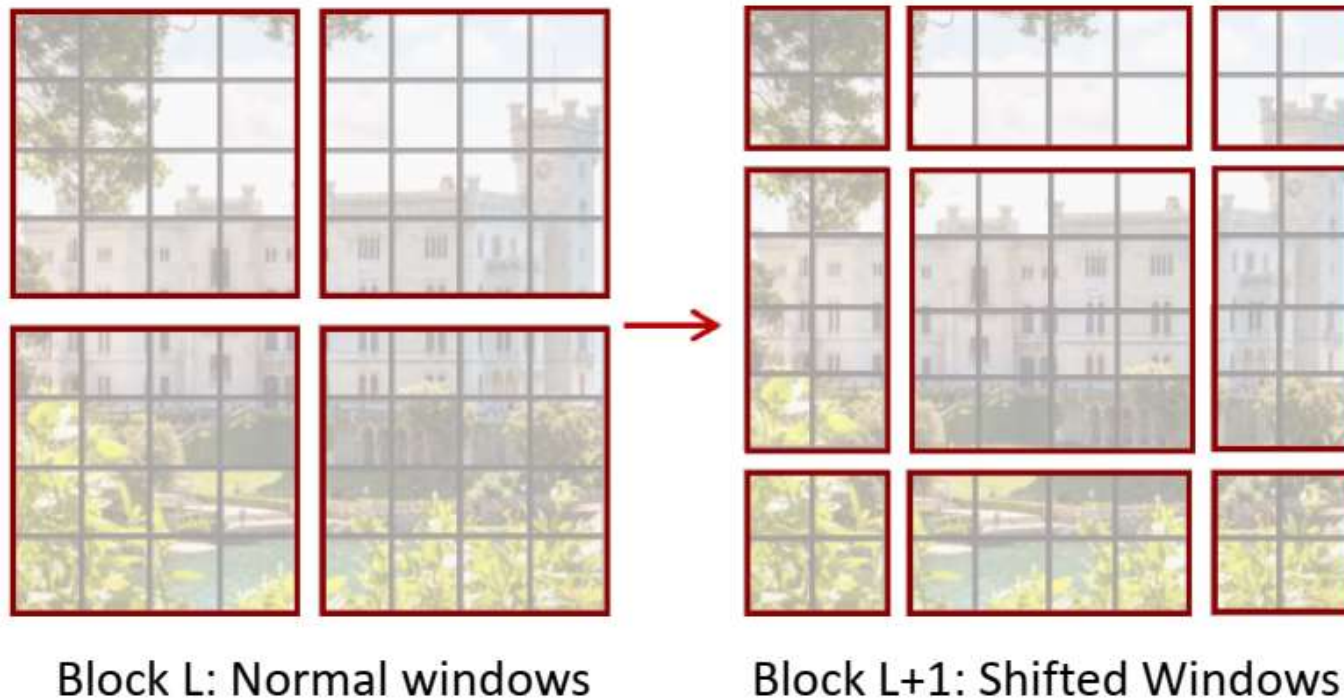
# Swin Transformer: Window Attention

Problema: tokens solo interactúan con otros tokens en la misma ventana, pero no hay comunicación entre ventanas



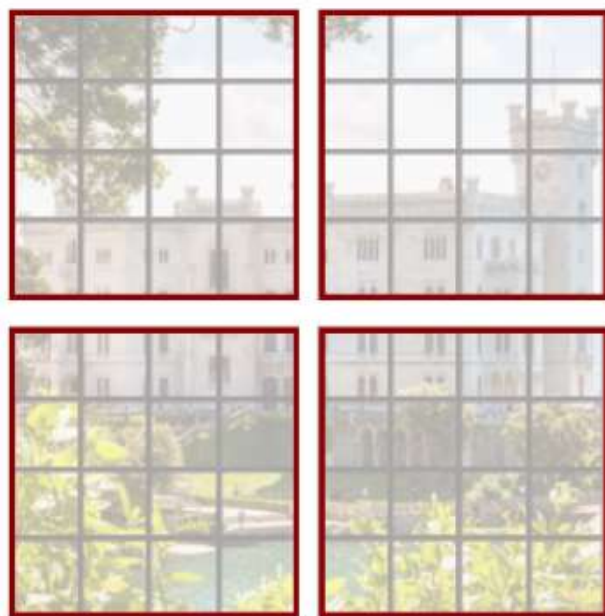
# Swin Transformer: Window Attention

Solución: alternar entre ventanas normales y ventanas desplazadas en bloques sucesivos

