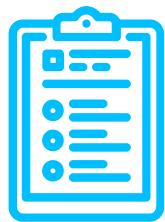


# Sesión 3.1

## *Visual Attention*

*ViT, Swin transformer, CrossViT*

**1.**



## **Atención** *humana*

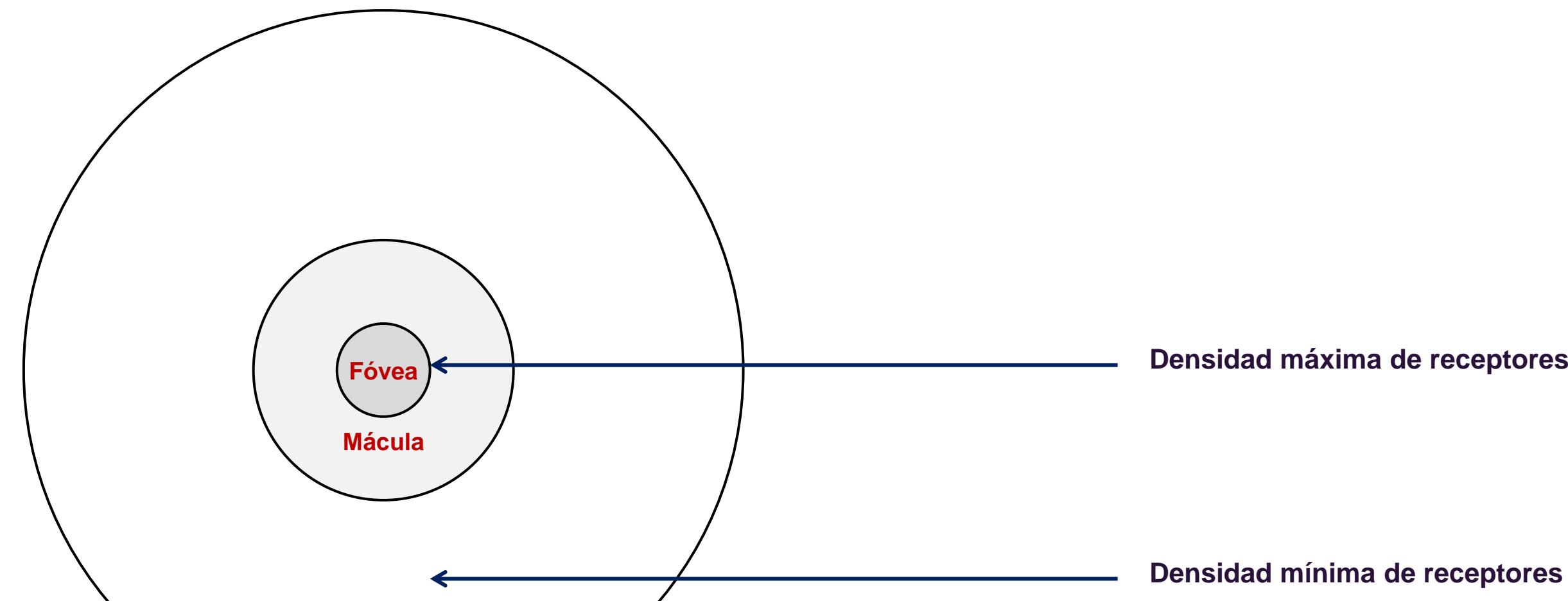




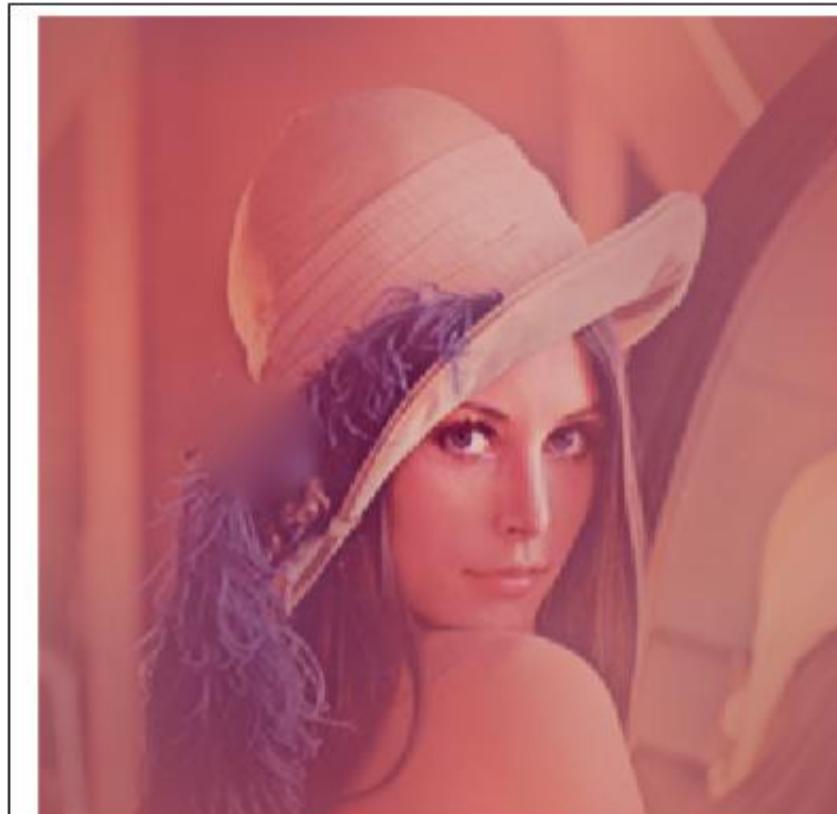




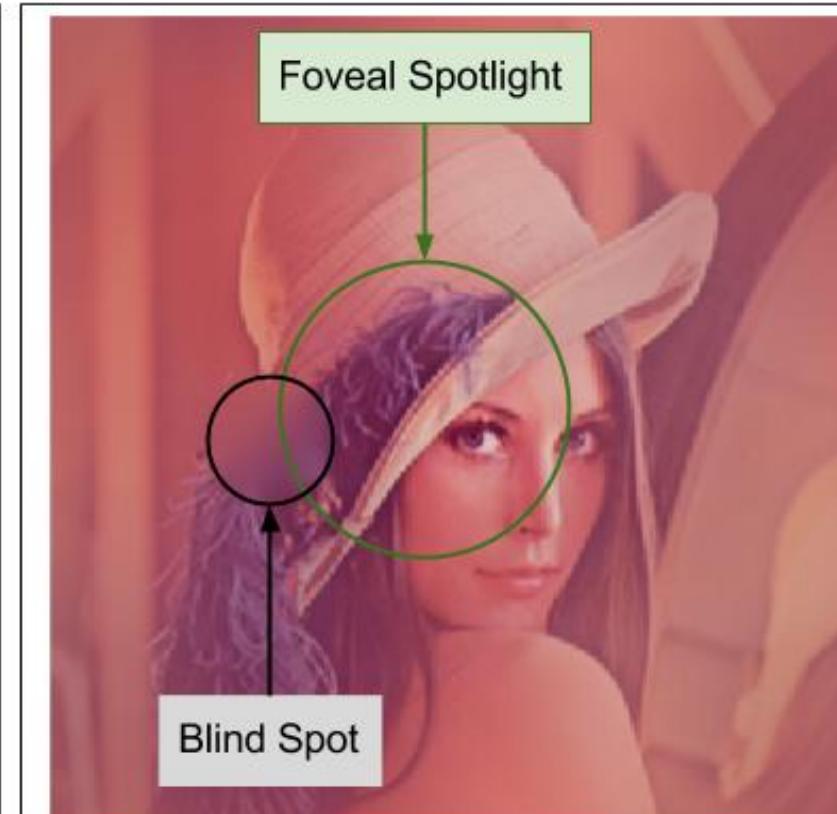
# Atención *visual humana*



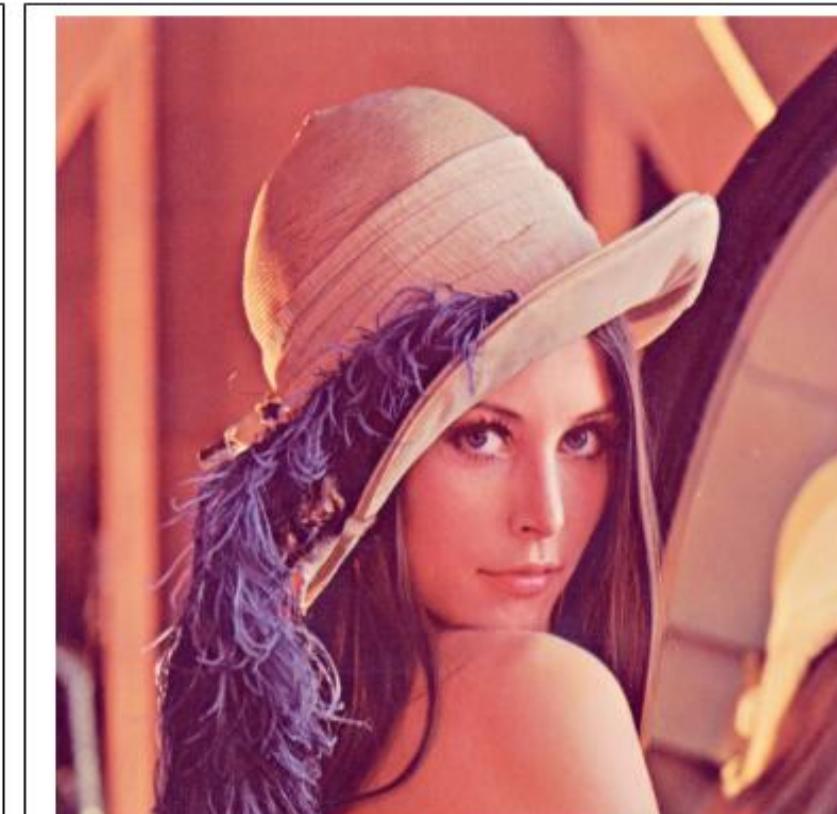
# Atención *visual humana*



What you *actually* see



What you *actually* see (annotated)



What you *think* you see



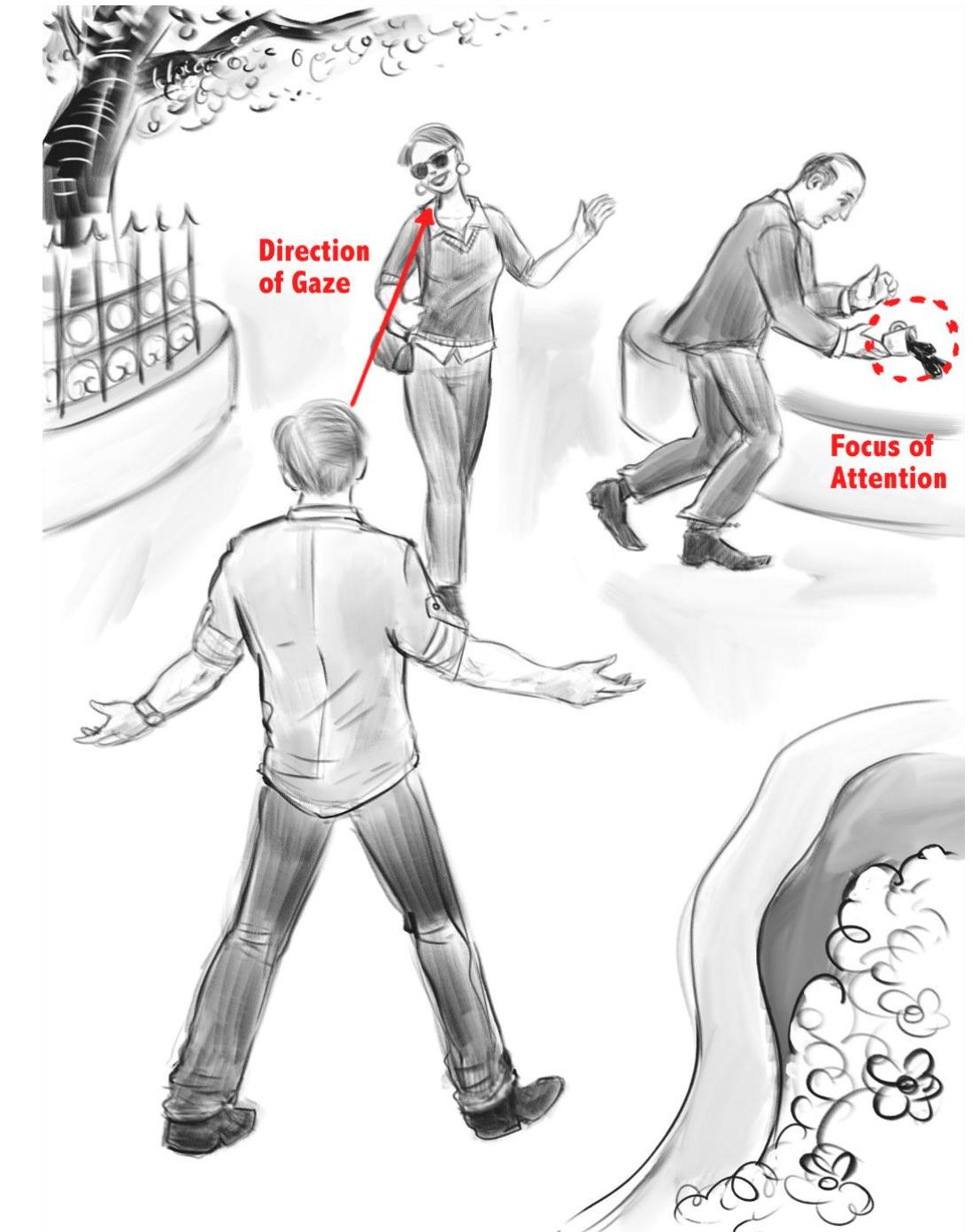
**TRANSFORMATEC**

# Modalidades de atención visual

Overt Attention

vs

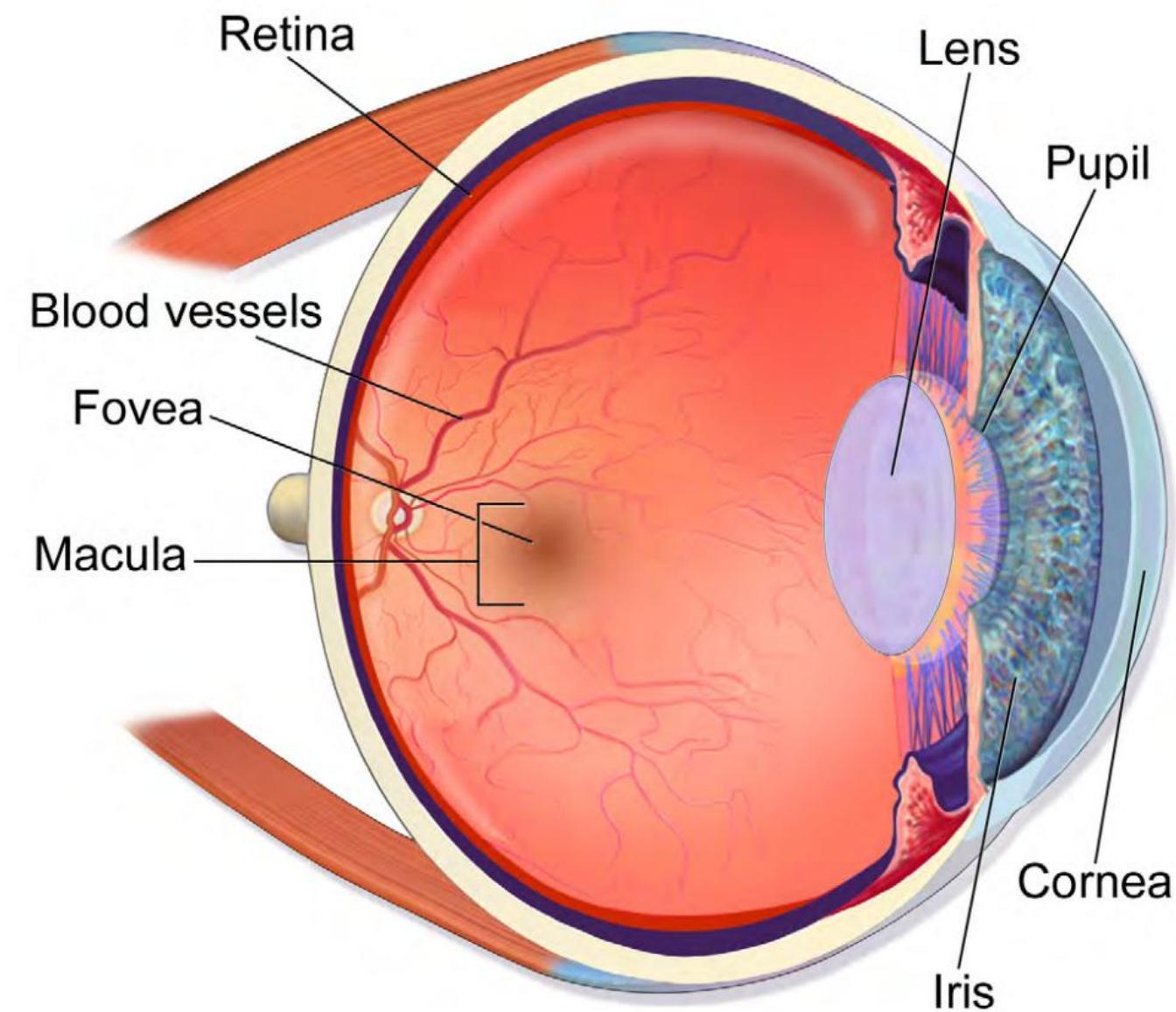
Covert Attention



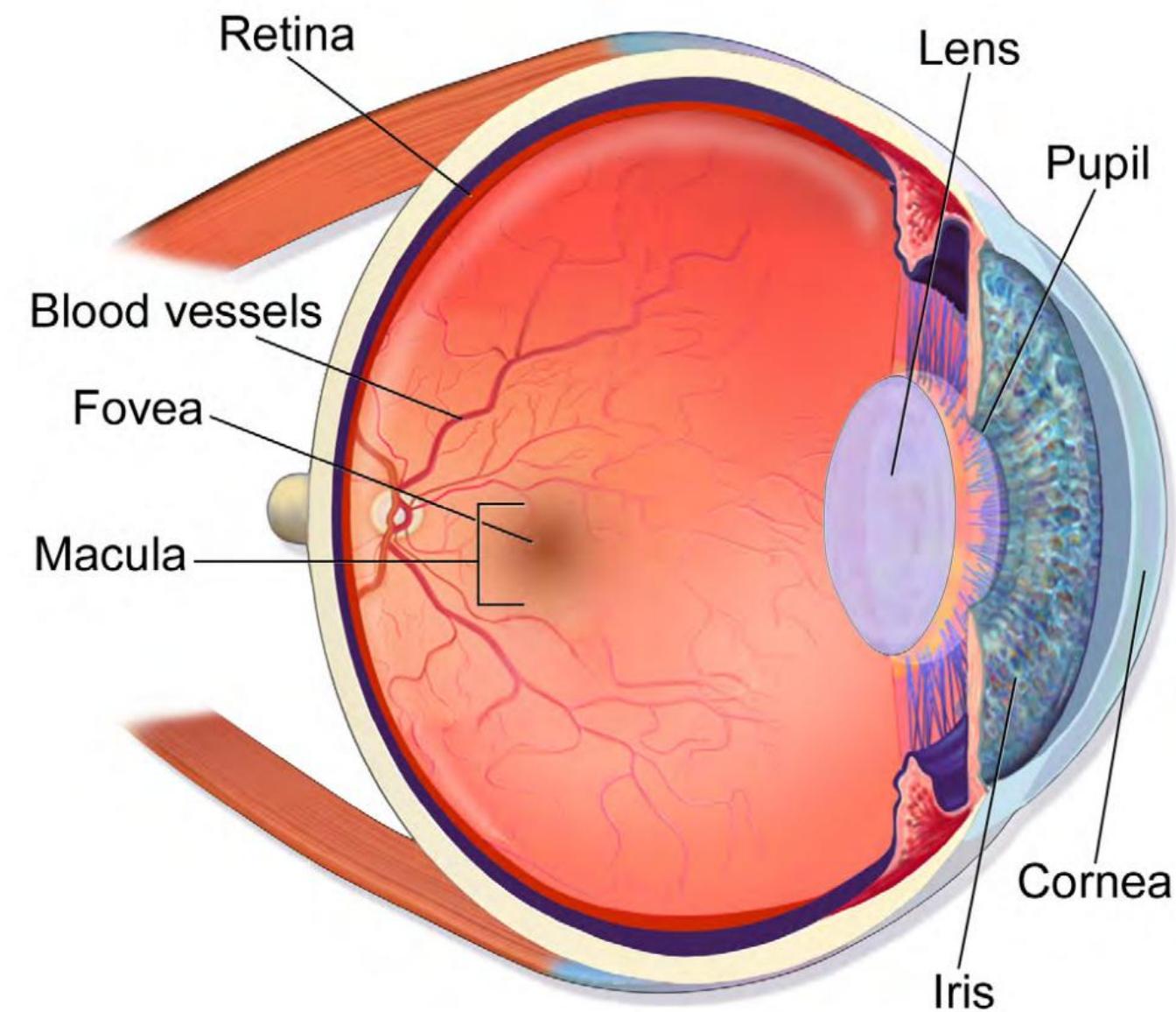
**TRANSFORMATEC**

Meng-Hao Guo et al. (2022) "Attention Mechanisms in Computer Vision: A Survey".  
Computational Visual Media. Springer.