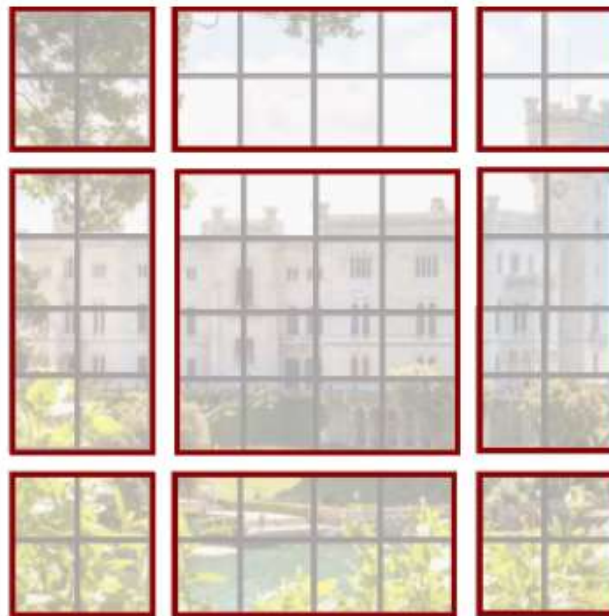


# Swin Transformer: Window Attention

Solución: alternar entre ventanas normales y ventanas desplazadas en bloques sucesivos



Block L: Normal windows

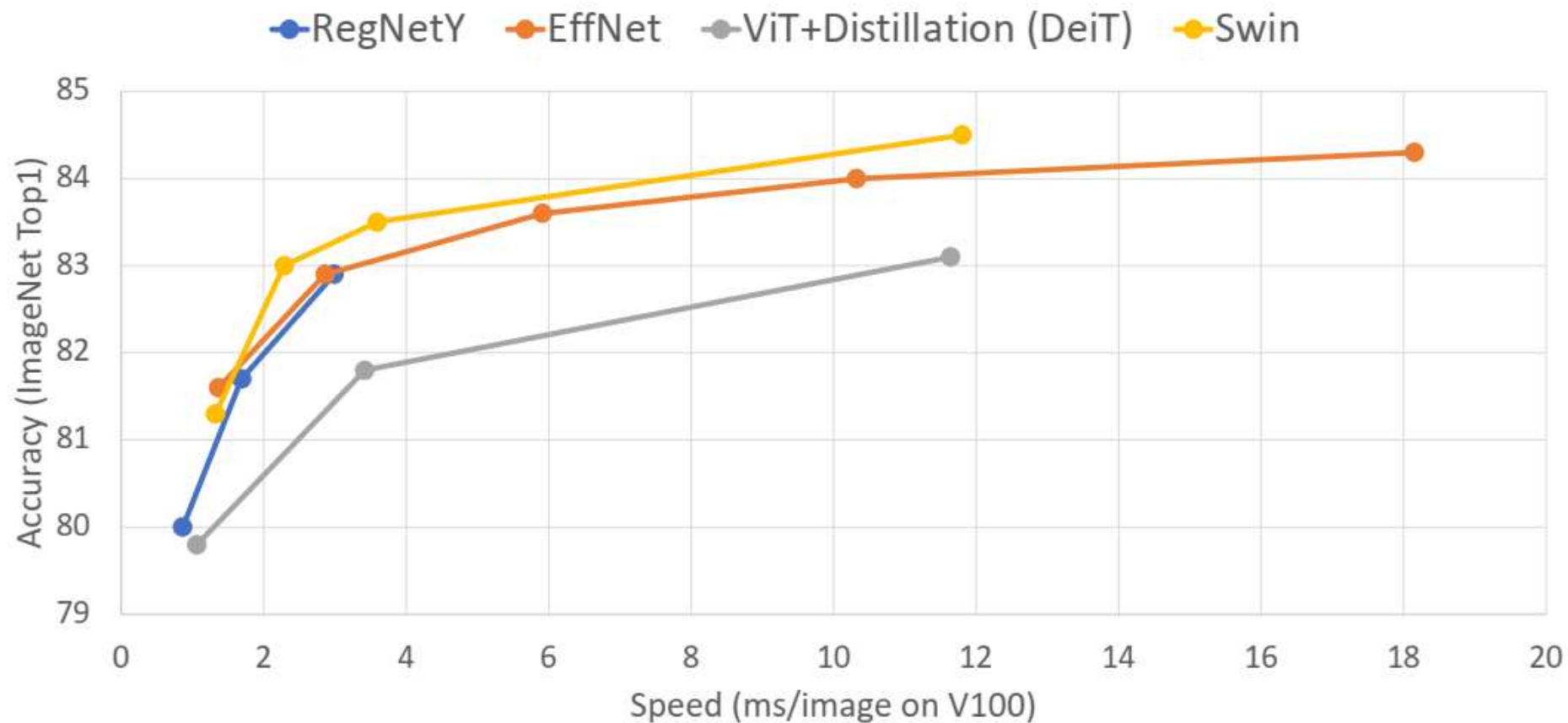


Block L+1: Shifted Windows

ViT añade embeddings posicionales a los tokens del input

Swin usa relative positions entre parches

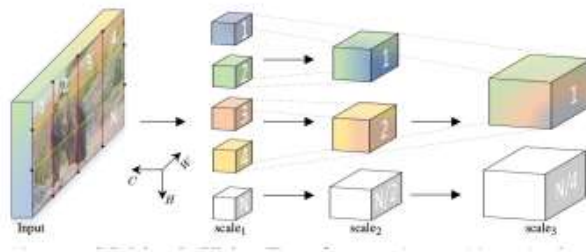
# Swin Transformers: Speed vs Accuracy





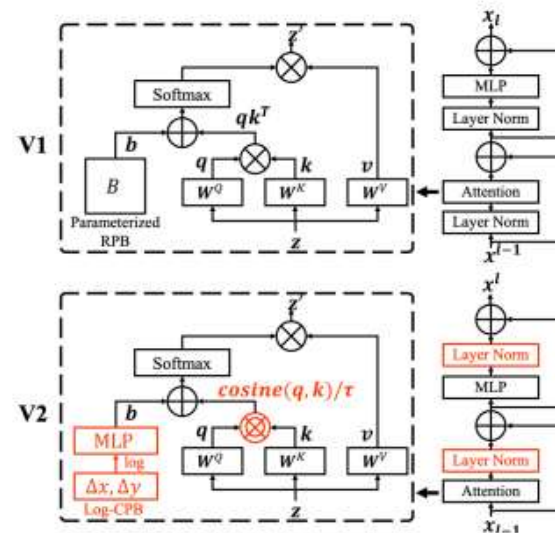
# Otros modelos jerárquicos

MViT



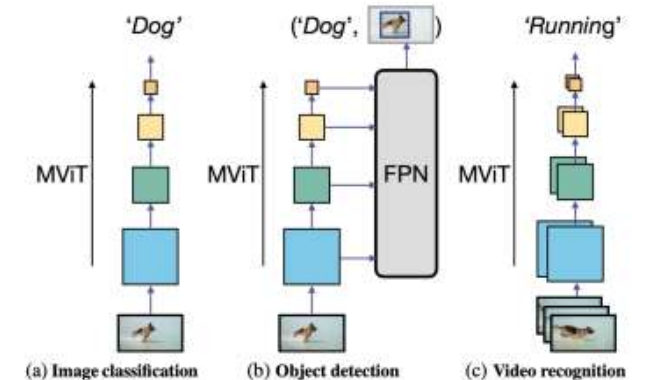
Fan et al, "Multiscale Vision Transformers", ICCV 2021

Swin-V2



Liu et al, "Swin Transformer V2: Scaling up Capacity and Resolution", CVPR 2022

Improved MViT



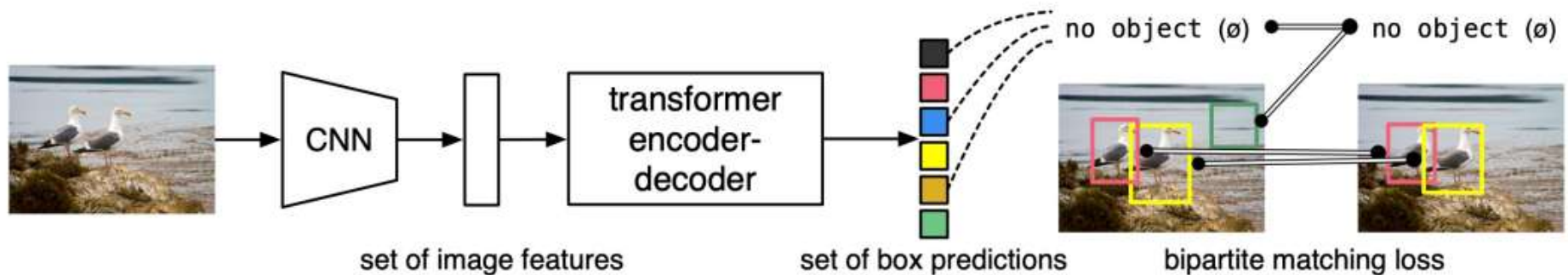
Li et al, "Improved Multiscale Vision Transformers for Classification and Detection", arXiv 2021