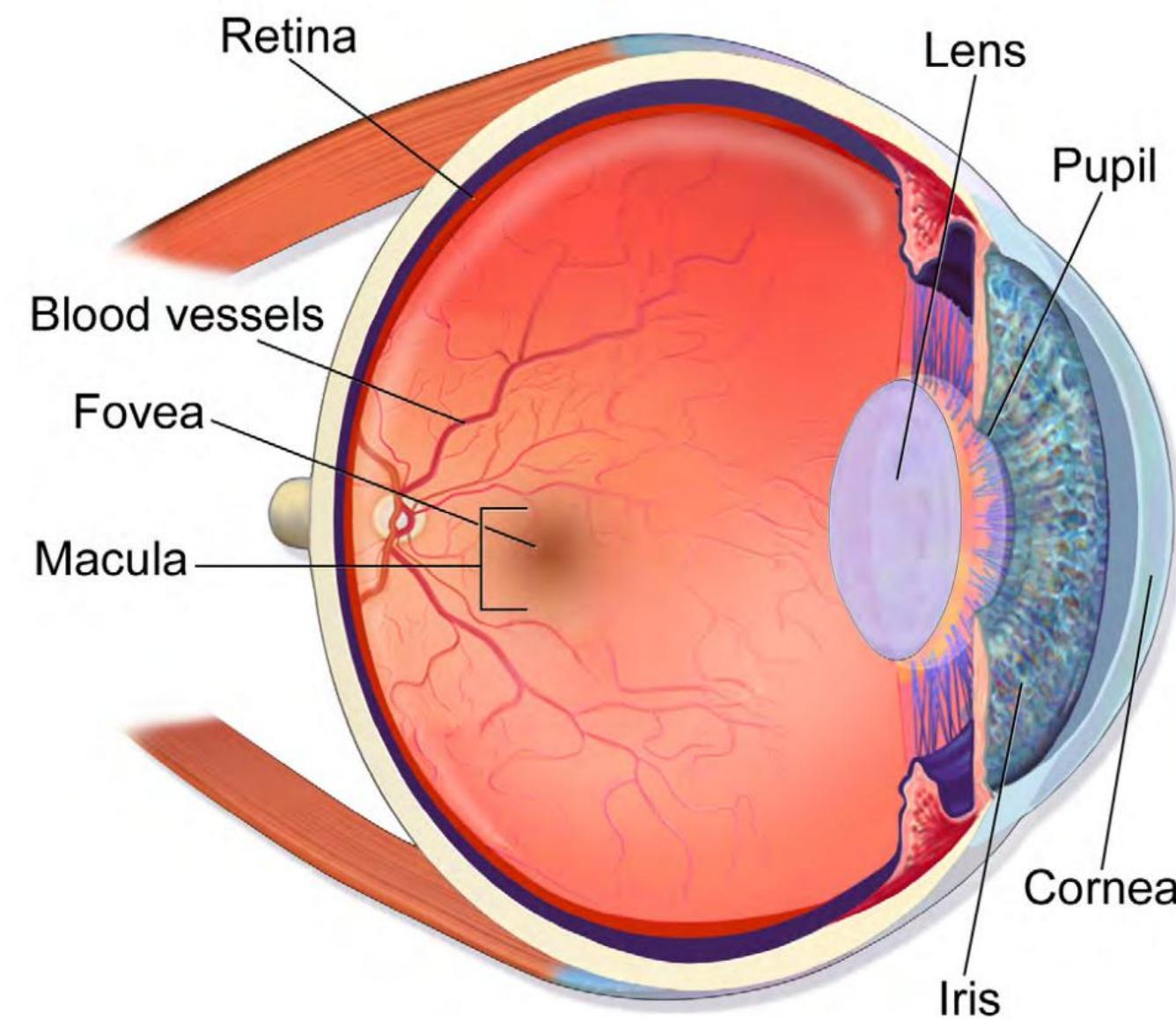
# Overt Attention



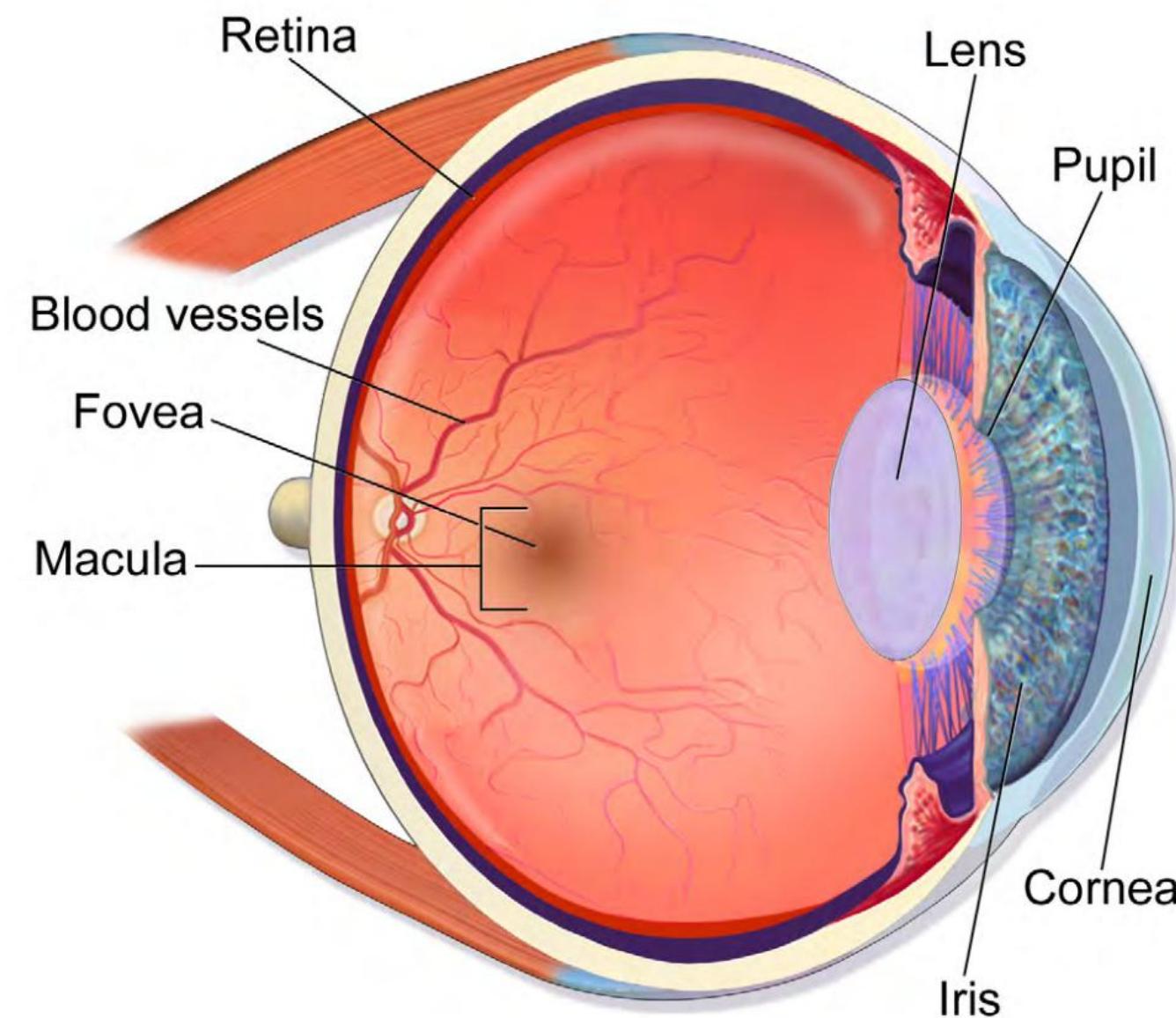
# Overt Attention



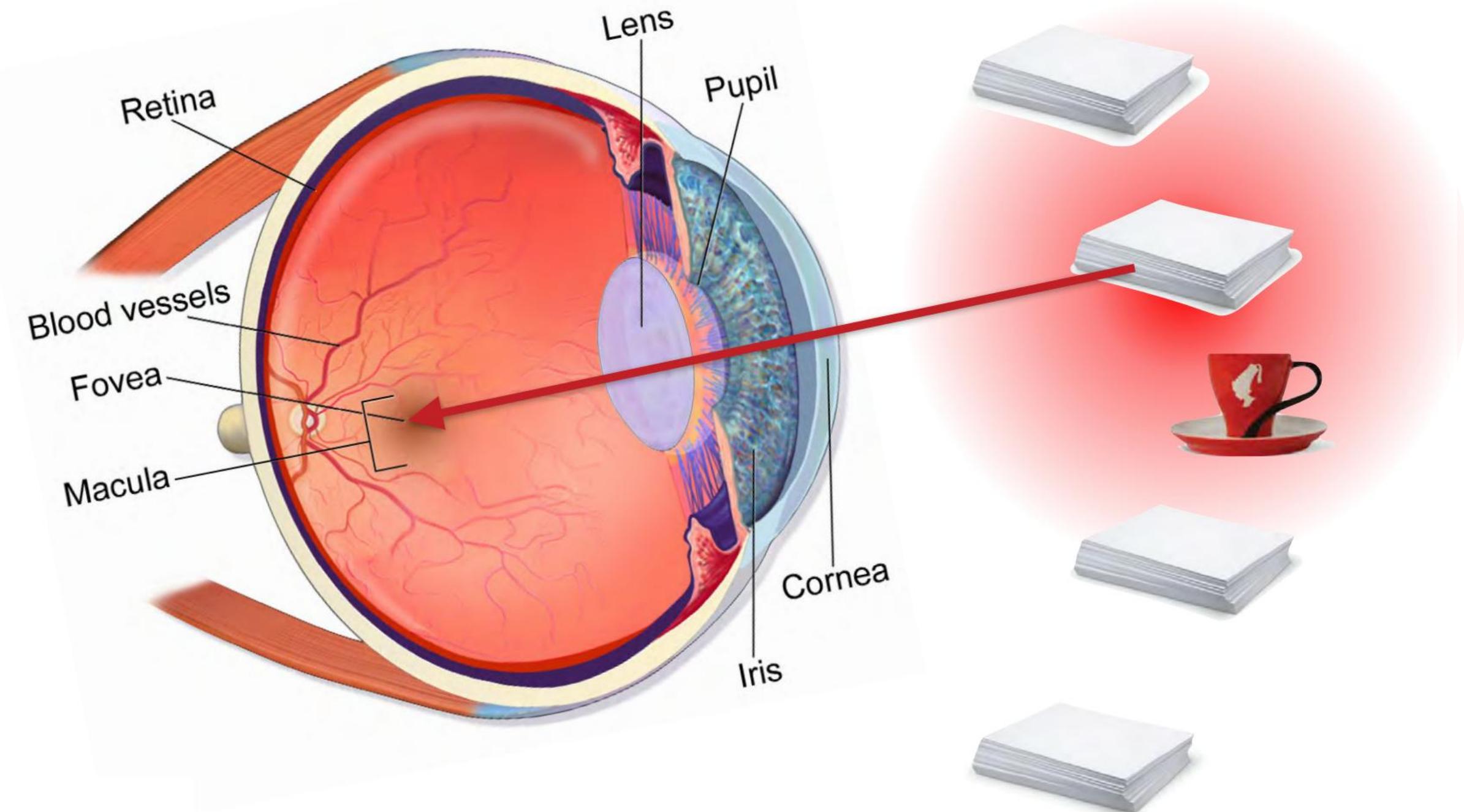
# Overt Attention



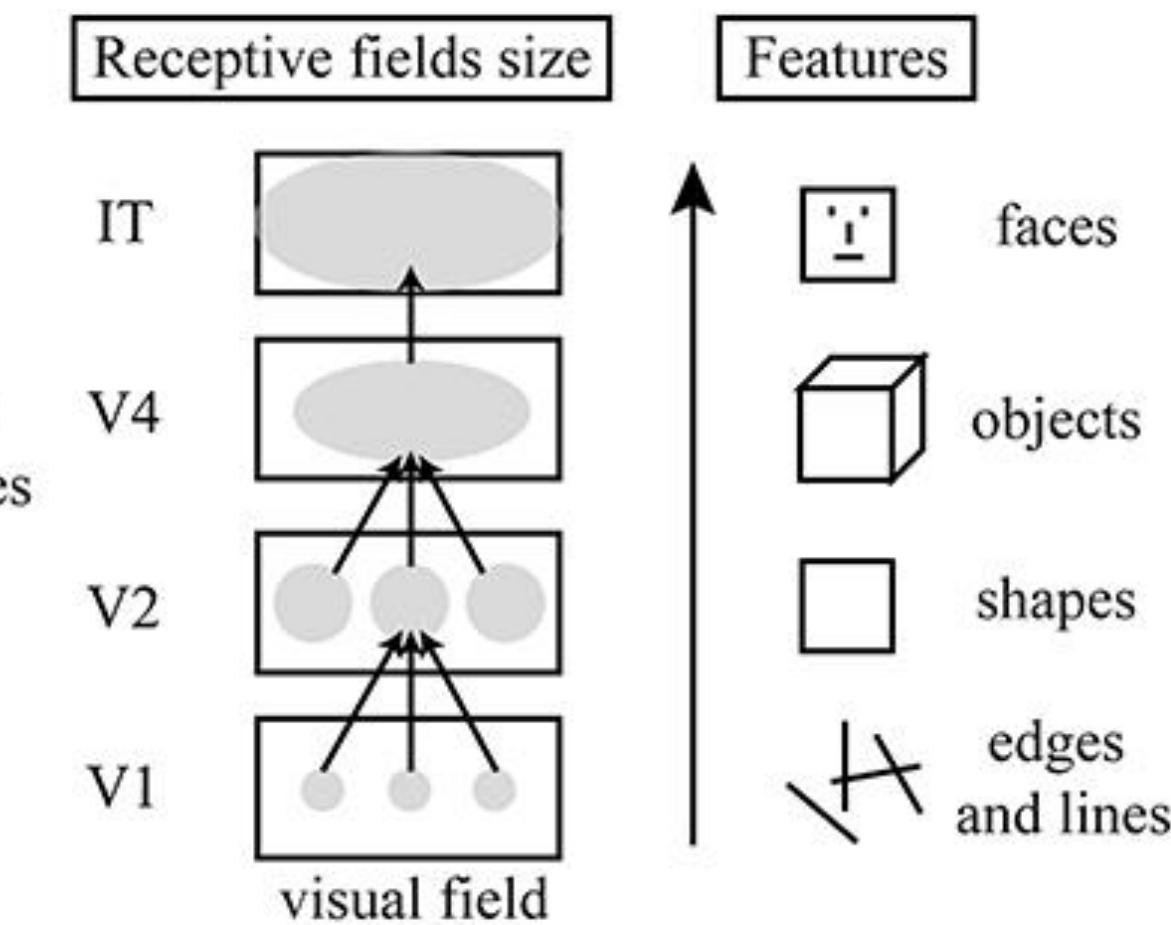
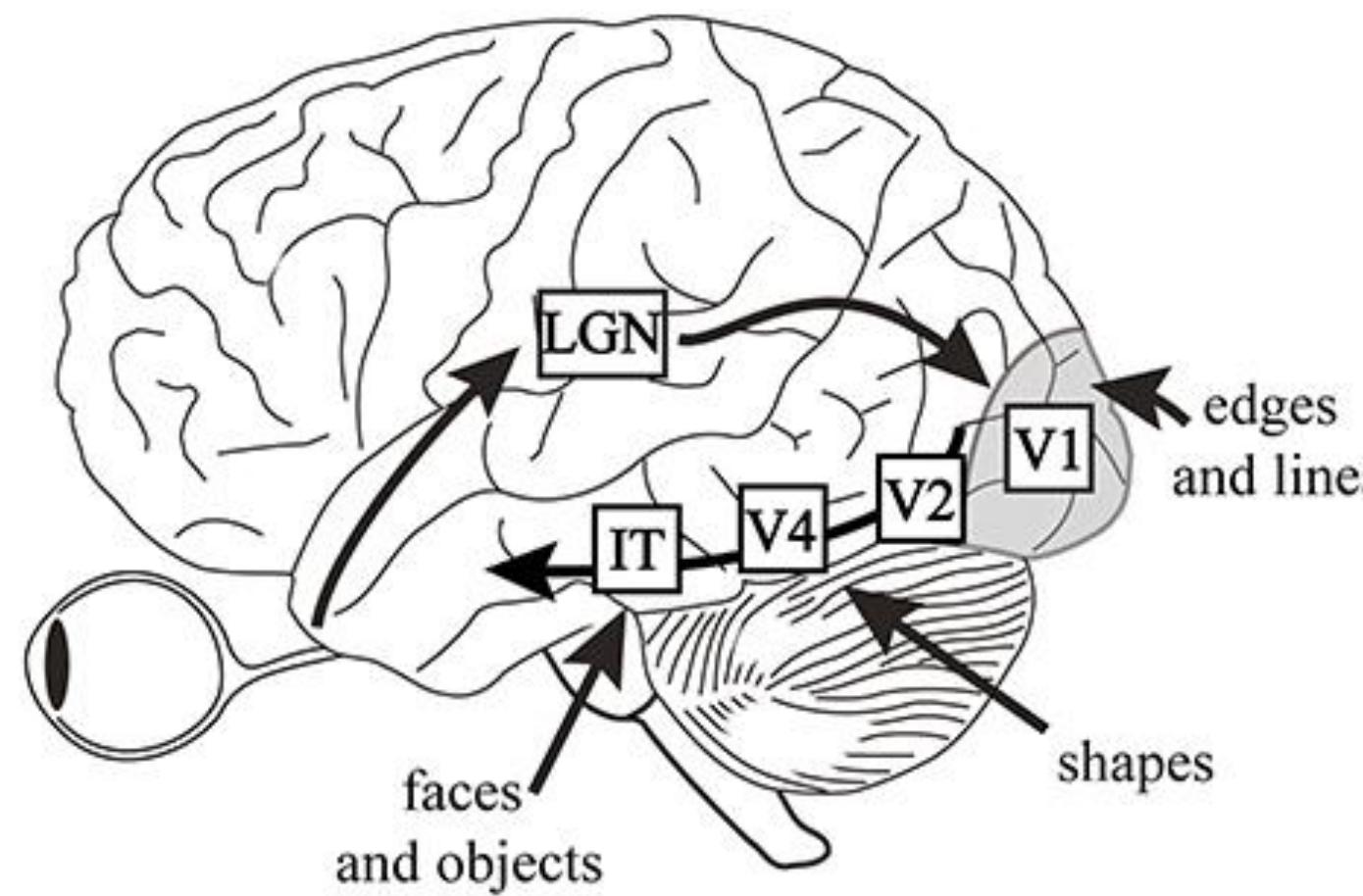
# Overt Attention



# Overt Attention

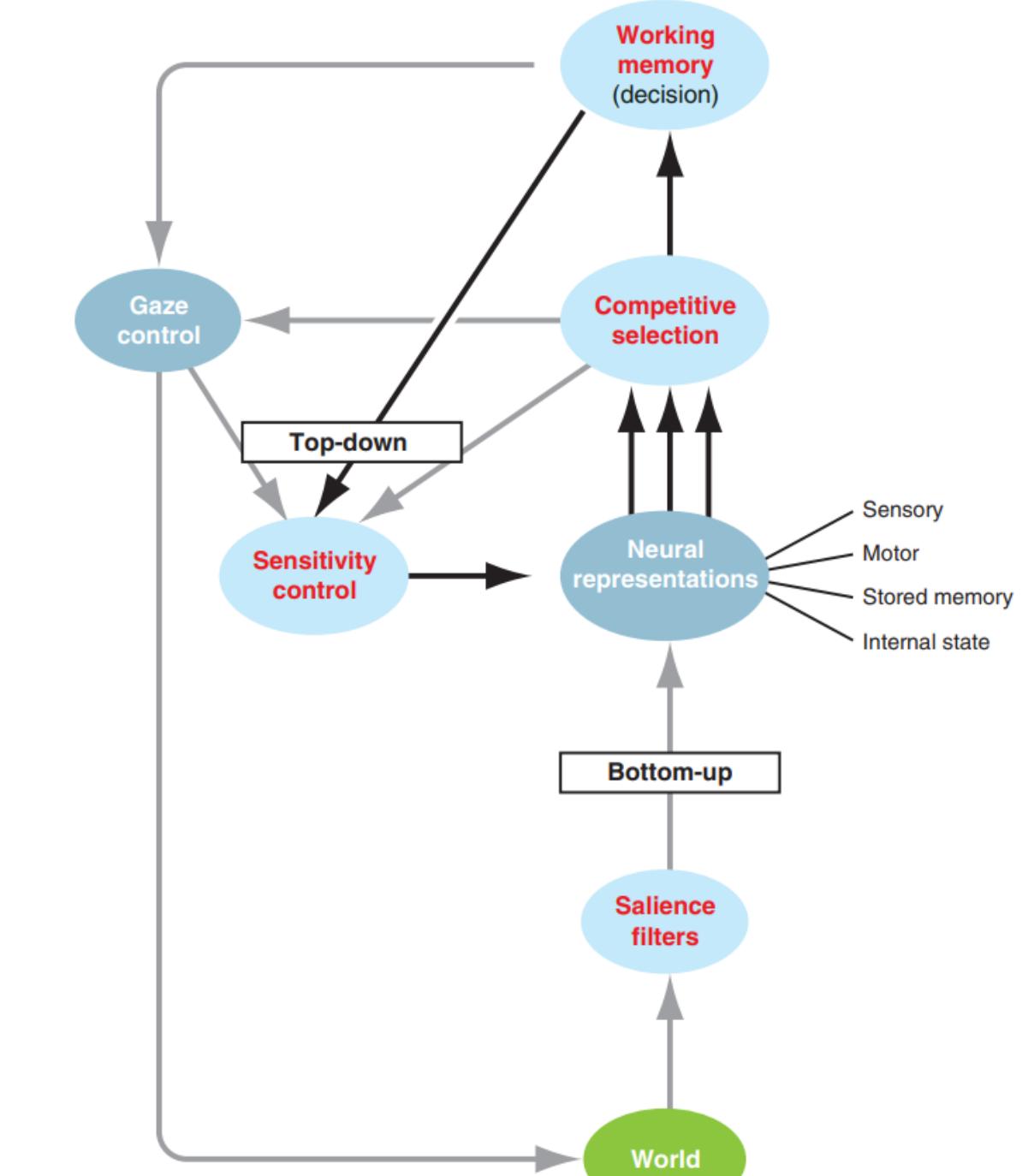
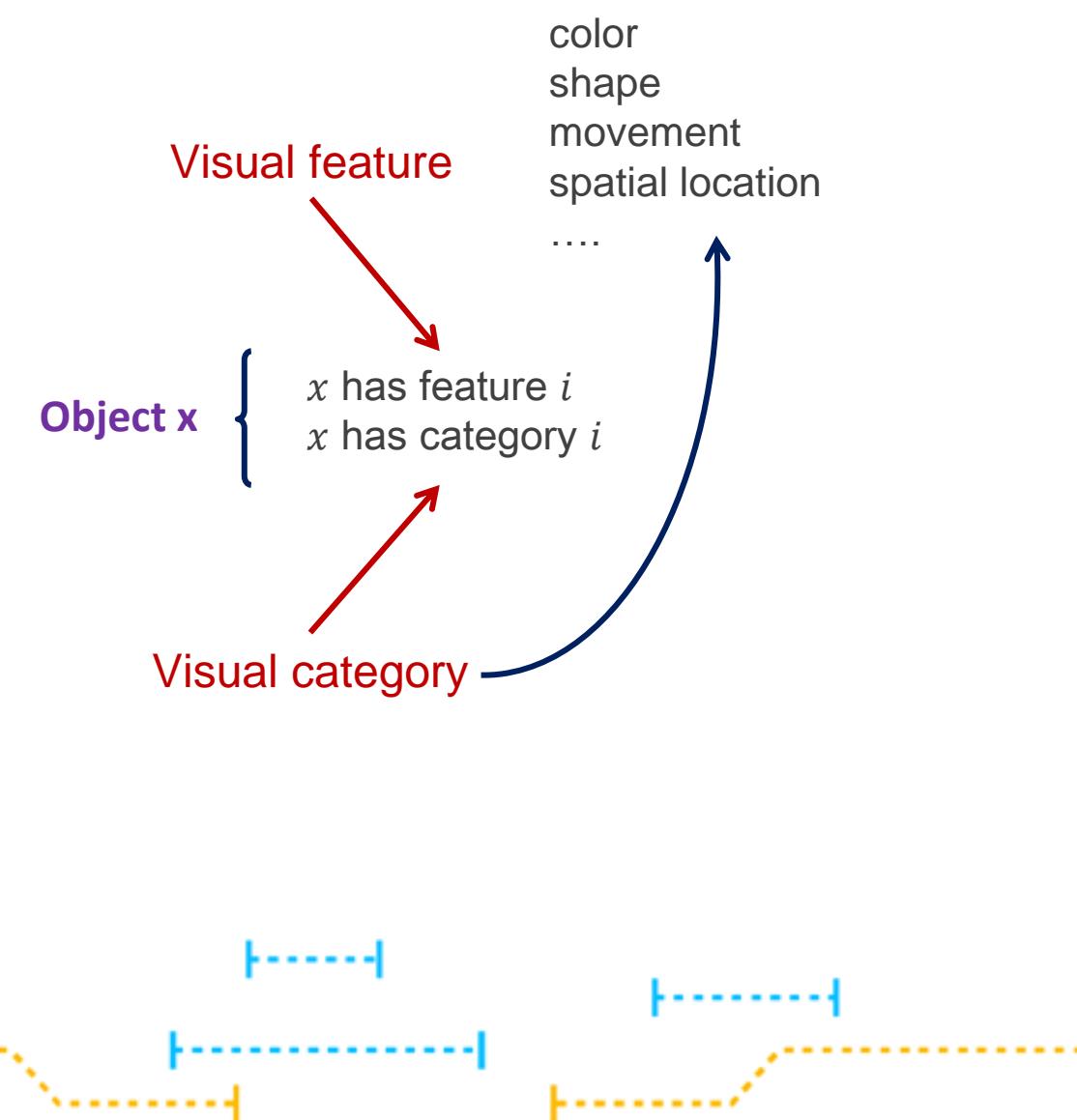


# Visual Cortex



# Covert Attention

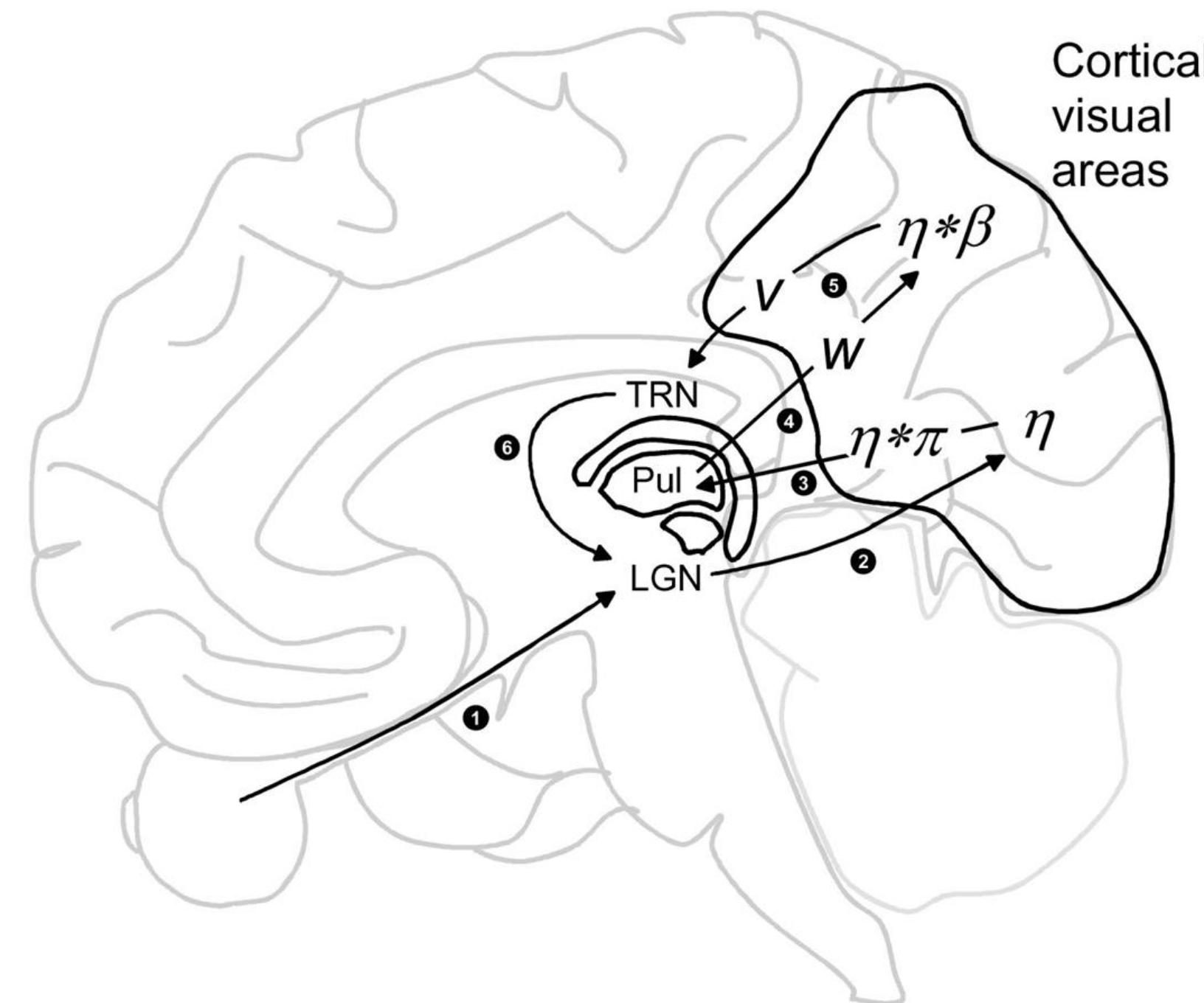
## Theory of Visual Attention



**TRANSFORMATEC**

Bundesen et al. (2011) "A neural theory of visual attention and short-term memory (NTVA)".  
Neuropsychologia, 49(6), 1446-1457.

# Covert Attention



**TRANSFORMATEC**

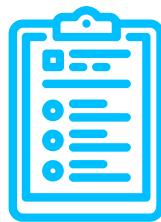
Bundesen et al. (2011) "A neural theory of visual attention and short-term memory (NTVA)".  
Neuropsychologia, 49(6), 1446-1457.

# Atención *Humana*

“La atención es un proceso cognitivo básico en el cerebro, y consiste en la selección dinámica de entidades de entrada mediante la ponderación adaptativa acorde a la importancia de cada entrada.”



**2.**



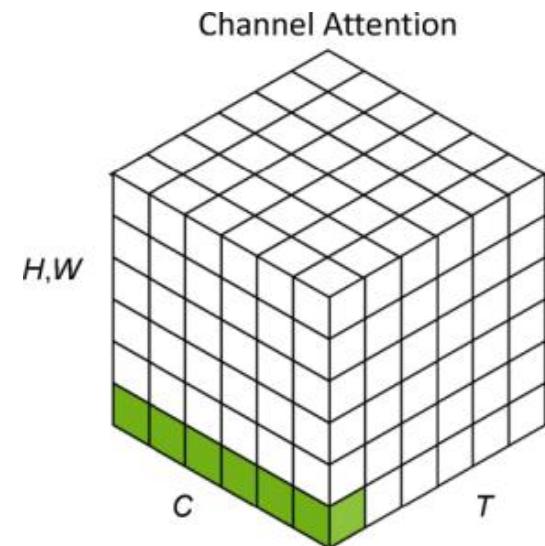
## **Beyond** *Transformers*

**TRANSFORMATEC**

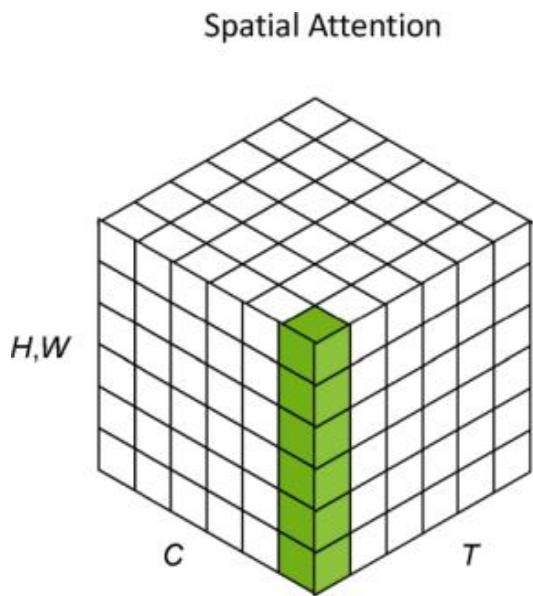
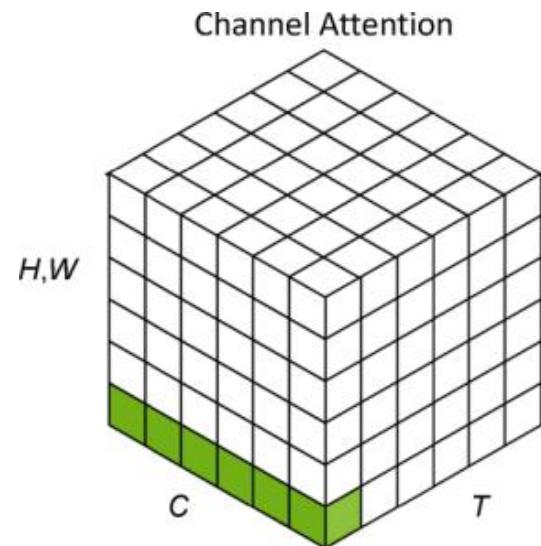
> Reinventa el mundo <



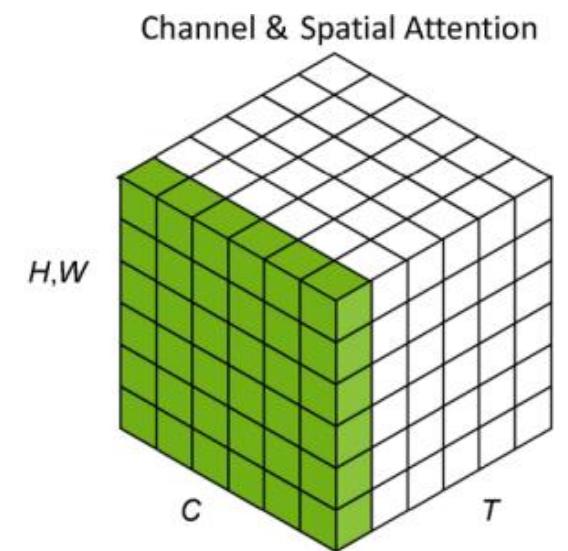
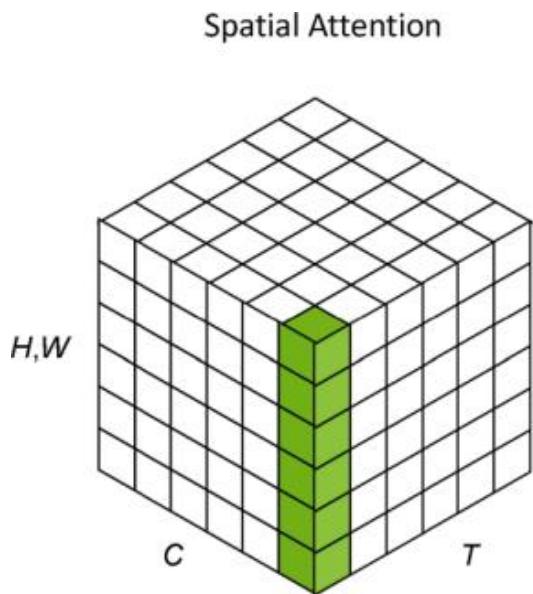
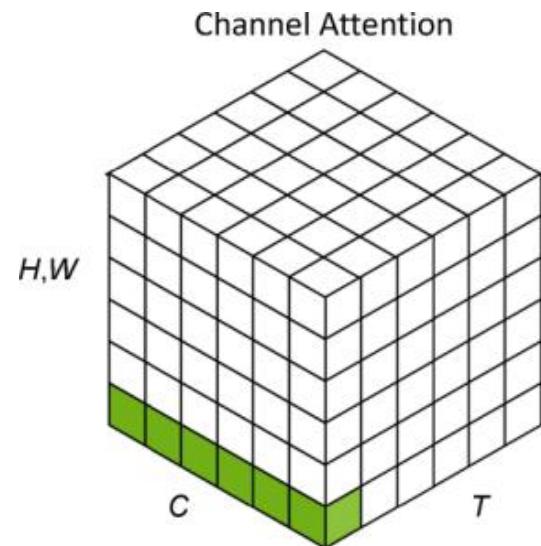
# Modalidades de atención



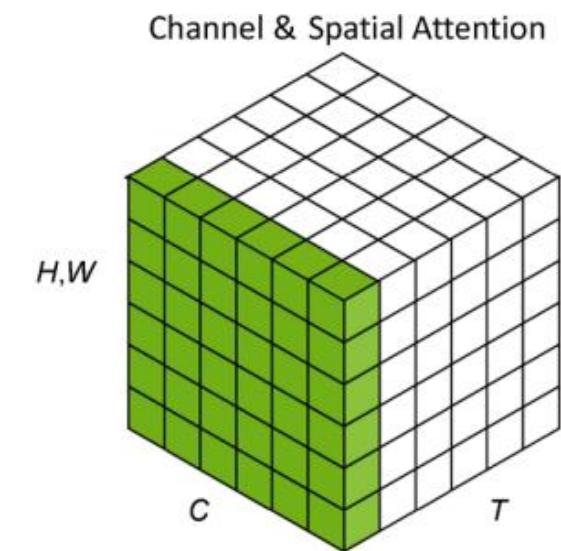
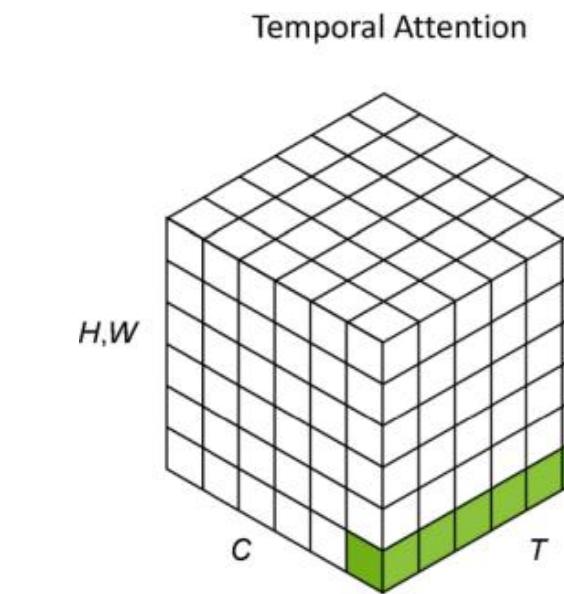
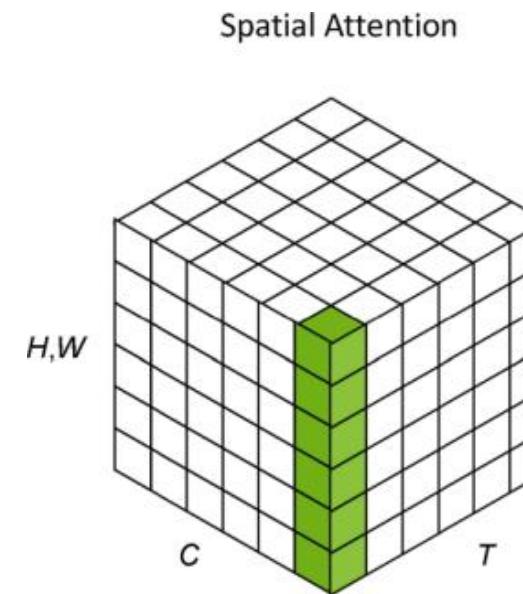
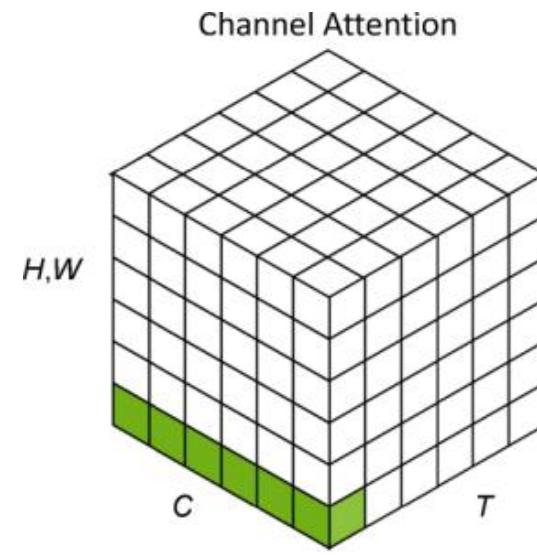
# Modalidades de atención



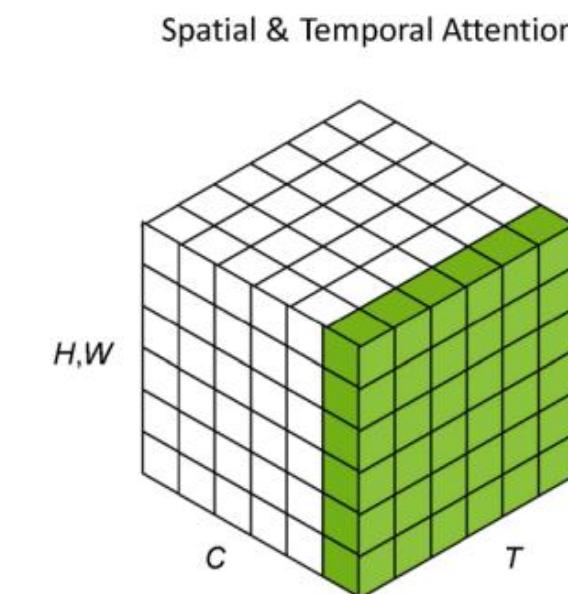
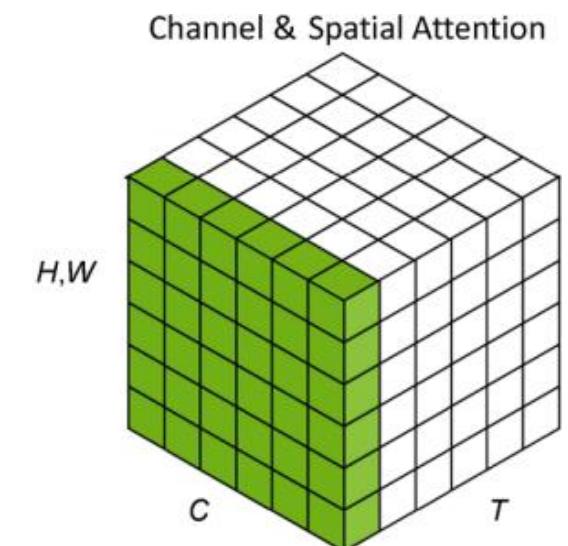
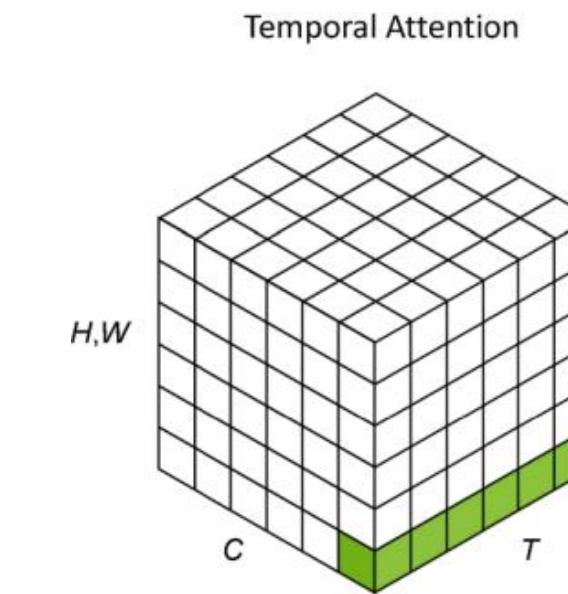
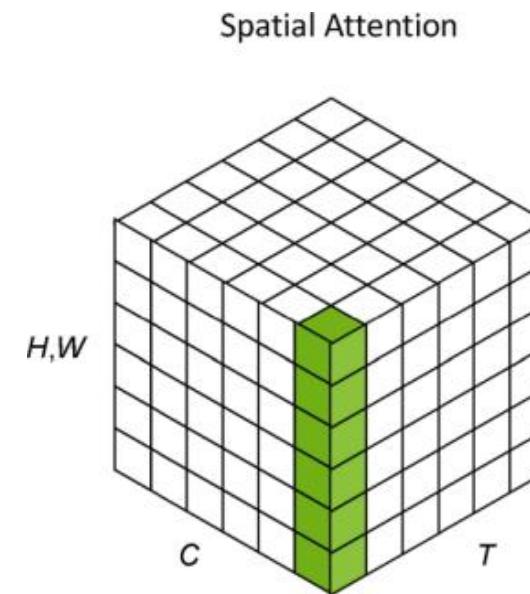
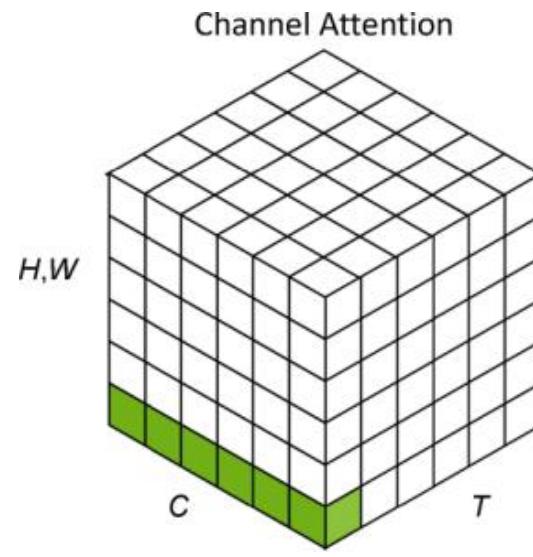
# Modalidades de atención



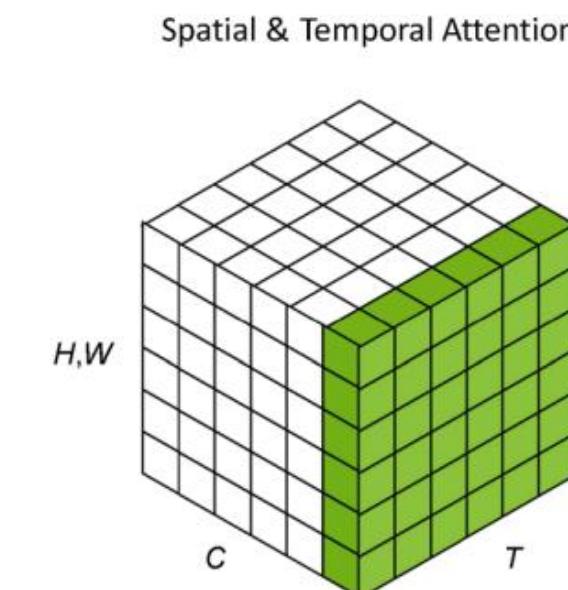
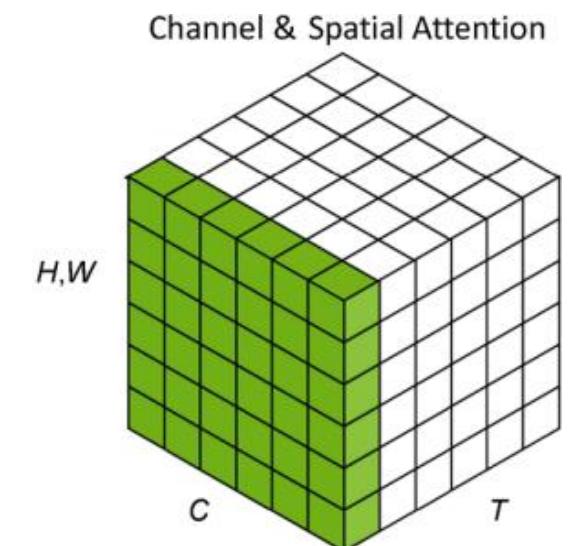
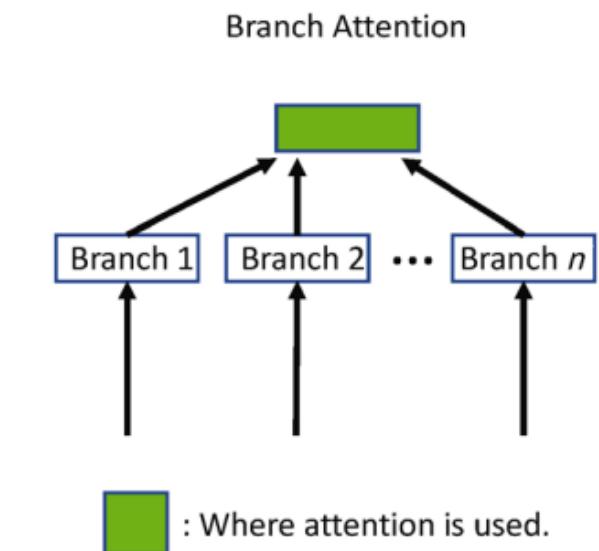
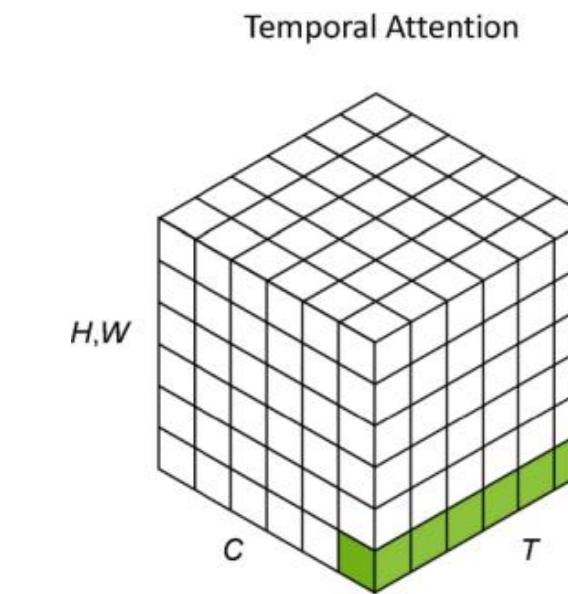
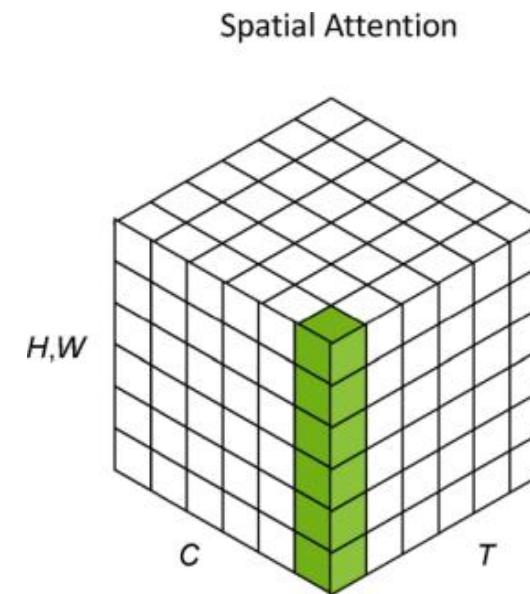
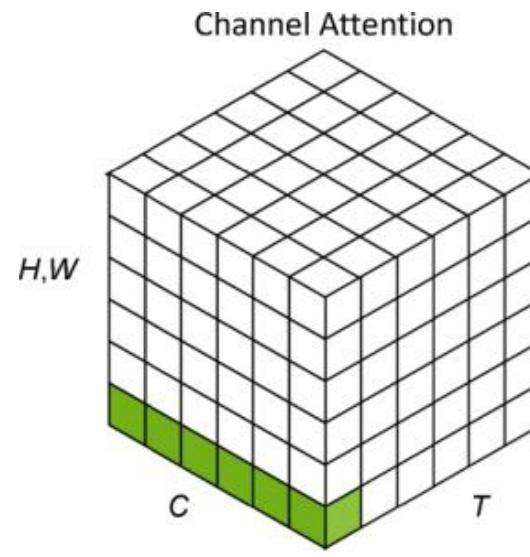
# Modalidades de atención