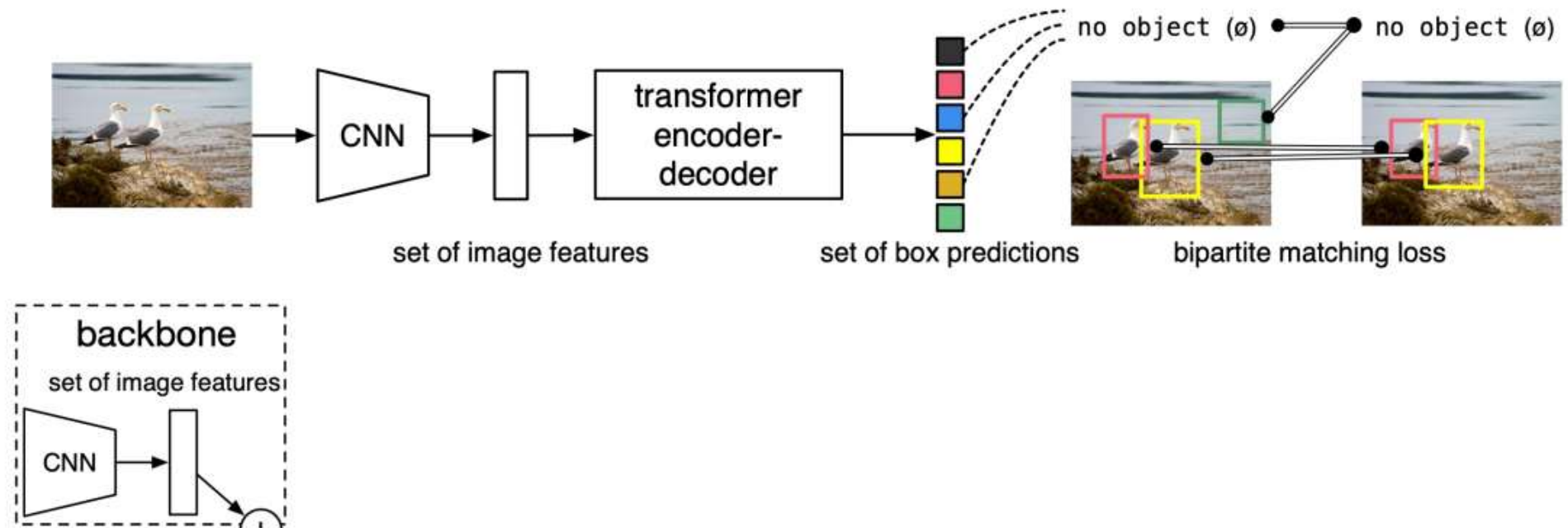
# Detección de objetos con DETR

Pipeline simple de detección de objetos: salida conjunto de bounding boxes del Transformer