# Modalidades de atención



# Modalidades de atención



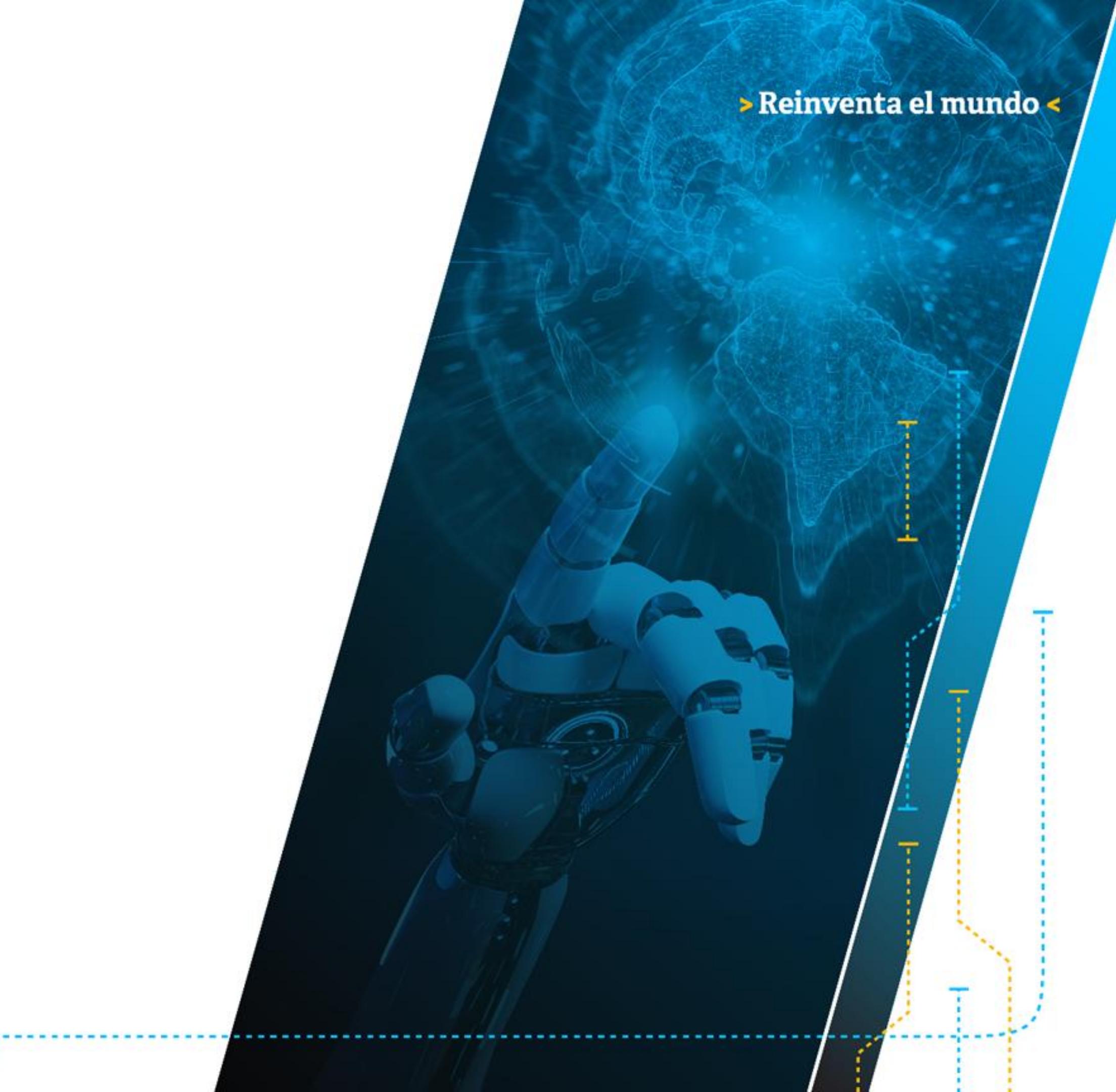
**3.**



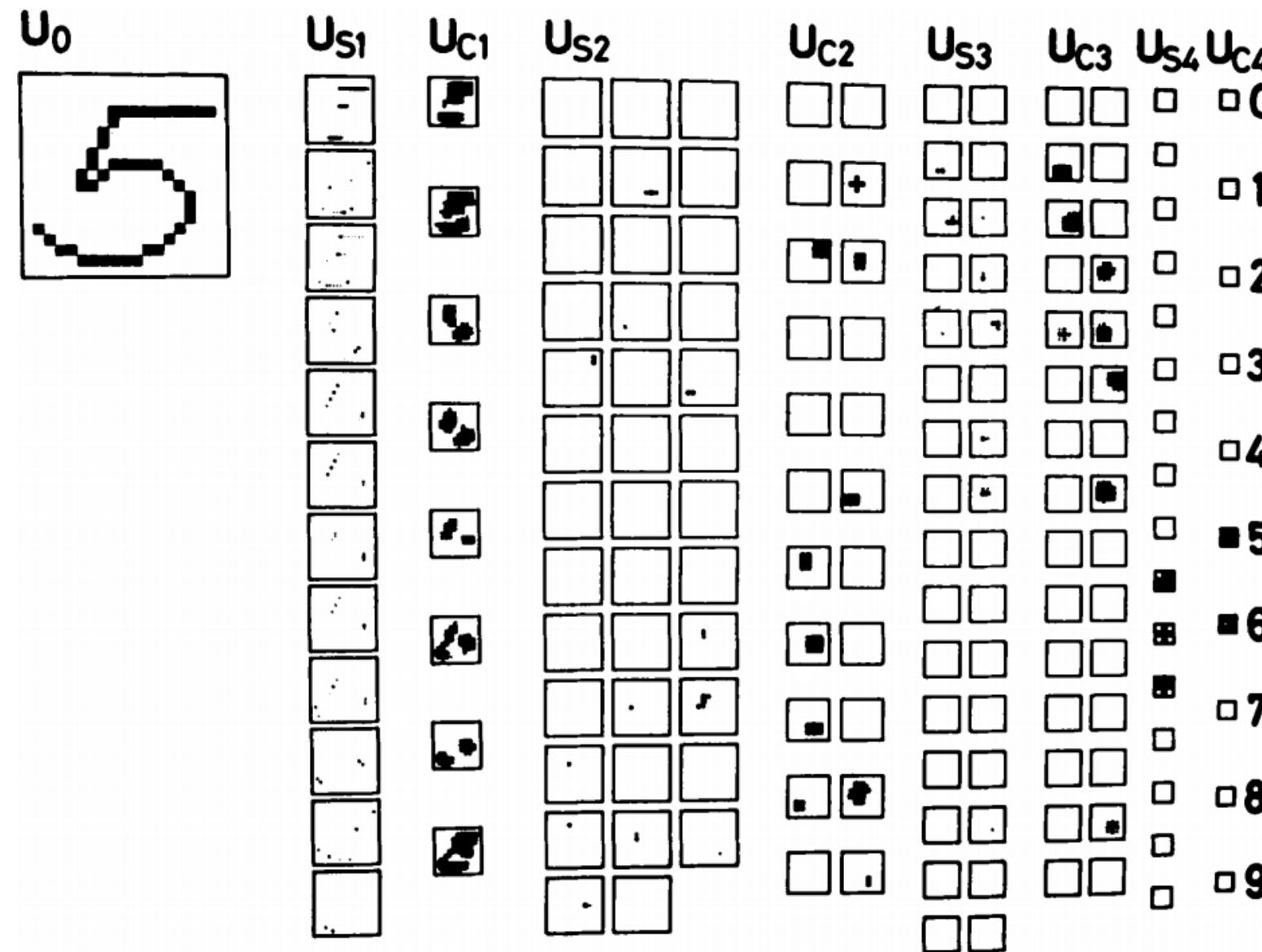
**Hard** attention

**TRANSFORMATEC**

> Reinventa el mundo <



# Neocognitron

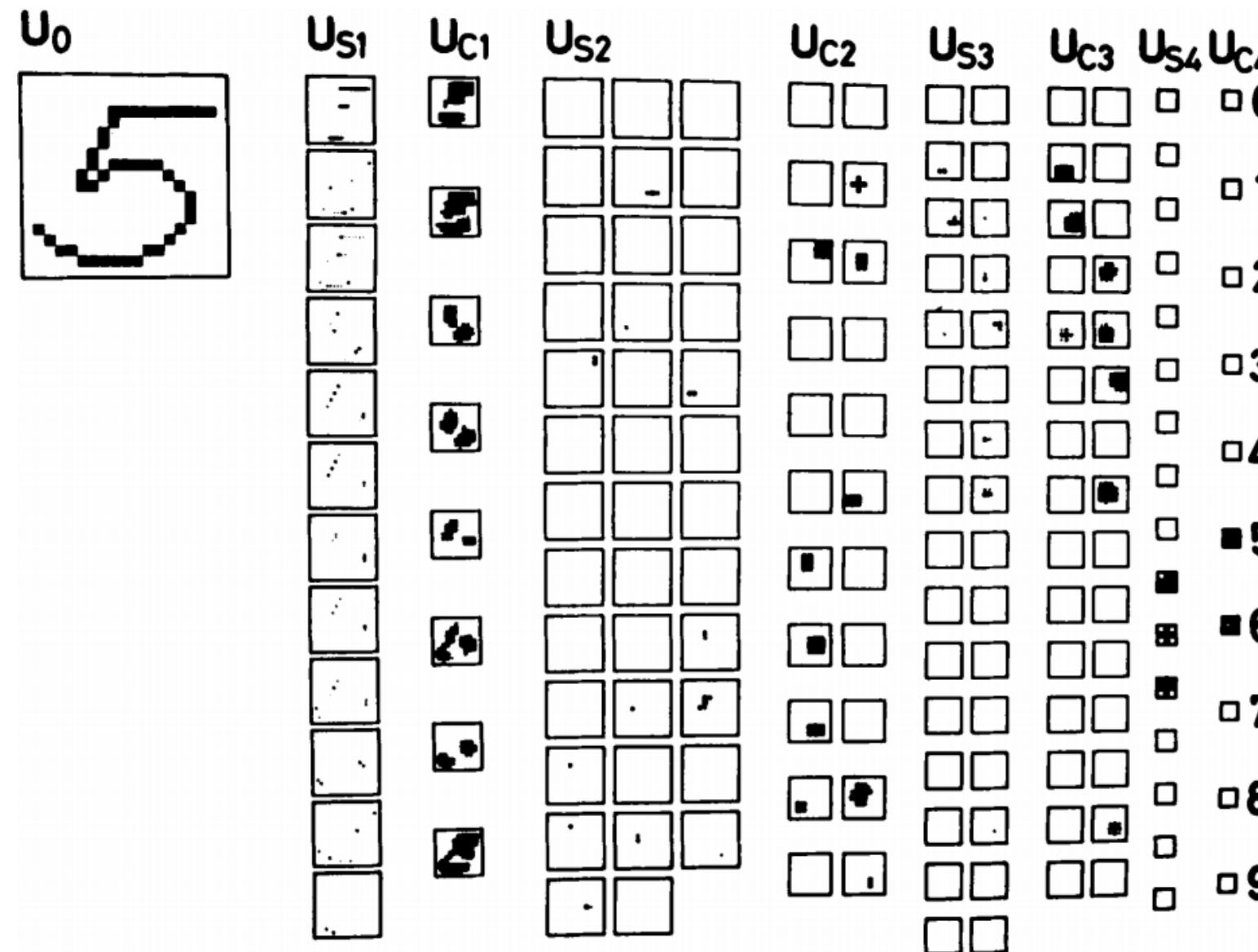


Kunihiko Fukushima



Fukushima, Kunihiko and Sei Miyake (1980) "Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position".  
Biological Cybernetics .

# Neocognitron



**S-cell:**

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

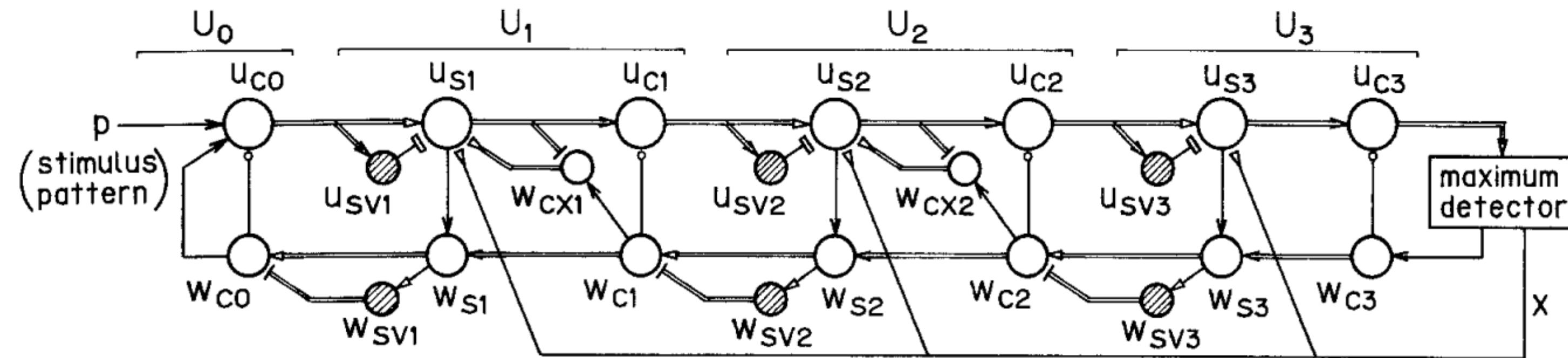
**C-cell:**

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

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



# Neocognitron



## connections

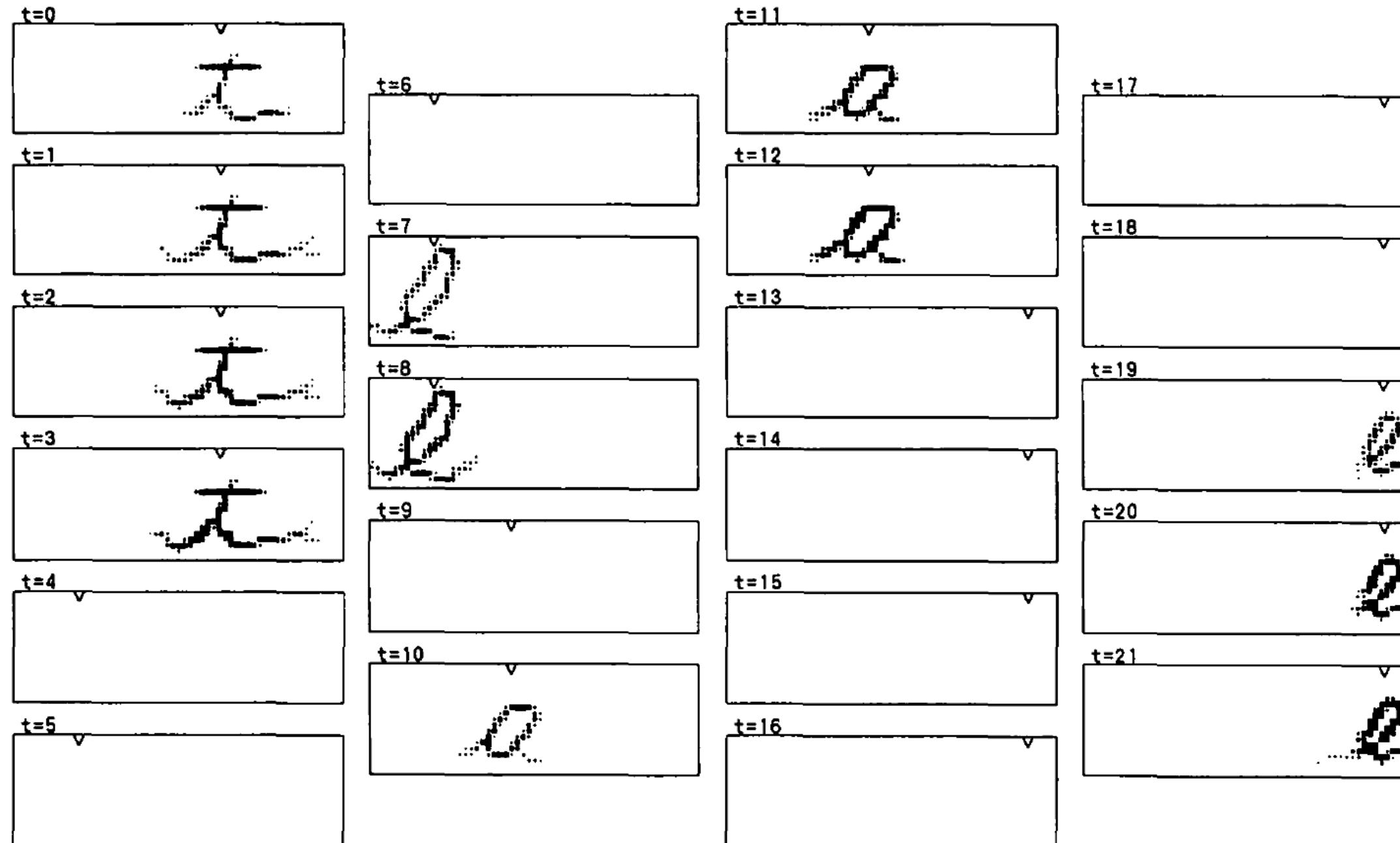
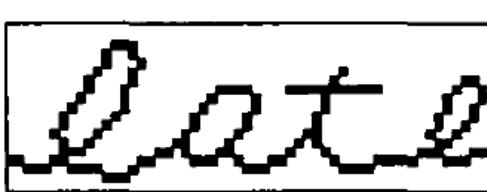
- converging or diverging (between two groups of cells)
- one-to-one connection (between two corresponding cells)

## synapses

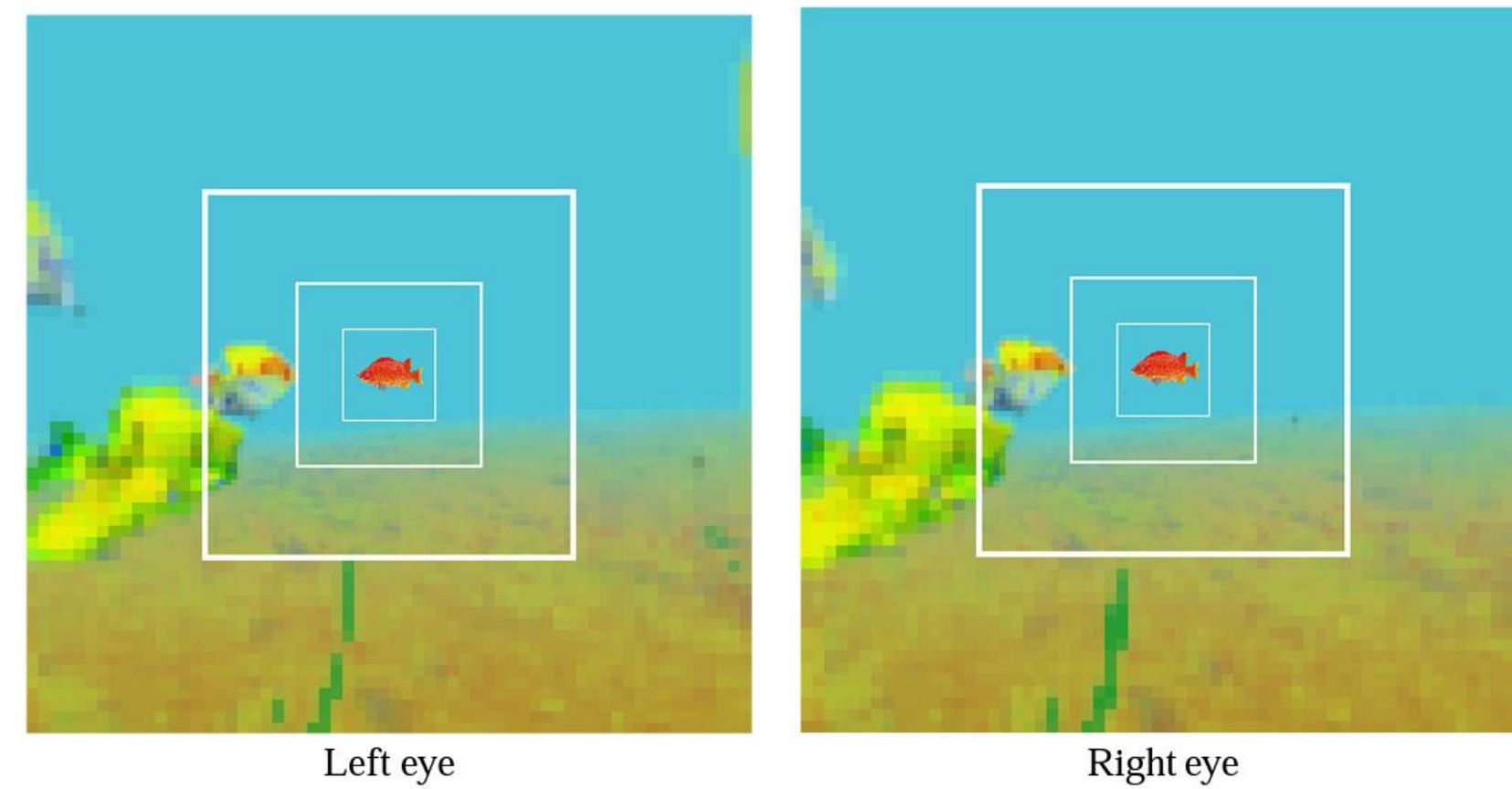
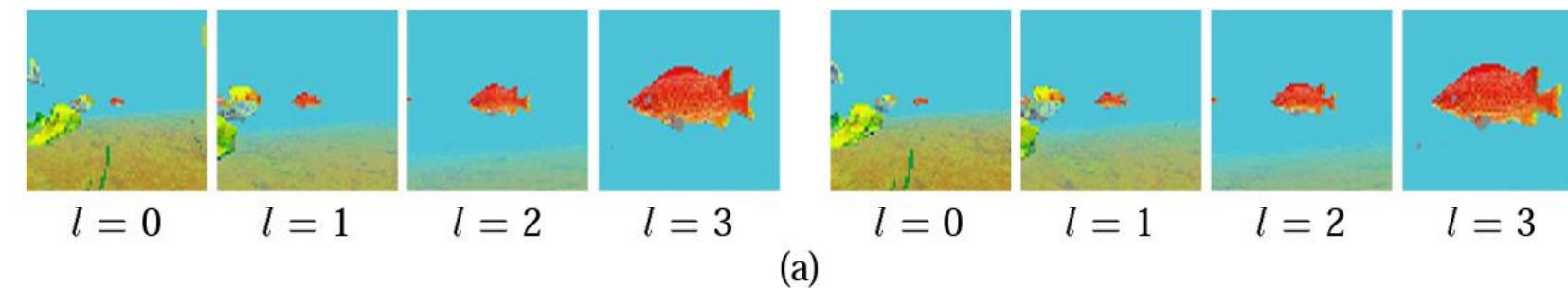
- |                |            |  |
|----------------|------------|--|
| → unmodifiable | excitatory | — heterosynaptic facilitation<br>(gain control of the target cell) |
| → modifiable   |            | ← weakening the inhibition<br>to the target cell                   |
| — unmodifiable | inhibitory |  |
| — modifiable   |            |  |



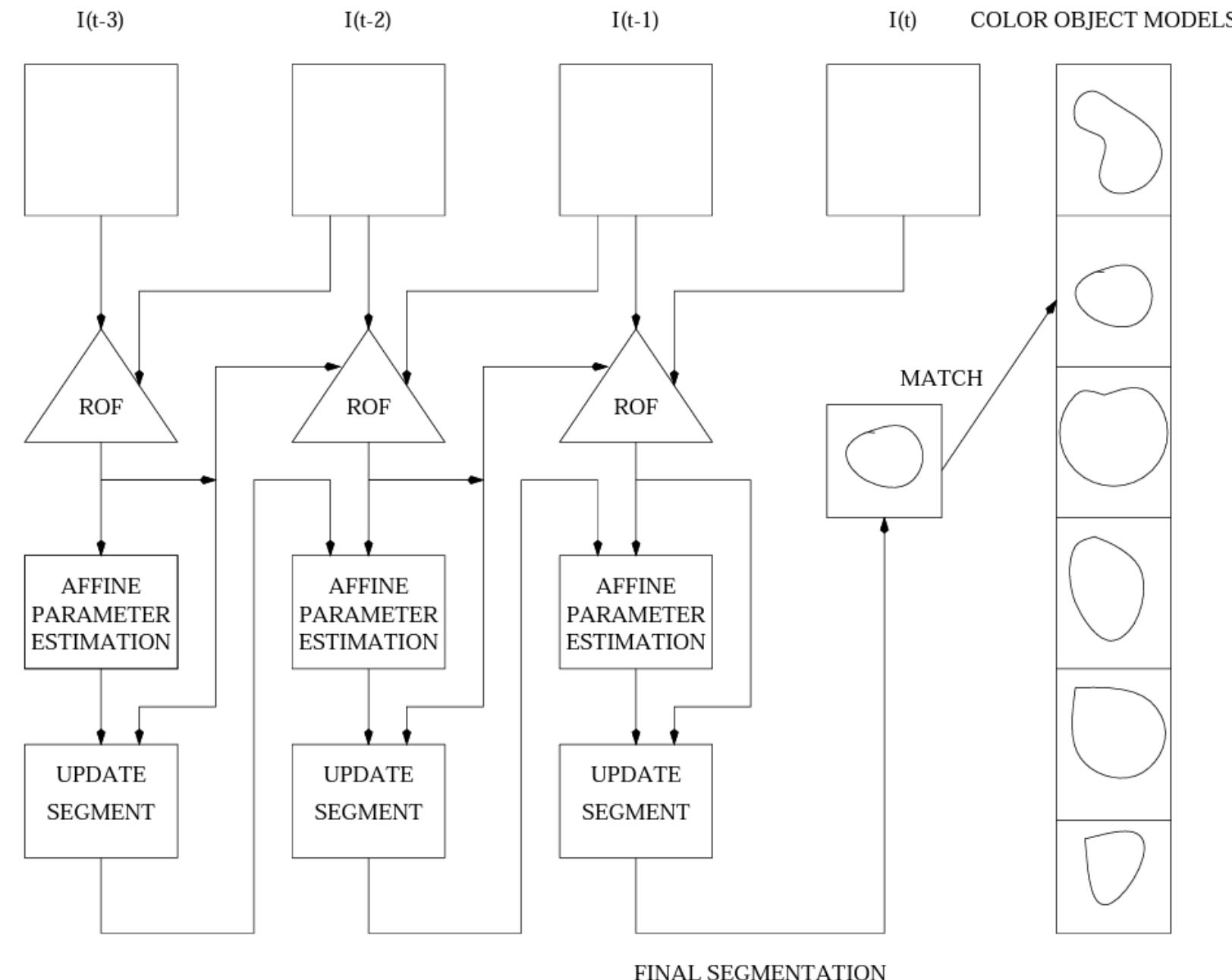
# Neocog<sup>n</sup>itron



# Overt Attention

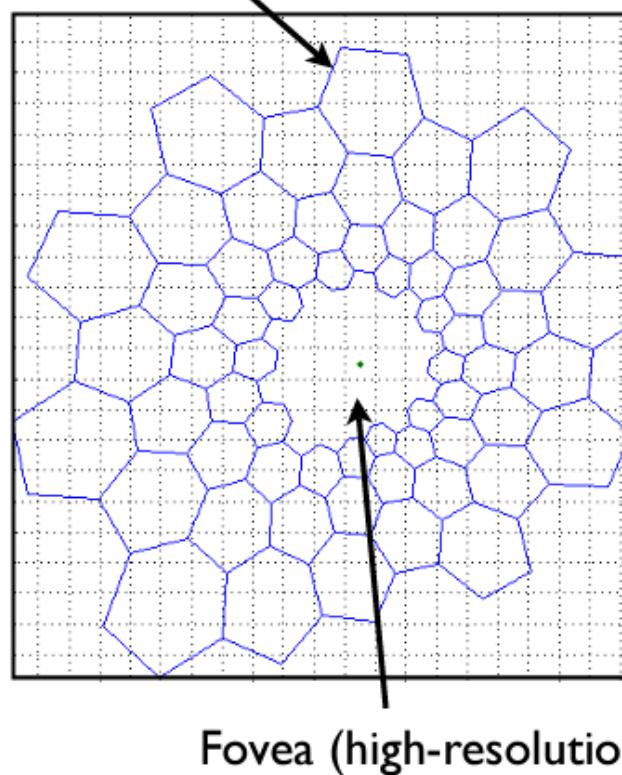


# Overt Attention

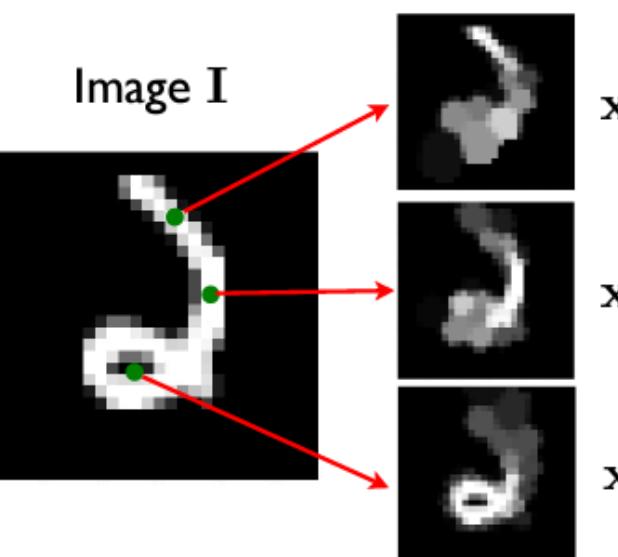


# Foveal Glimpse

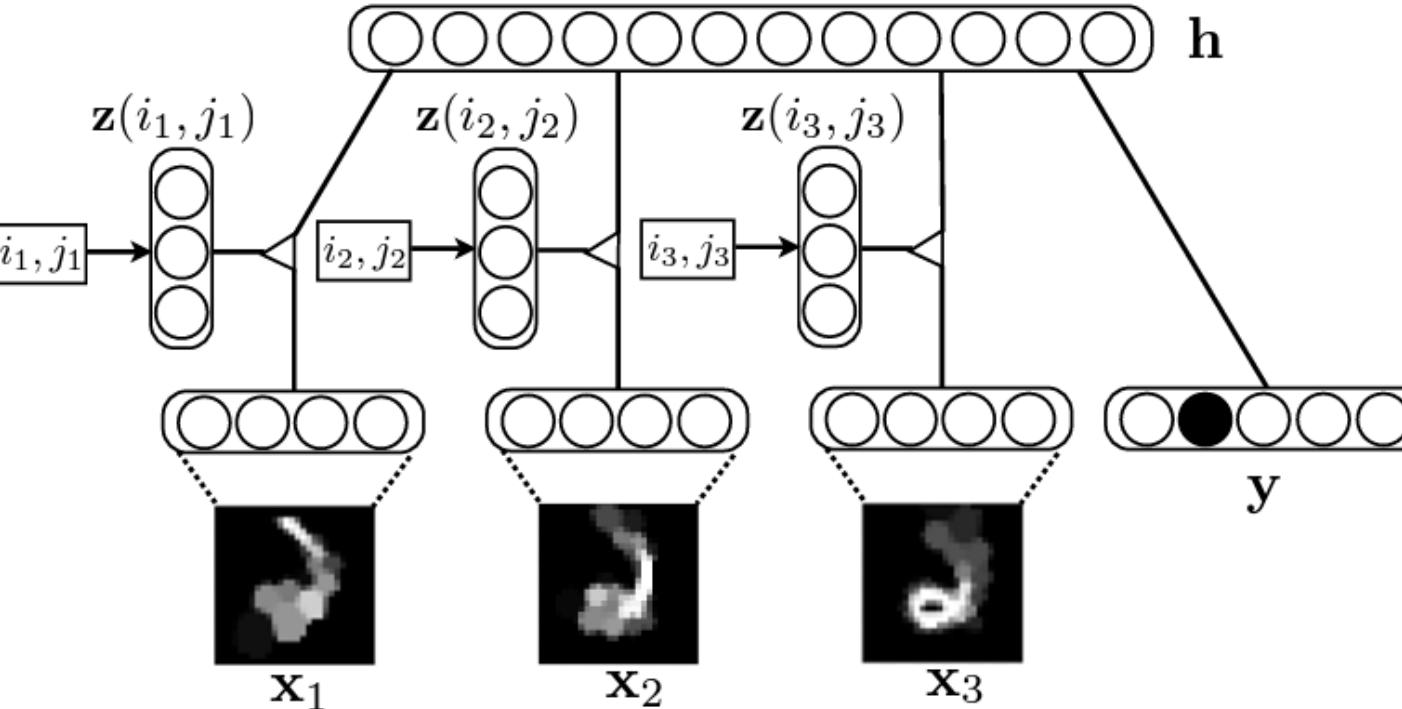
Periphery (low-resolution)



Retinal transformation  
(reconstruction from  $x_k$ )



A



C

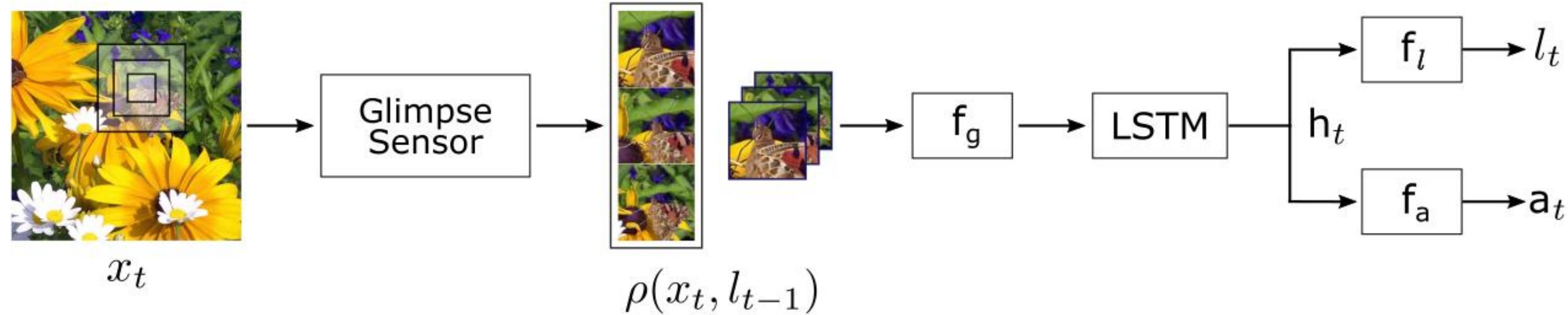
B



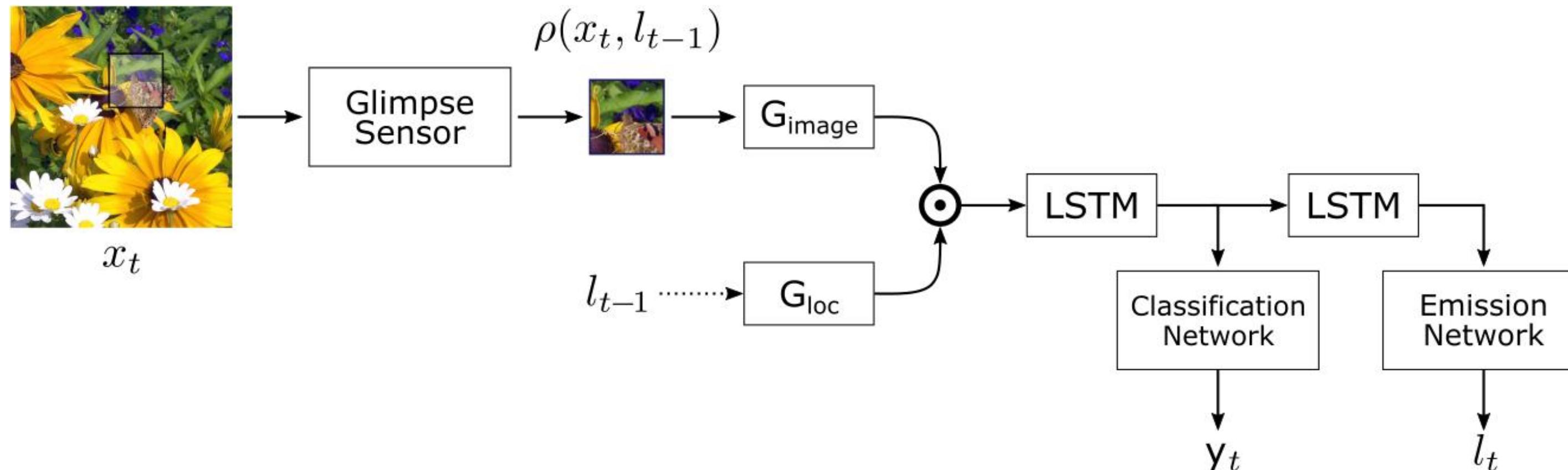
# Recurrent Attention Model (RAM)



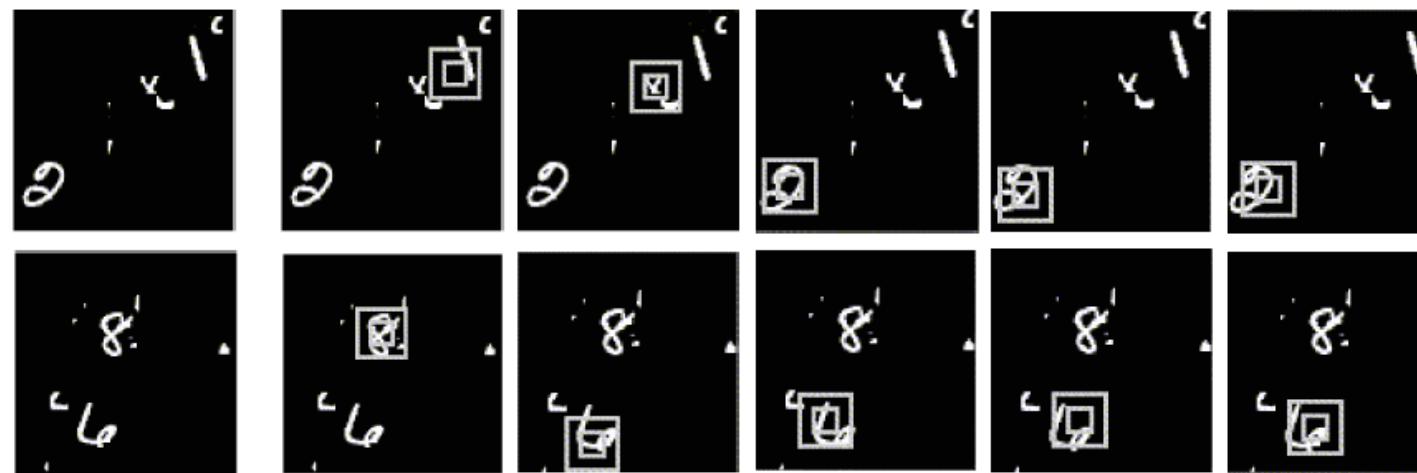
# Recurrent Attention Model (RAM)



# Deep Recurrent Visual Attention Model (DRAM)



# Deep Recurrent Visual Attention Model (DRAM)



Model	Test Err.
RAM Mnih et al. (2014)	9%
DRAM w/o context	7%
DRAM	<b>5%</b>

Model	Test Err.
ConvNet 64-64-64-512	3.2%
DRAM	<b>2.5%</b>



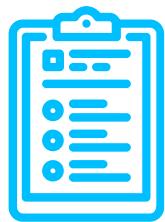
Model	Test Err.
11 layer CNN Goodfellow et al. (2013)	3.96%
10 layer CNN	4.11%
Single DRAM	5.1%
Single DRAM MC avg.	4.4%
forward-backward DRAM MC avg.	<b>3.9%</b>



**TRANSFOR**  
**MA**  
**TEC**

Ba et al. (2015) "Multiple object recognition with visual attention".  
arXiv preprint arXiv:1412.7755.

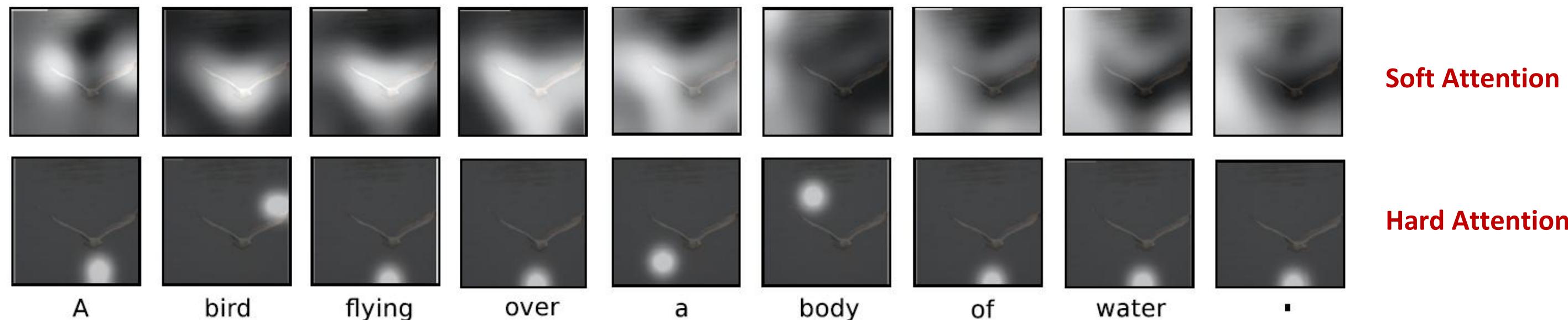
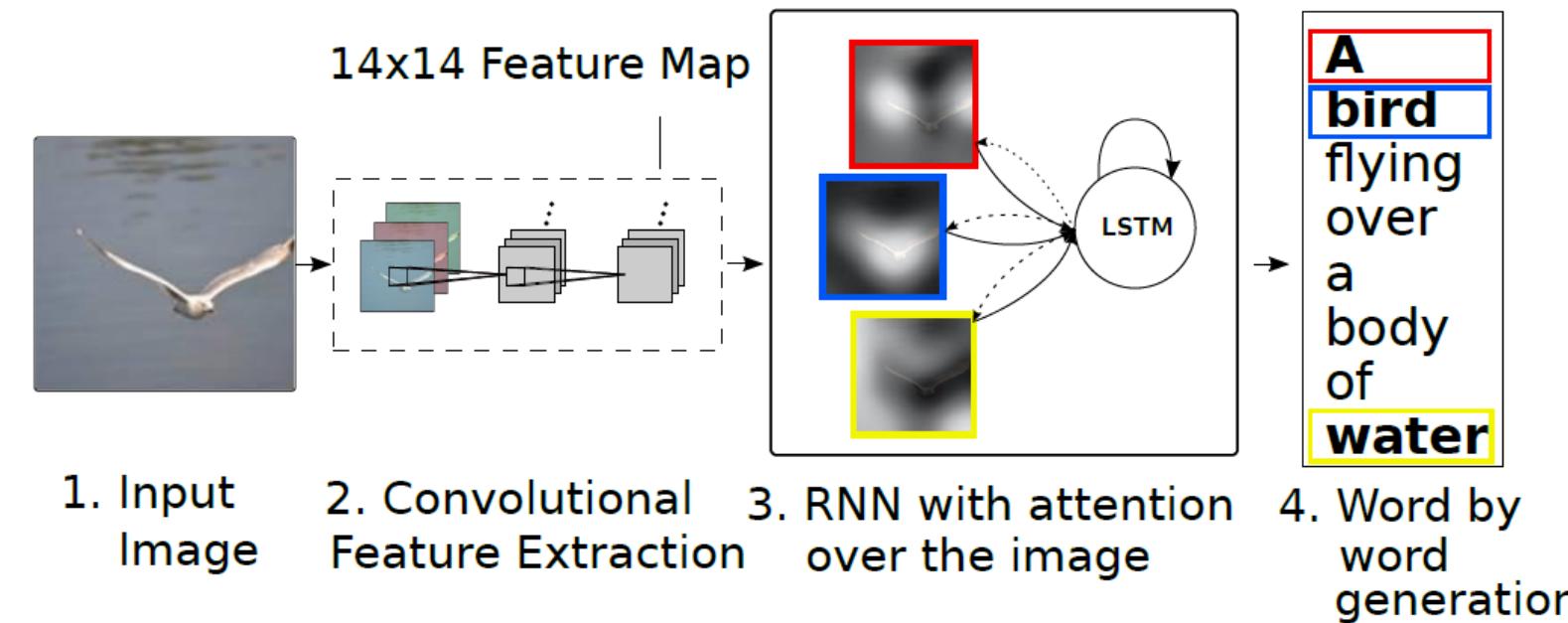
4.