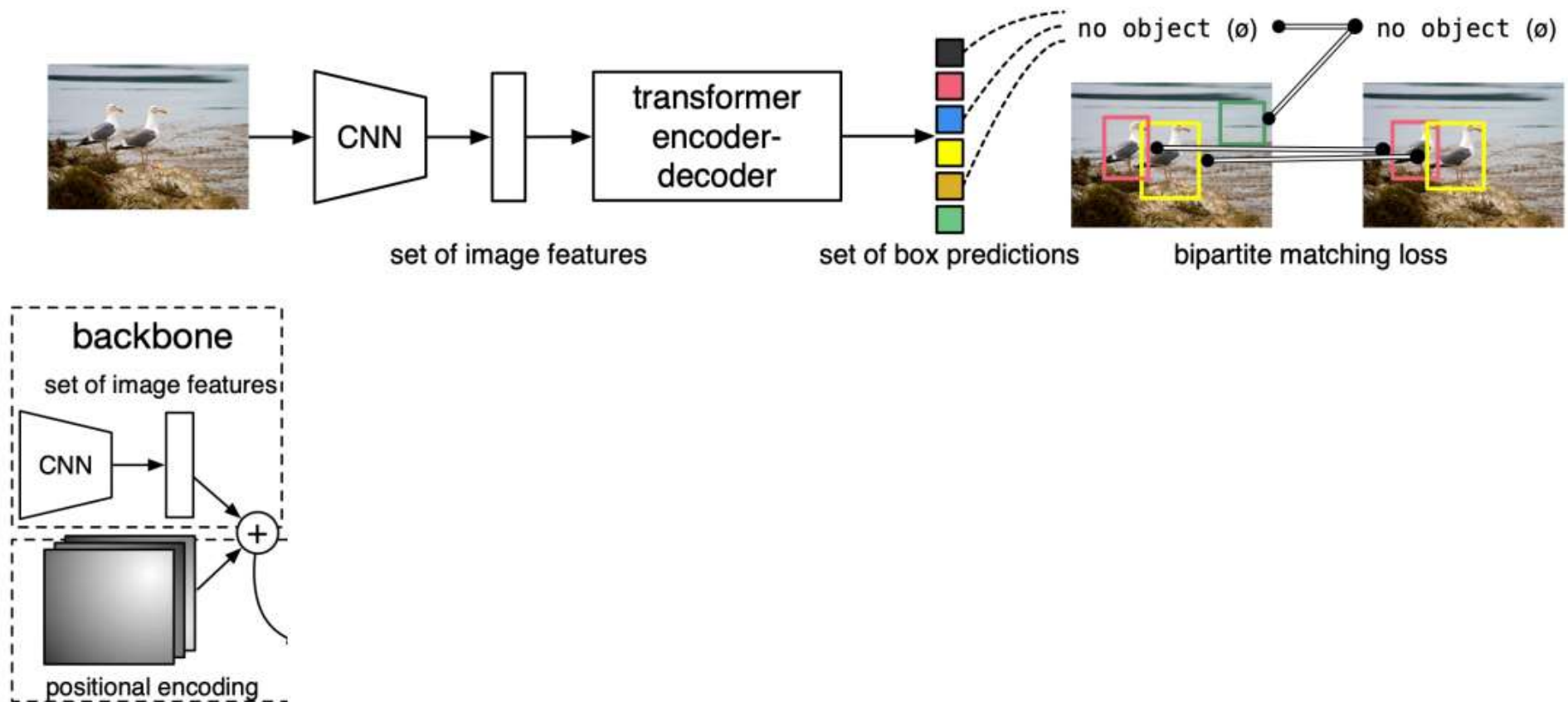
Matchear los boxes predichos a los boxes del groundtruth. Entrenar para hacer regresión de los boxes



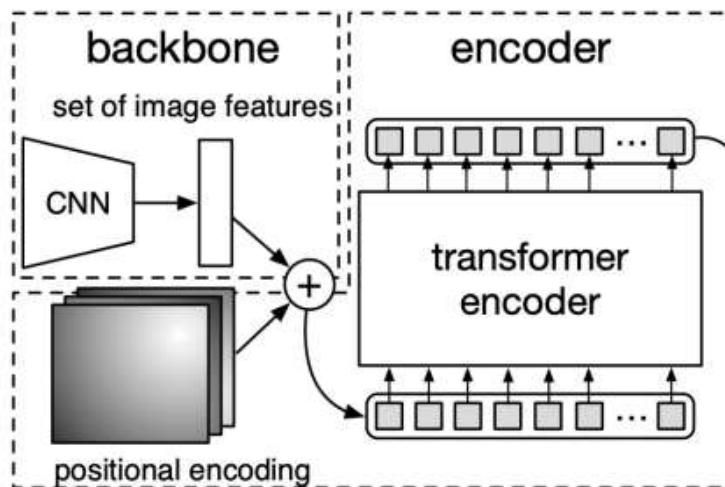
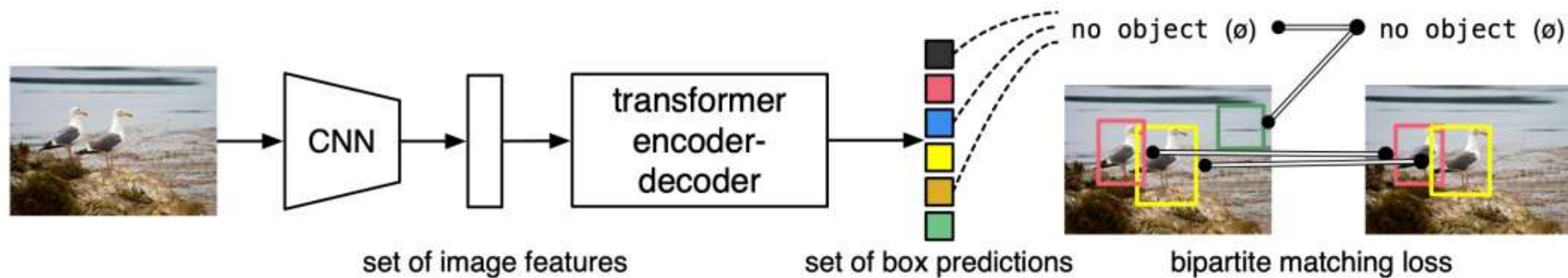
# Detección de objetos con DETR