## **Soft** attention

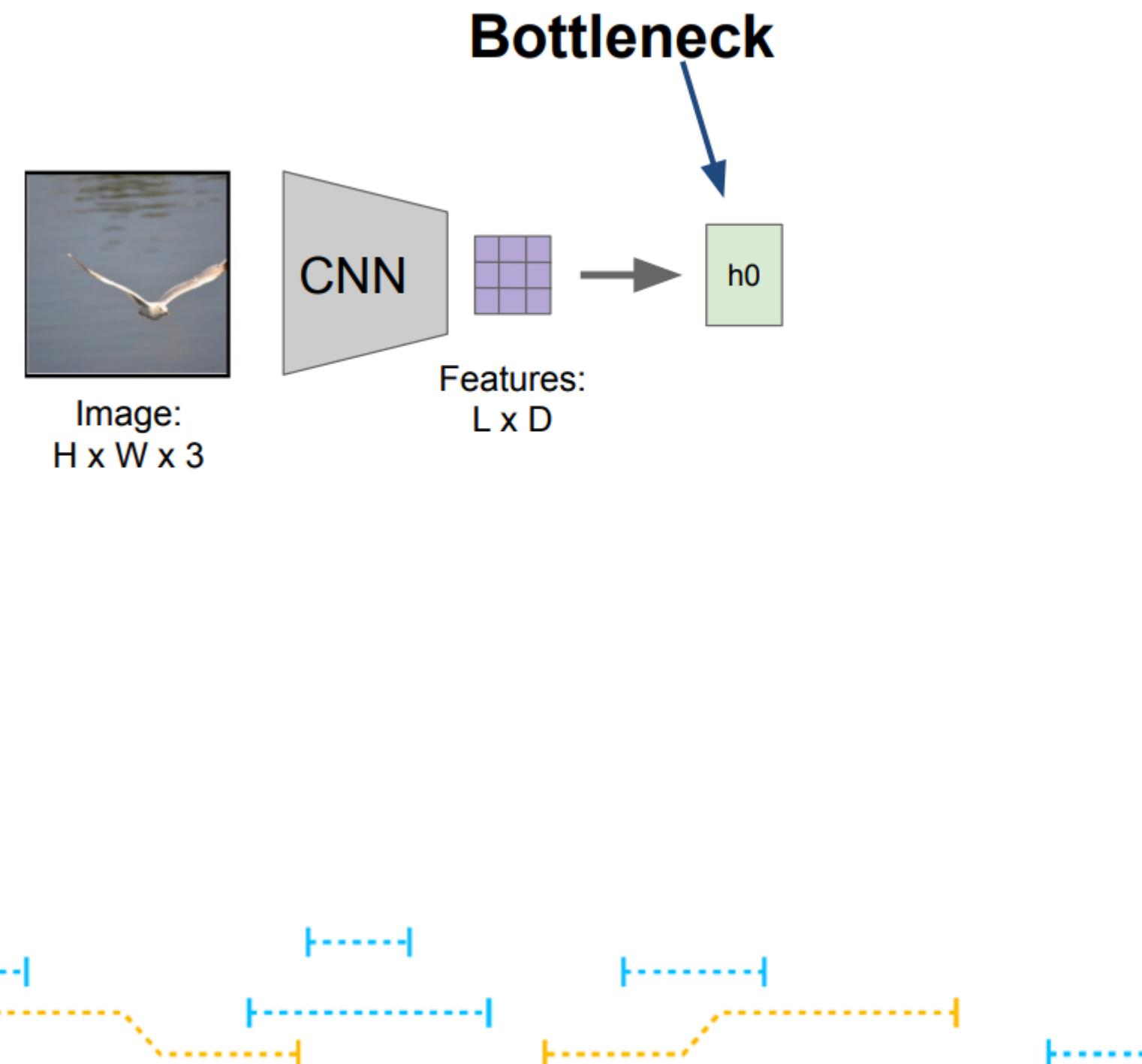


# Show, Attend and Tell

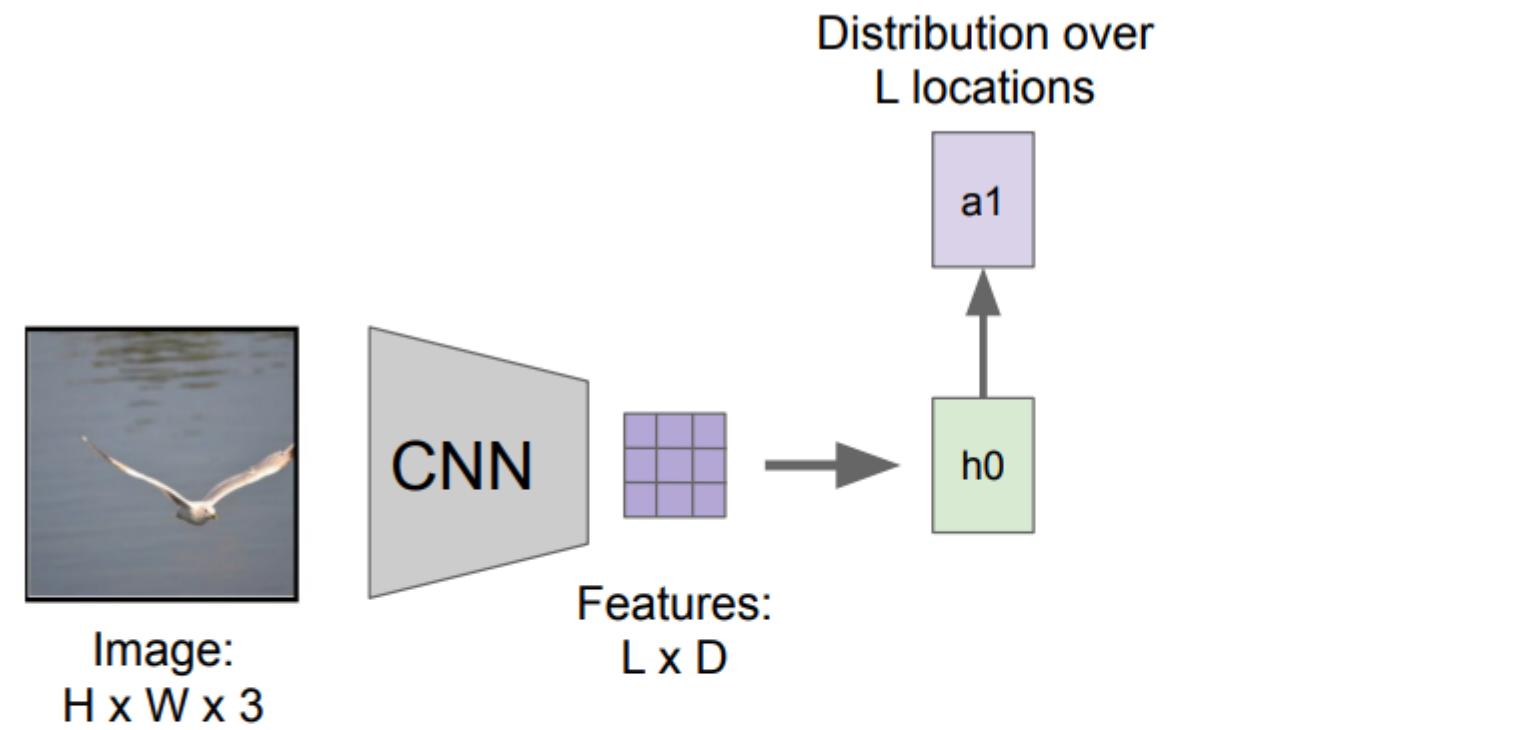


**TRANSFORMATEC**

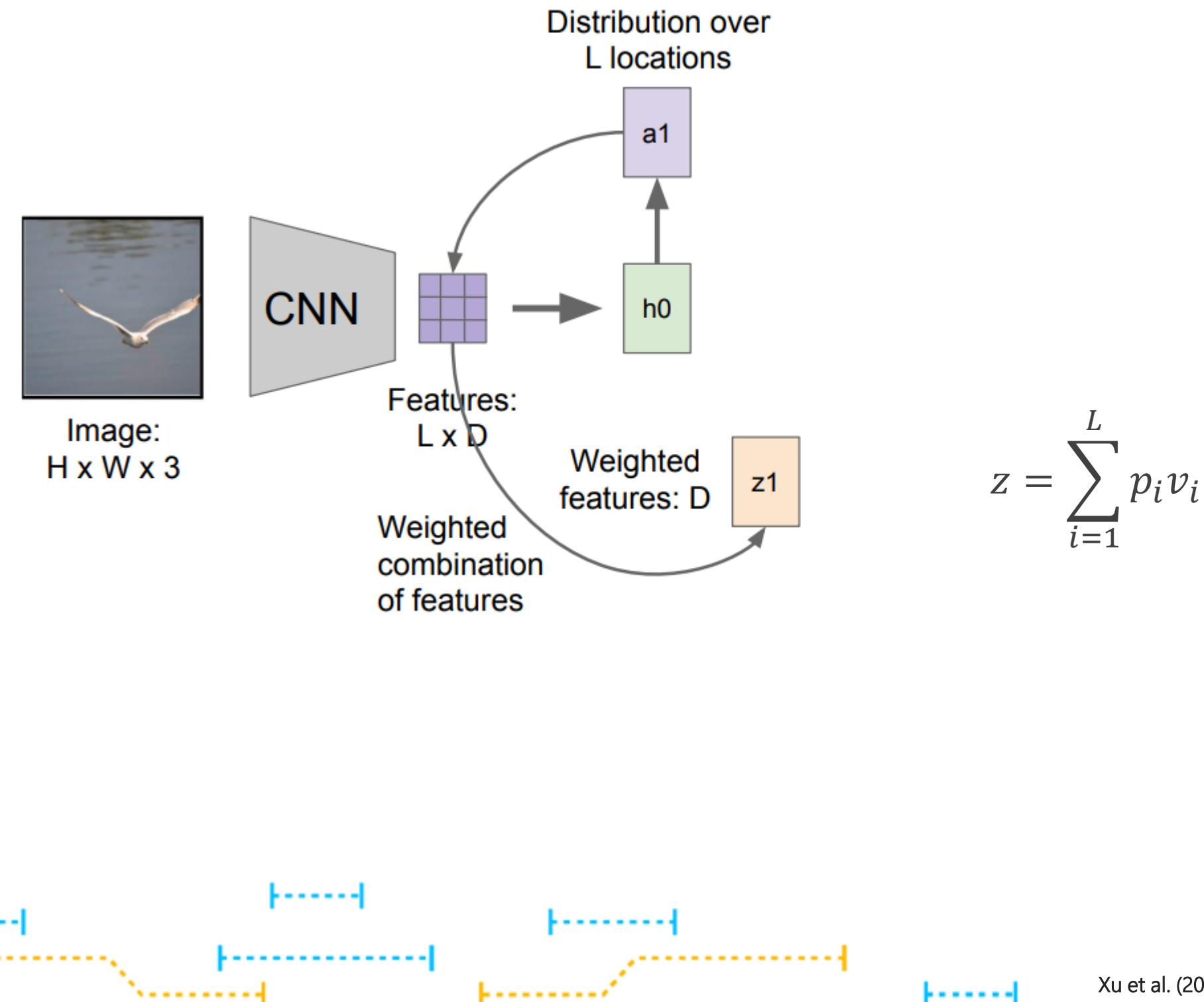
# Show, Attend *and Tell*



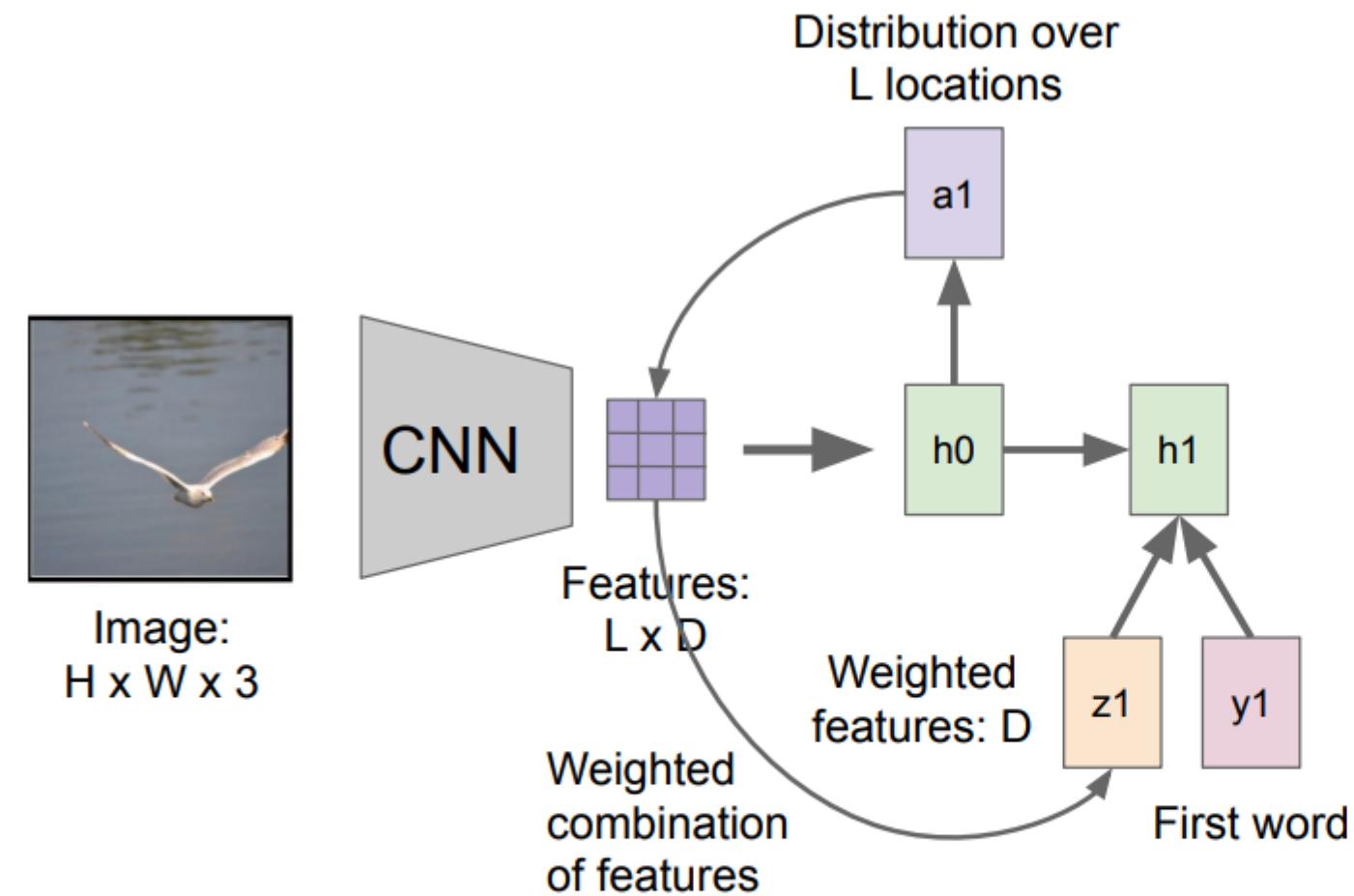
# Show, Attend and Tell



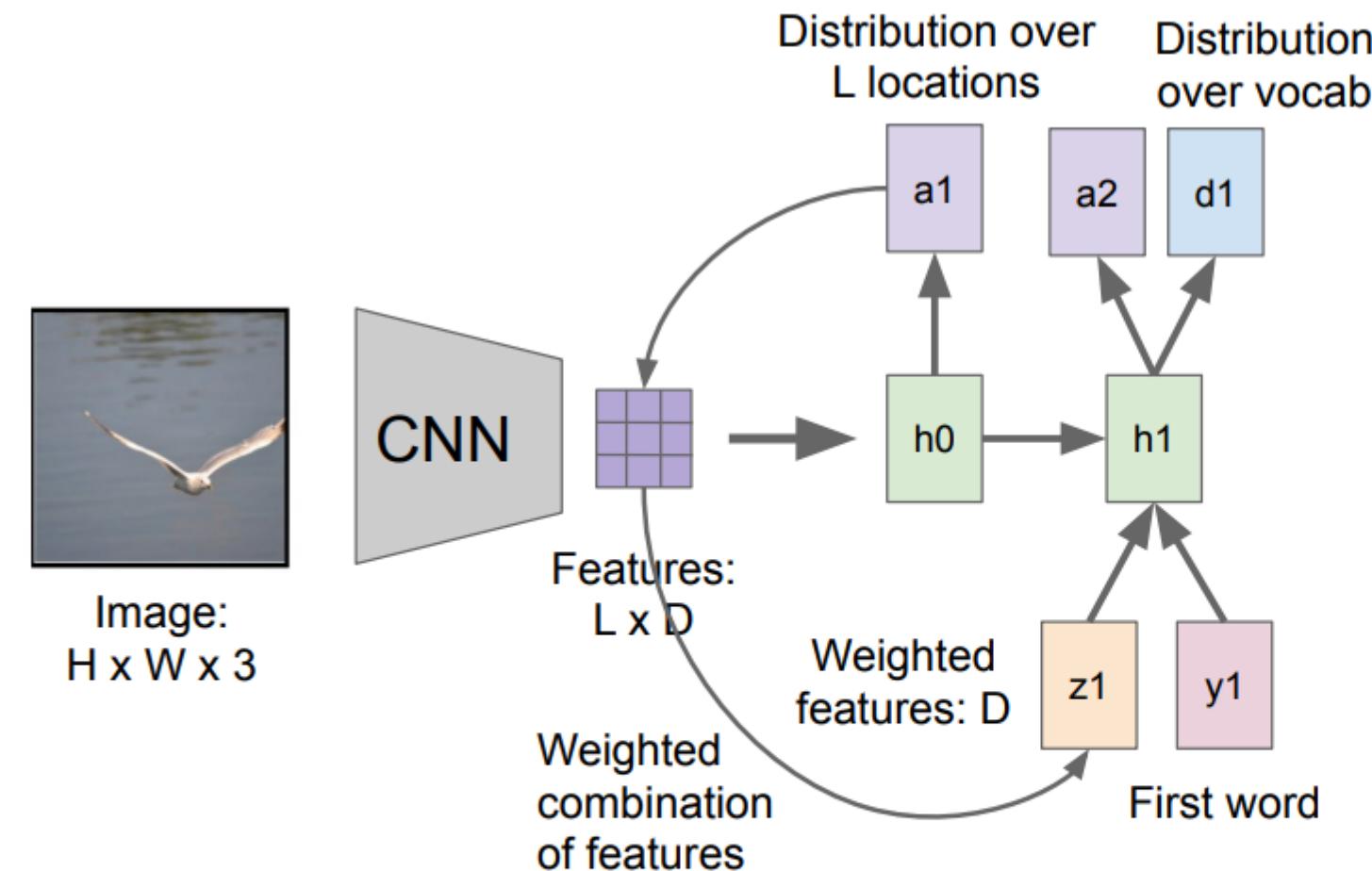
# Show, Attend and Tell



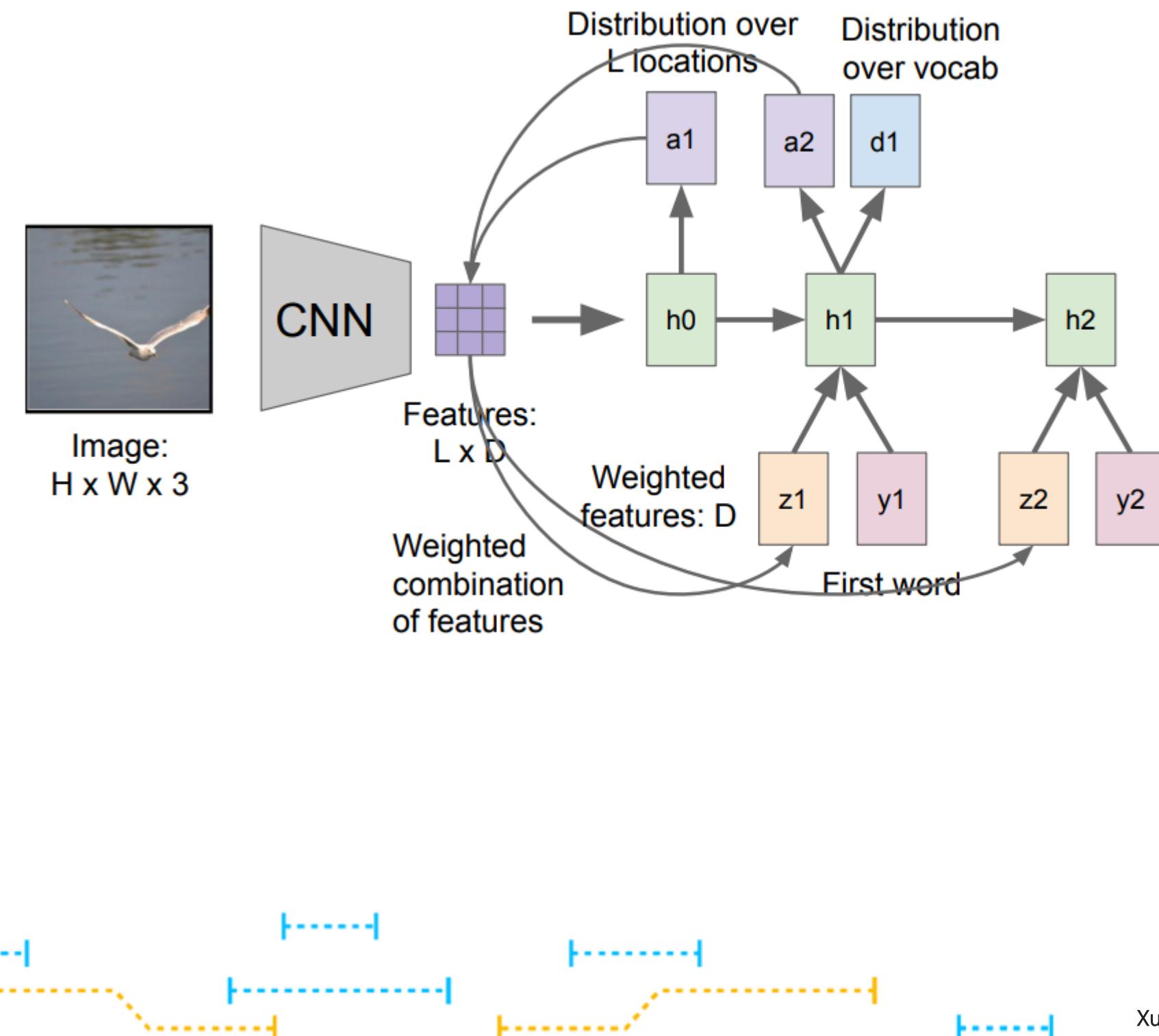
# Show, Attend and Tell



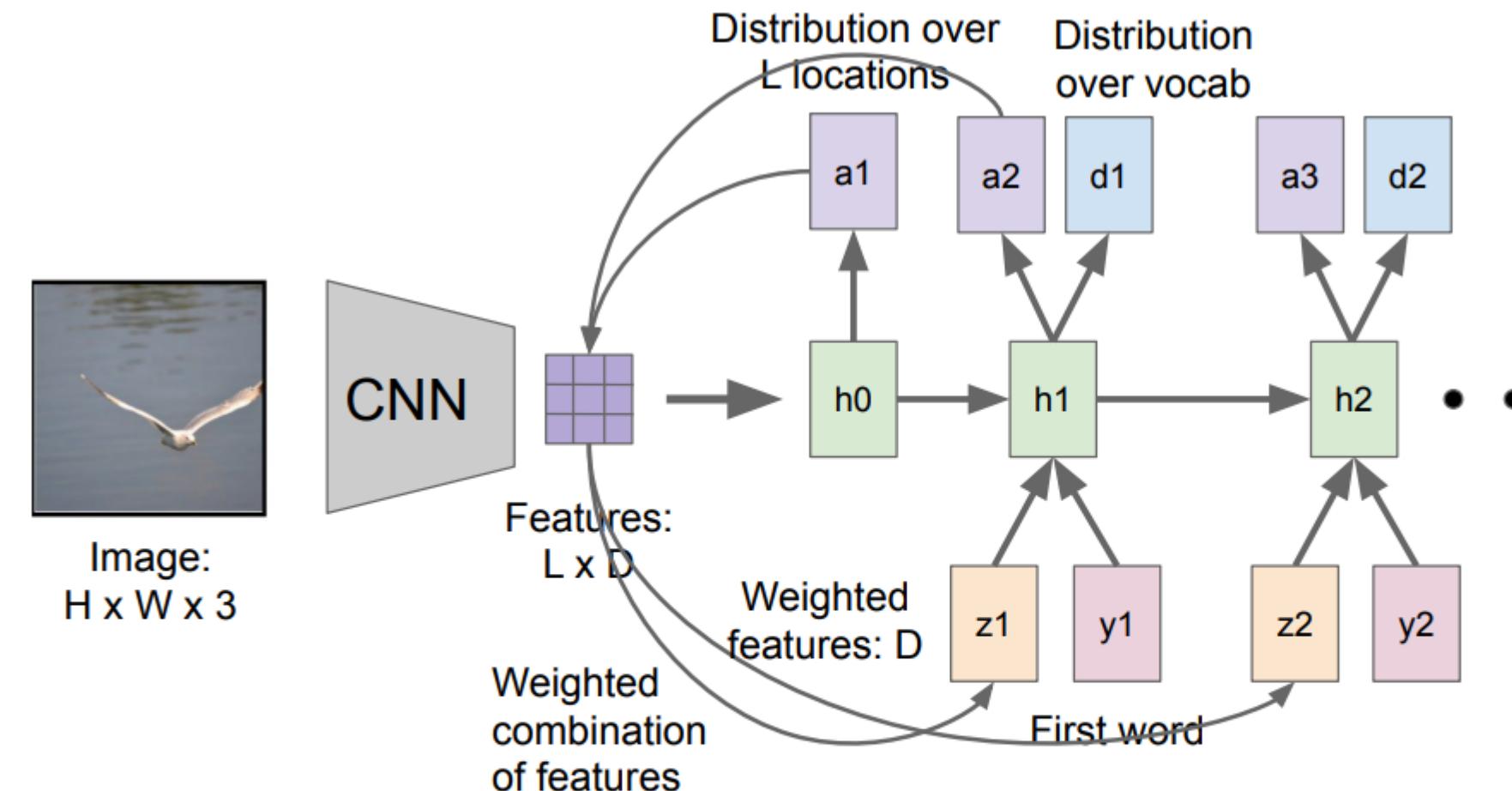
# Show, Attend and Tell



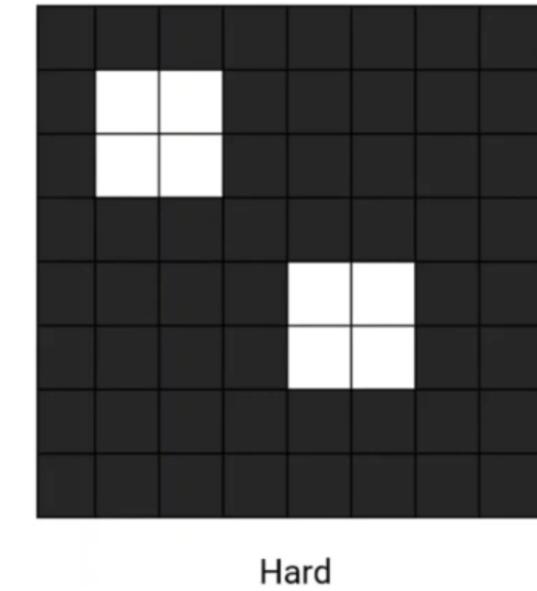
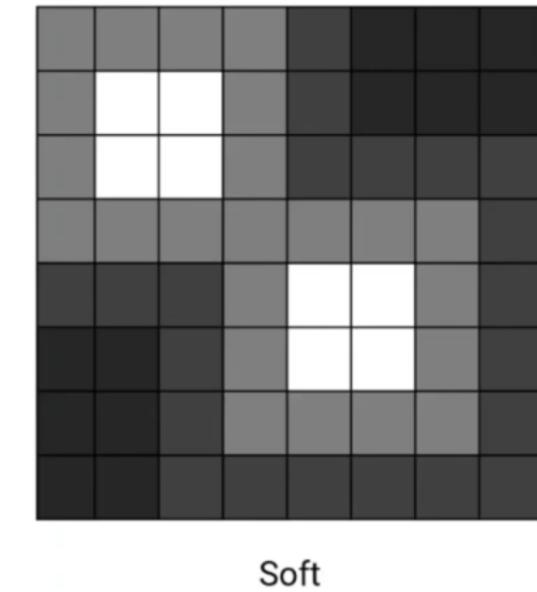
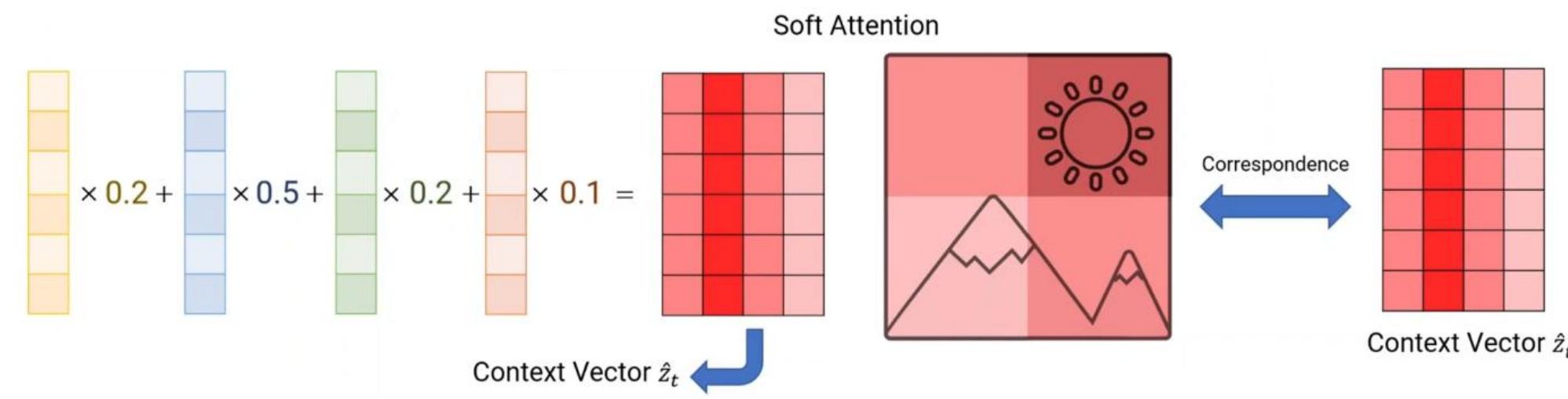
# Show, Attend and Tell



# Show, Attend and Tell



# Show, Attend and Tell



# Show, Attend and Tell

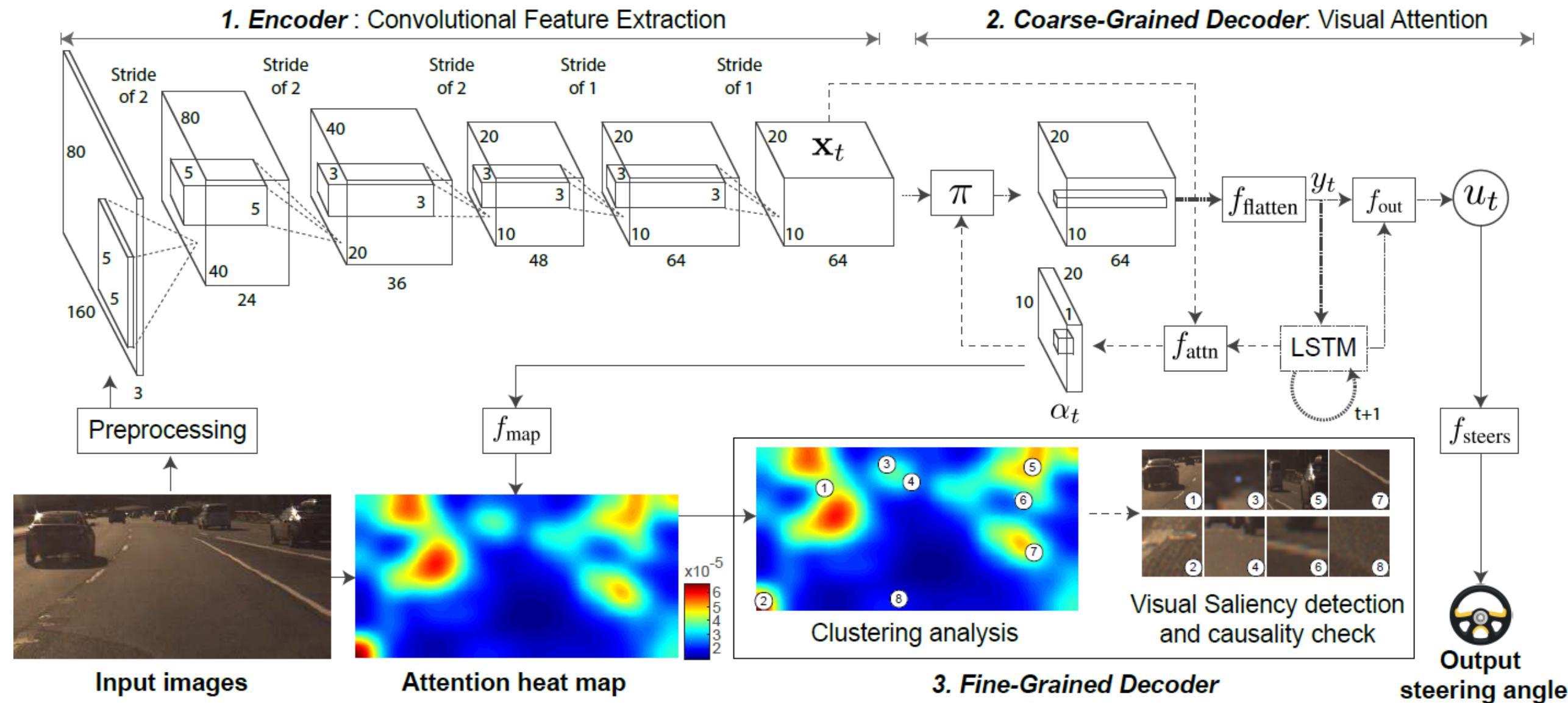


(a) A man and a woman playing frisbee in a field.

(b) A woman is throwing a frisbee in a park.

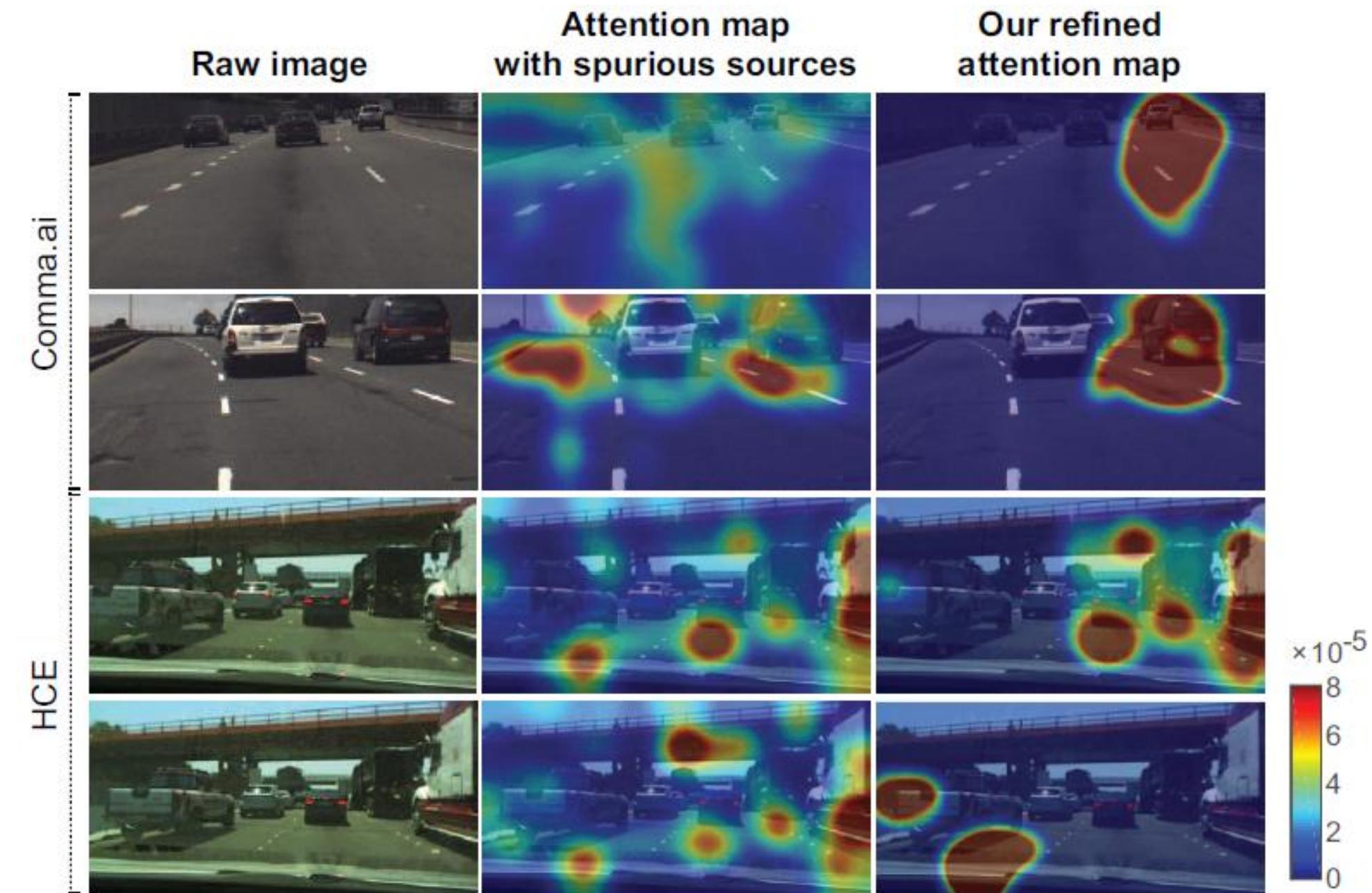


# Self-driving car



**TRANSFORMATEC**

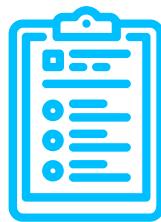
# Self-driving car



**TRANSFORMATEC**

Kim et al. (2017) "Interpretable learning for self-driving cars by visualizing causal attention".  
Proceedings of the IEEE international conference on computer vision (pp. 2942-2950).

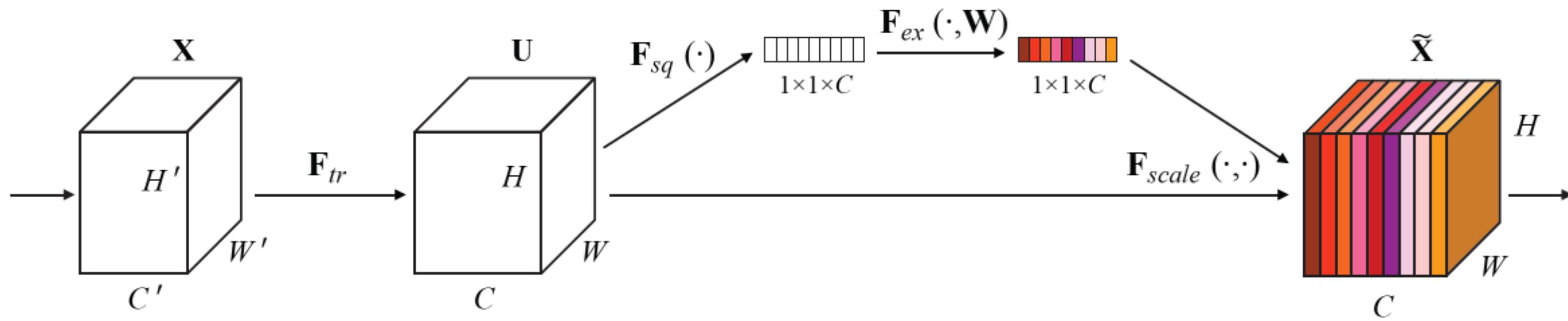
**5.**



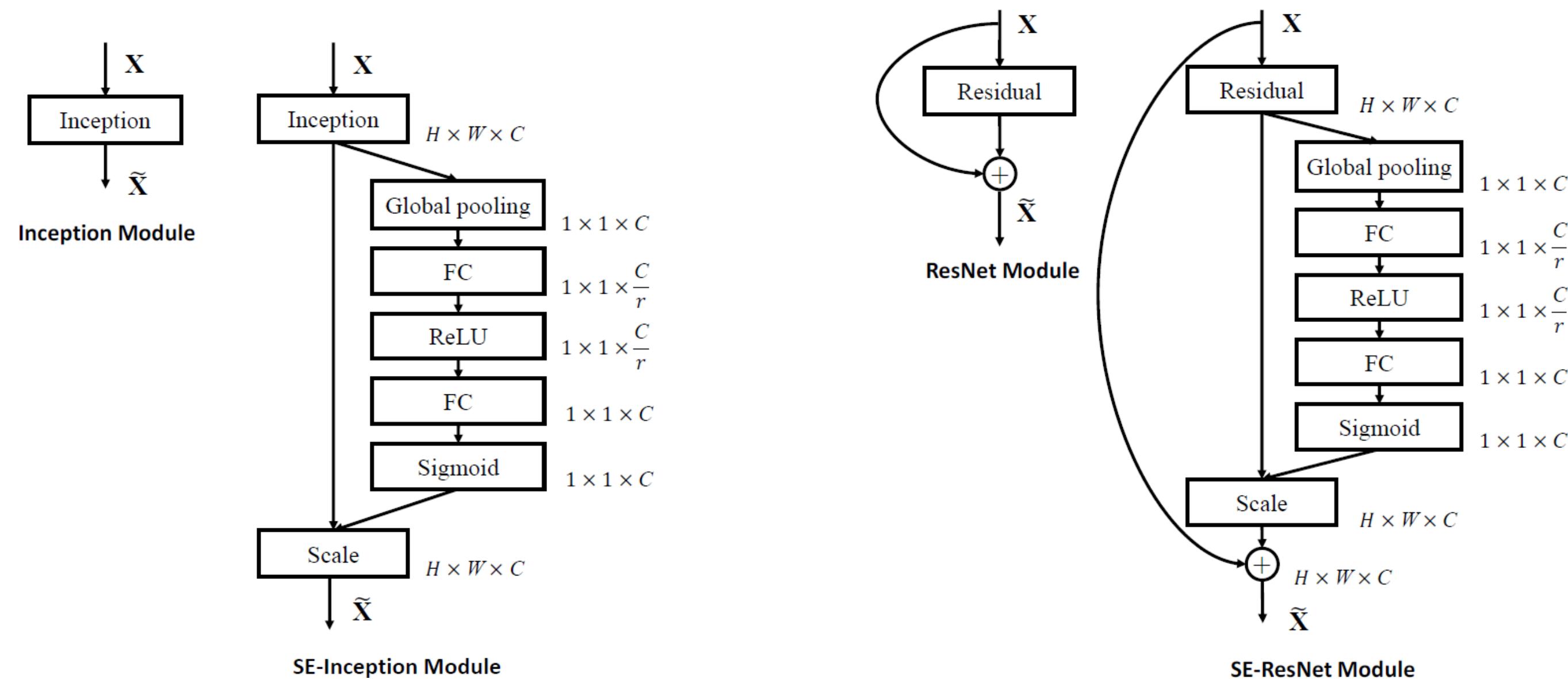
## **channel** *attention*



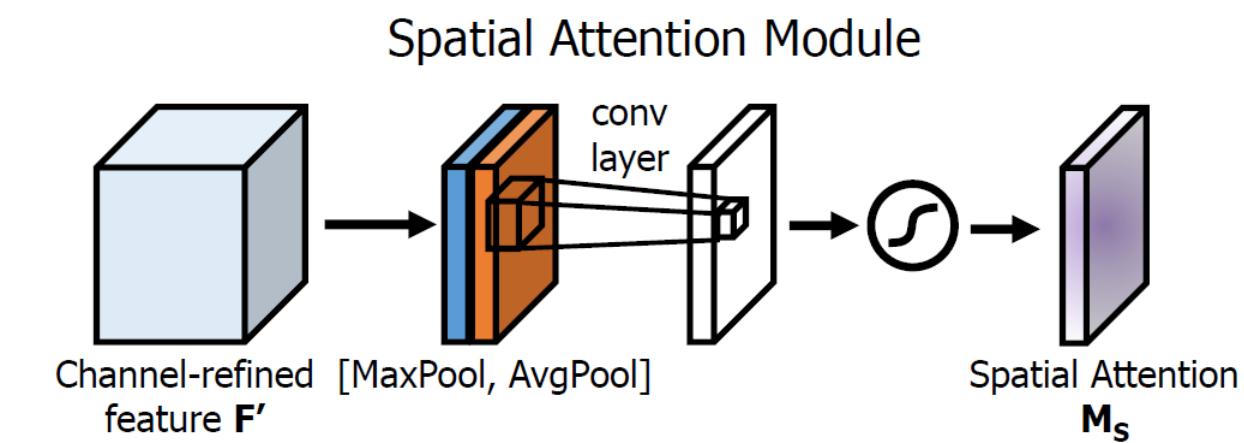
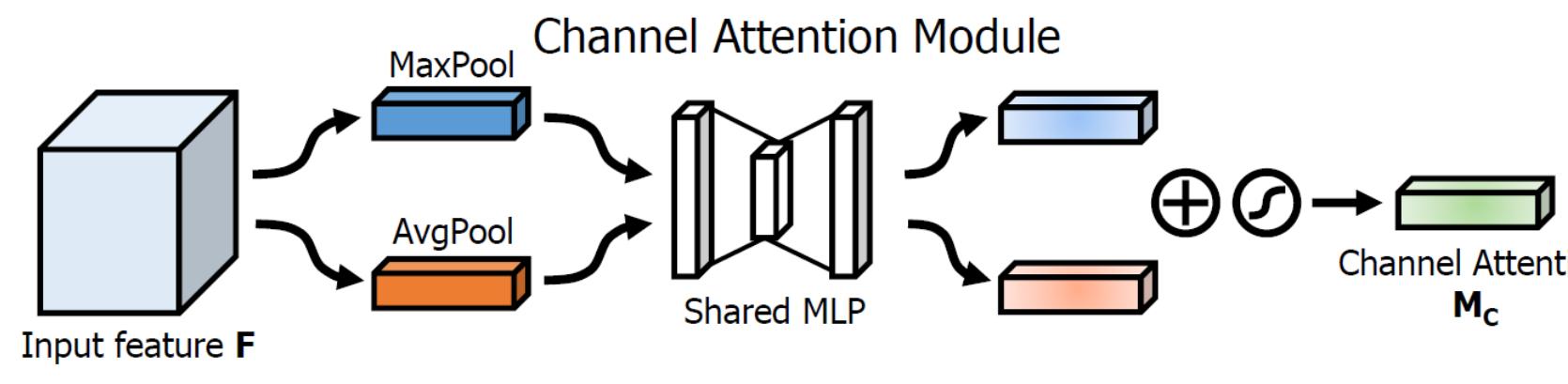
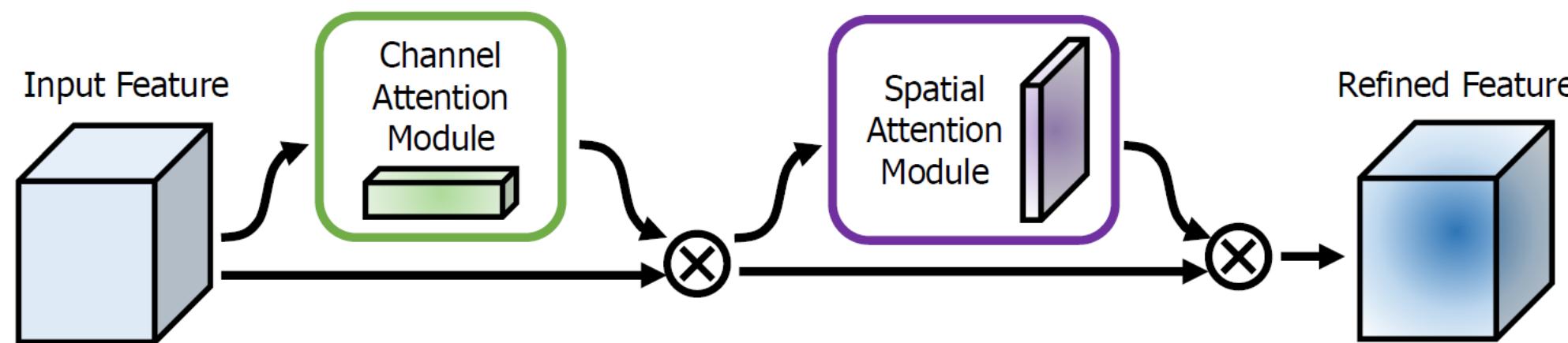
# Squeeze-and-Excitation



# Squeeze-and-Excitation



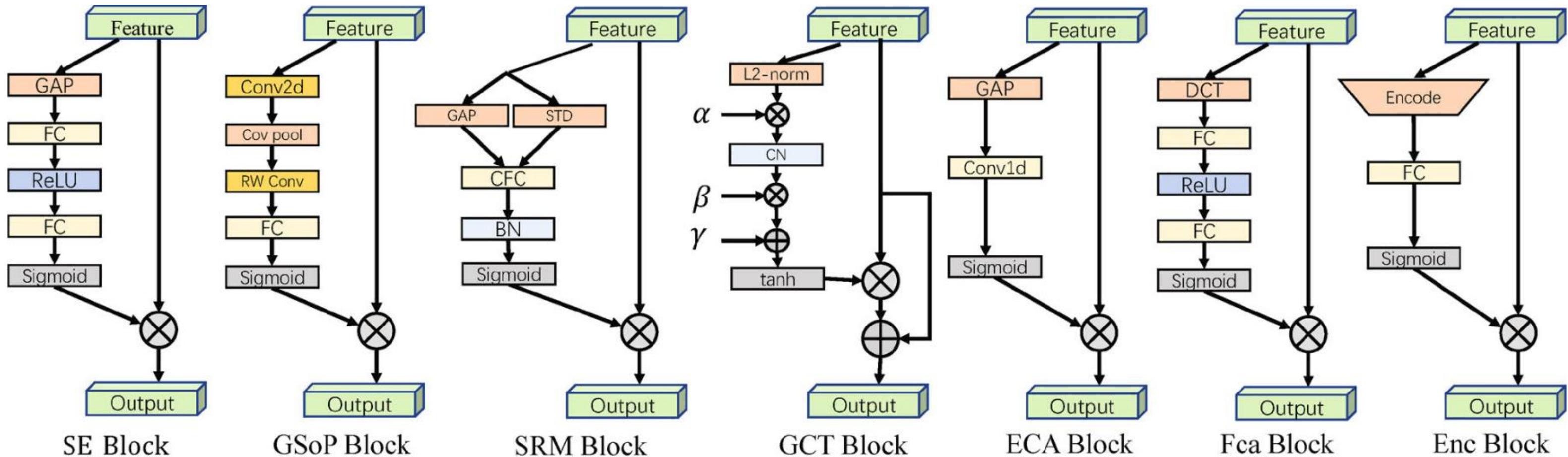
# Convolutional Block Attention Module



**TRANSFORMATEC**

Sanghyun Woo et al. (2018) "CBAM: Convolutional Block Attention Module".  
Proceedings of the European conference on computer vision (ECCV). 2018. p. 3-19.