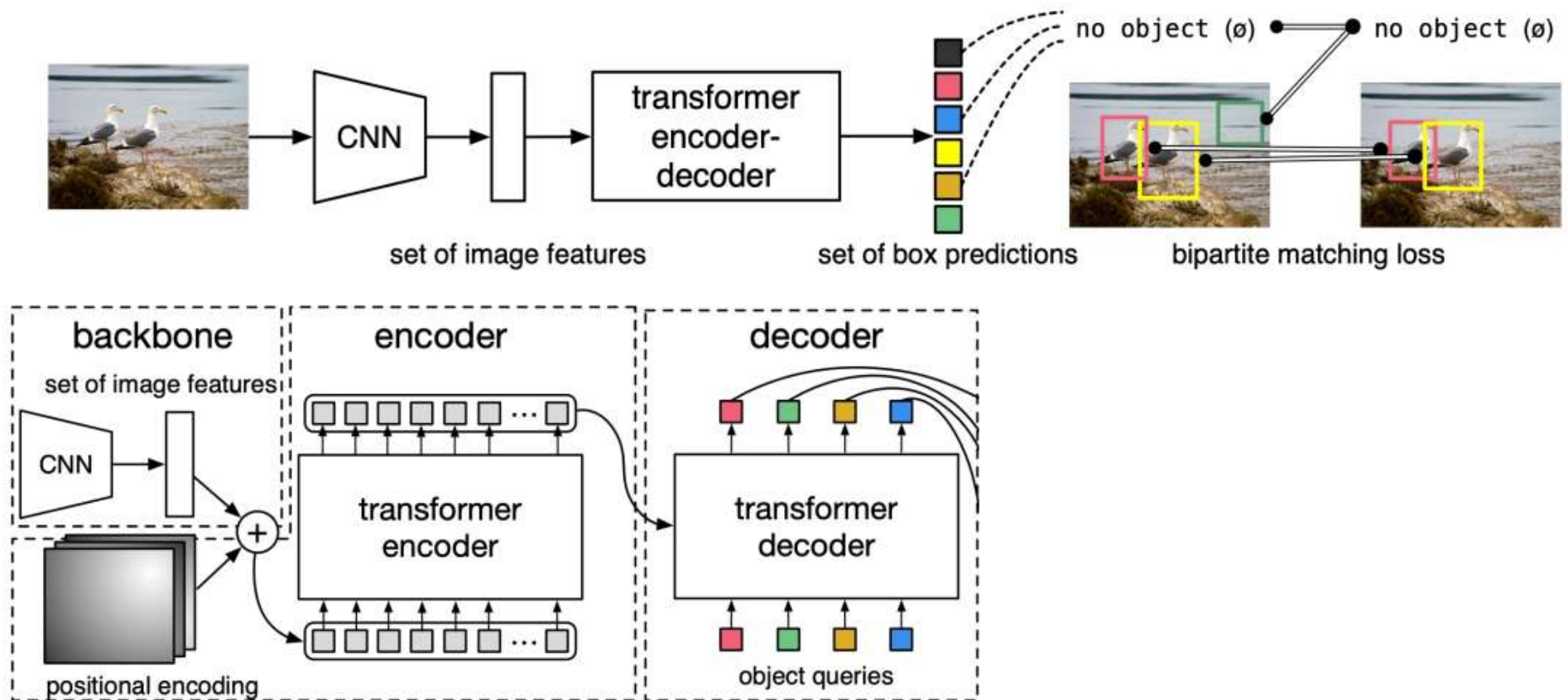
# Detección de objetos con DETR



# Detección de objetos con DETR



# Detección de objetos con DETR



# Detección de objetos con DETR