# Channel *attention*



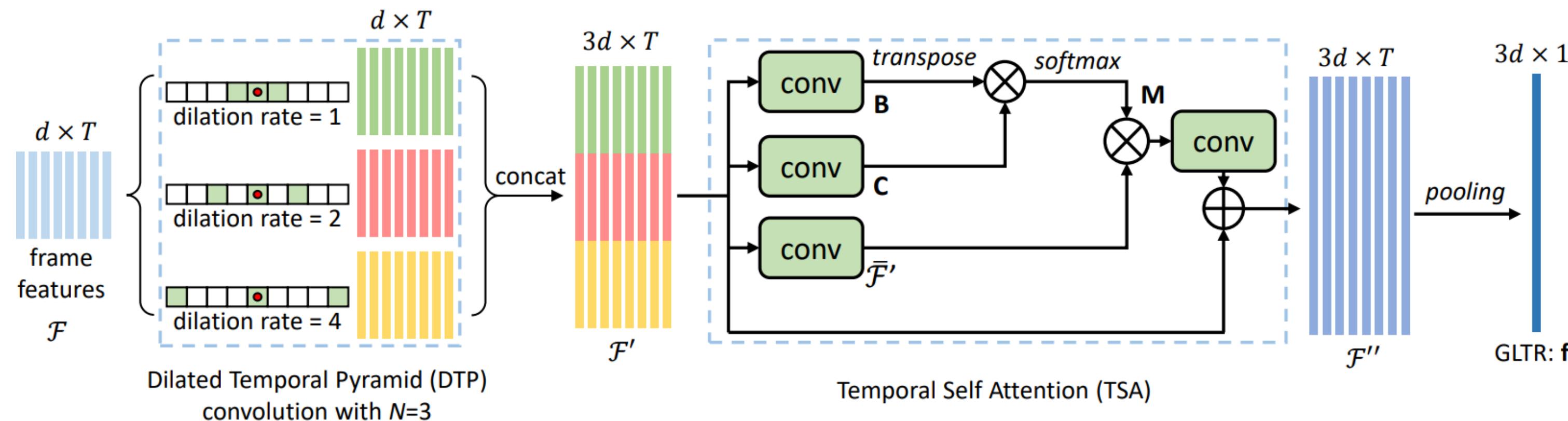
# 6.



## Temporal *attention*

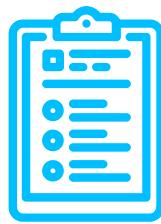


# Dilated Temporal Pyramid



**TRANSFORMATEC**

**7.**



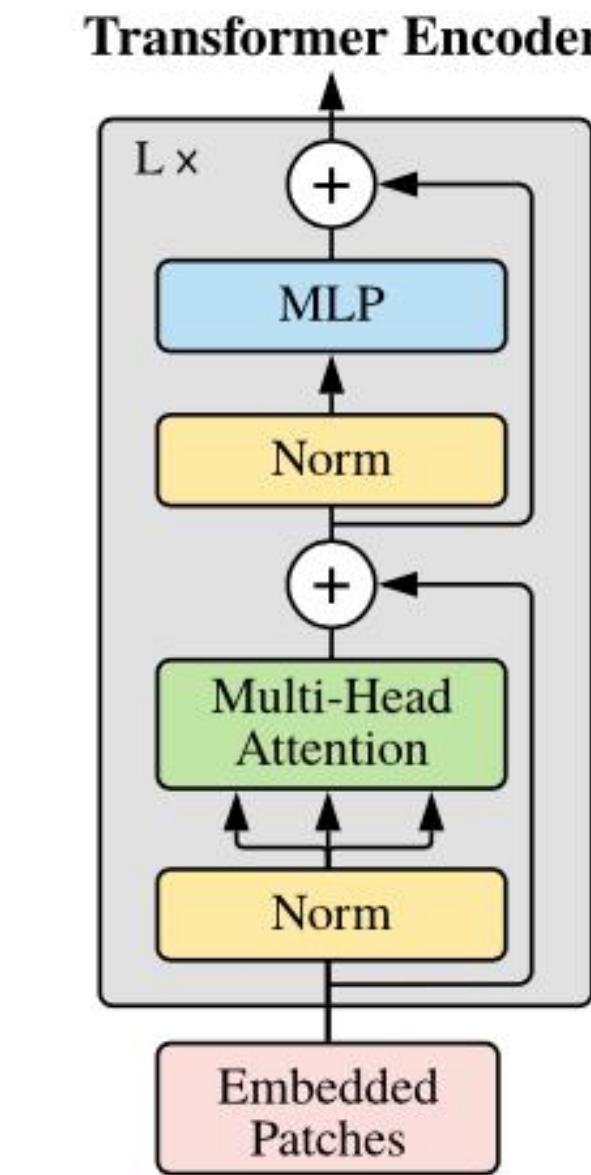
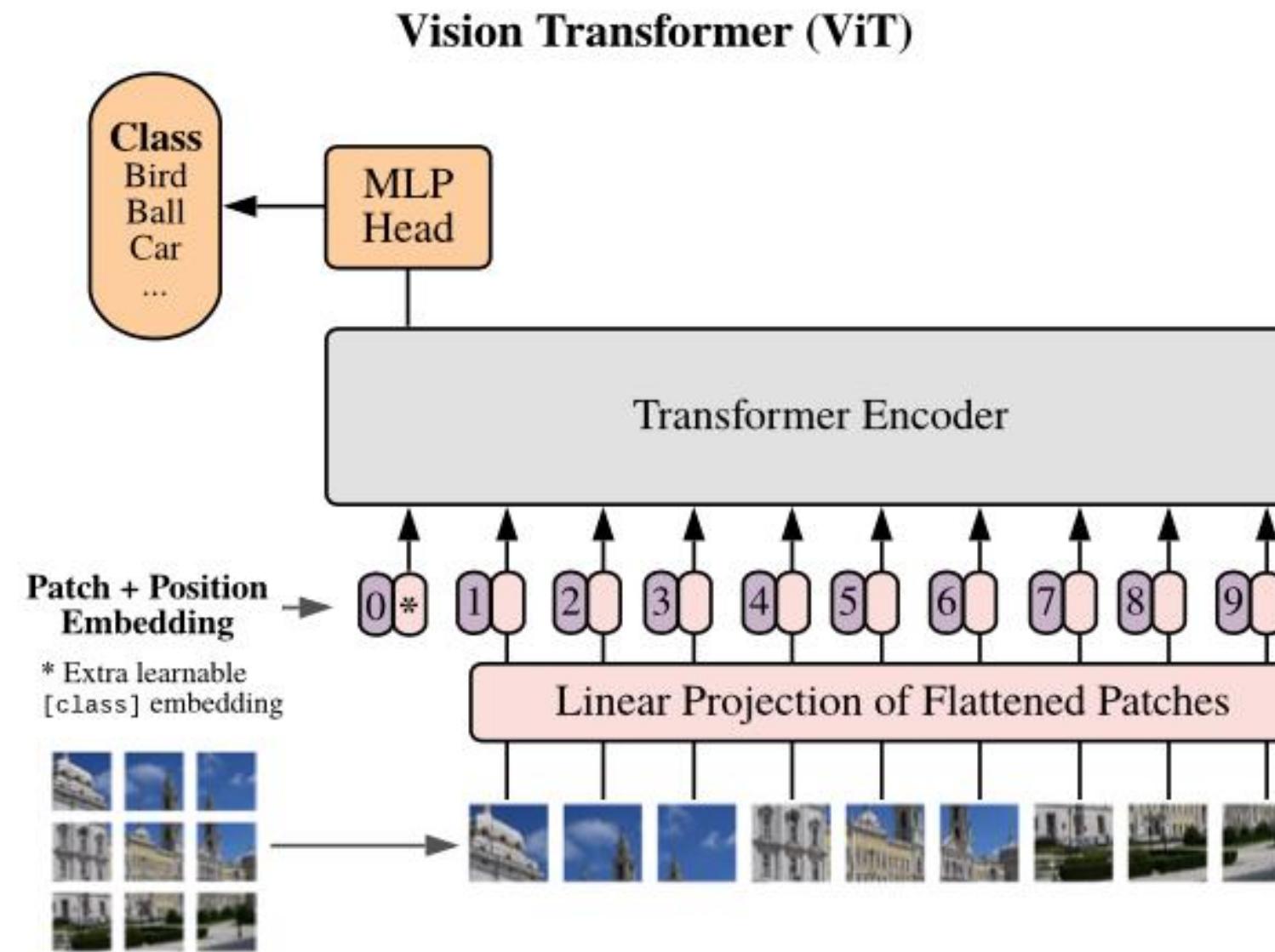
## **Visual Transformer**

**TRANSFORMATEC**

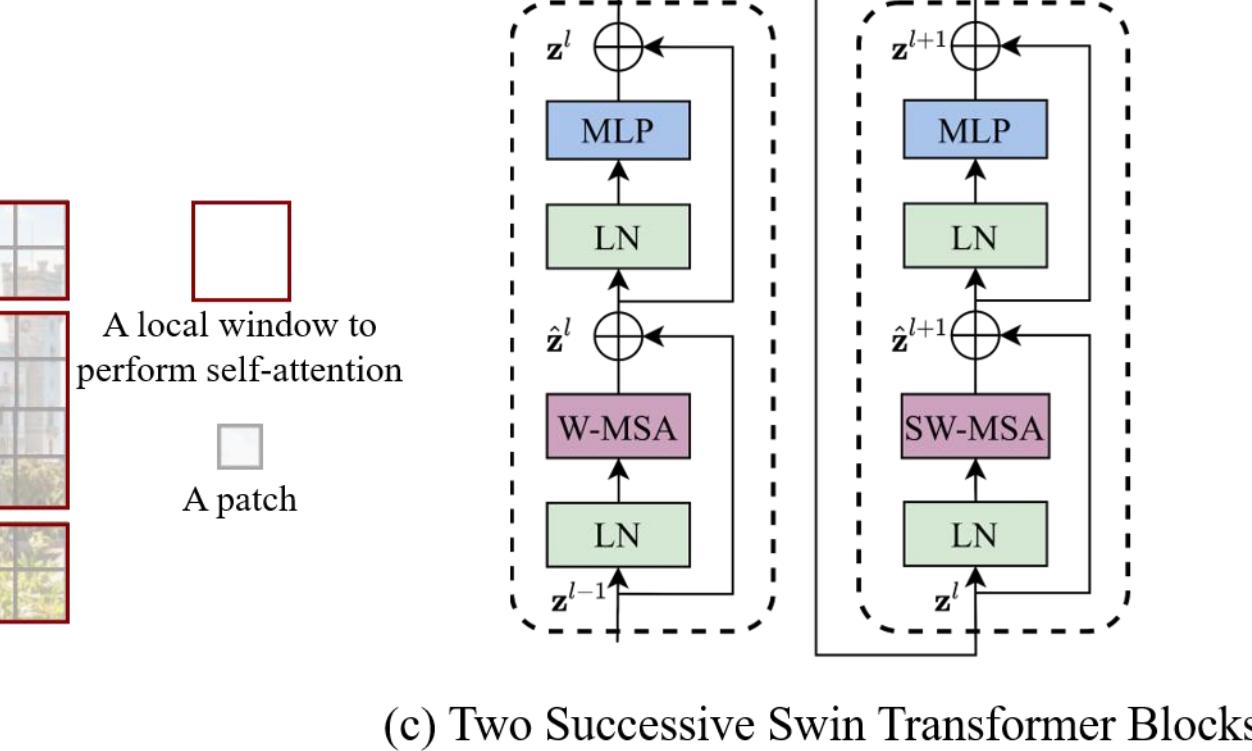
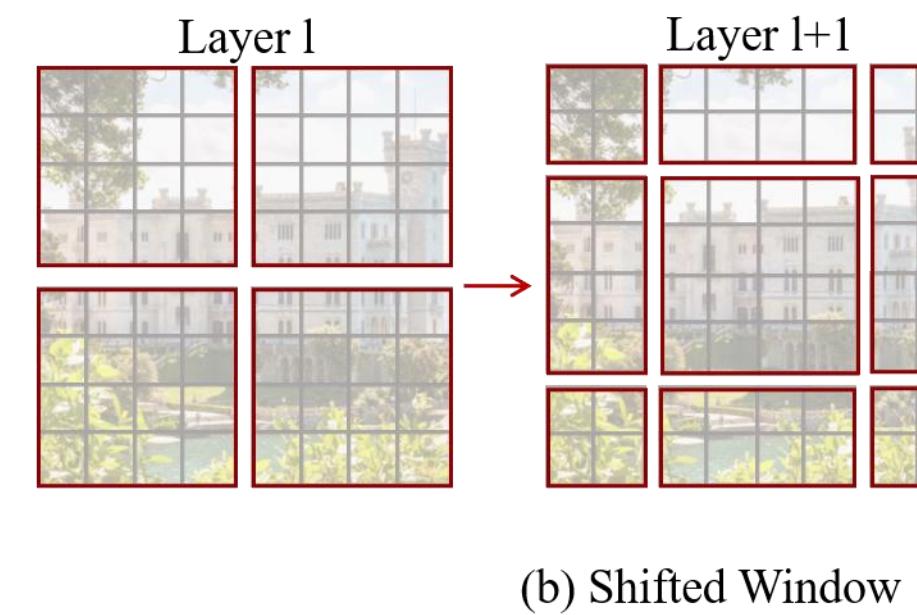
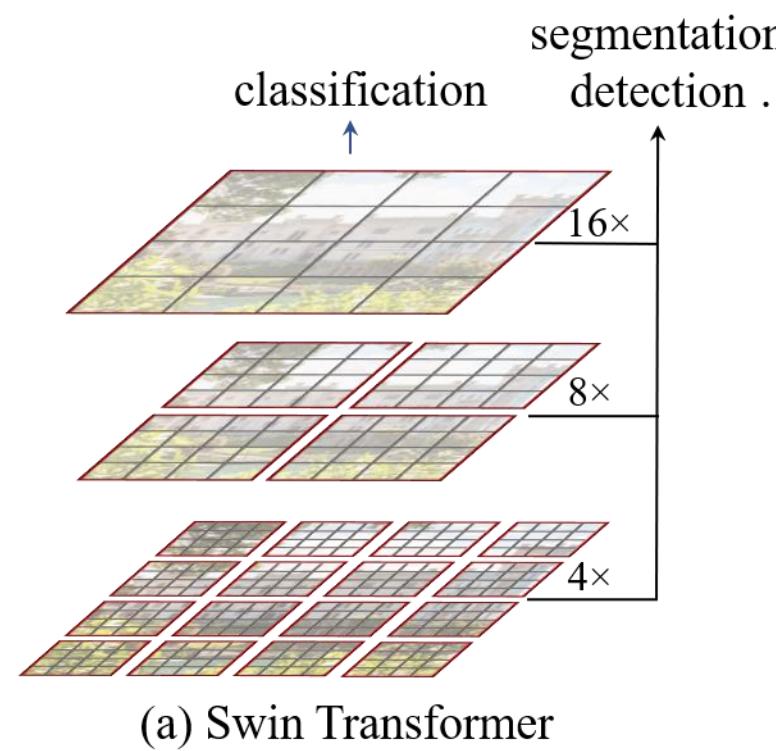
> Reinventa el mundo <



# ViT

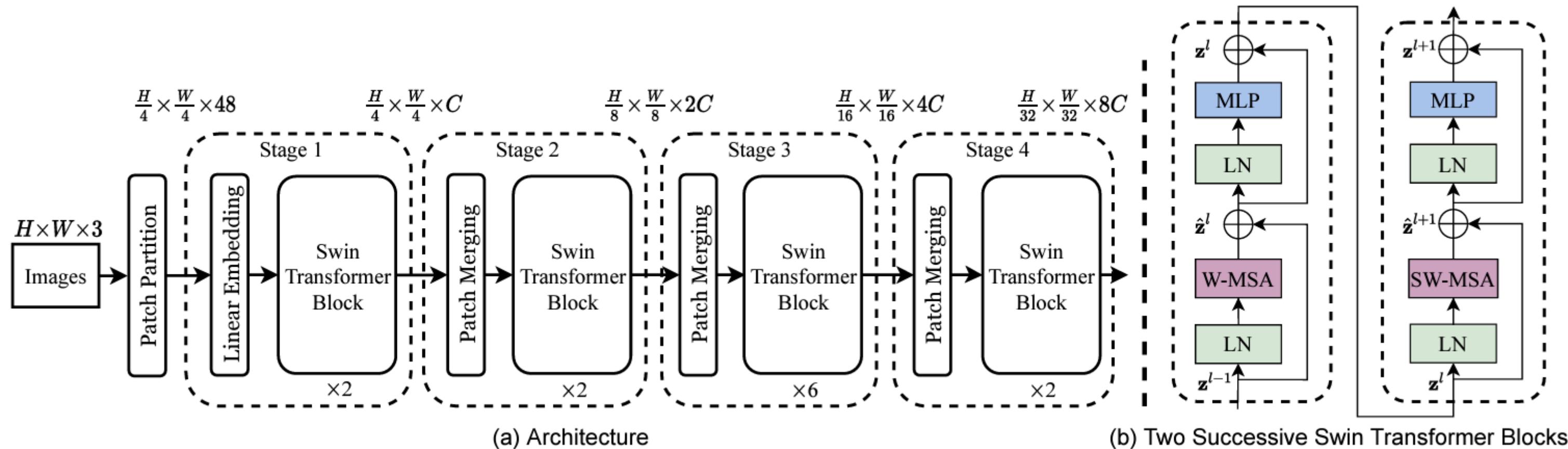


# Swin Transformer



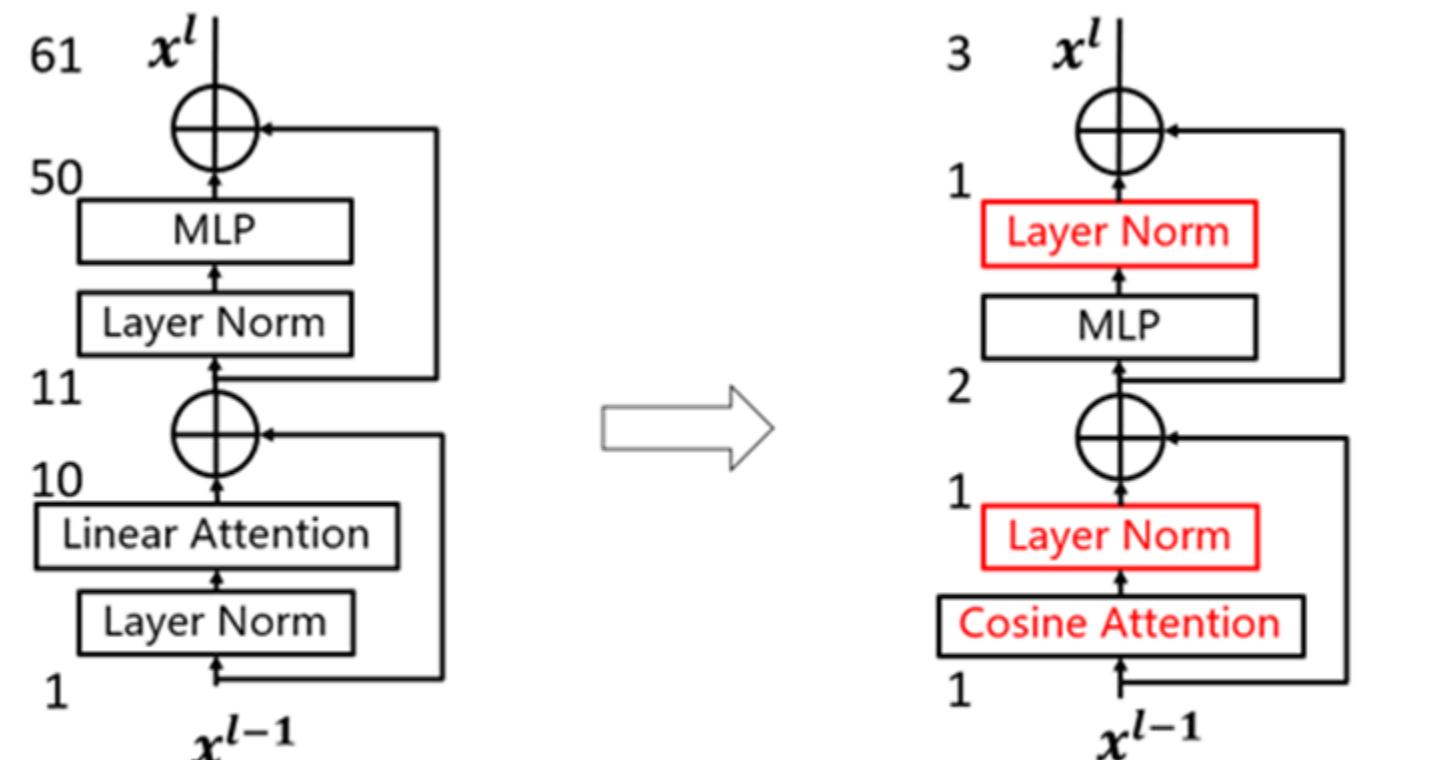
**TRANSFORMATEC**

# Swin Transformer



**TRANSFORMATEC**

# Swin Transformer v2



Swin V1

(pre-norm + linear attention)

Swin V2

(res-post-norm + cosine attention)

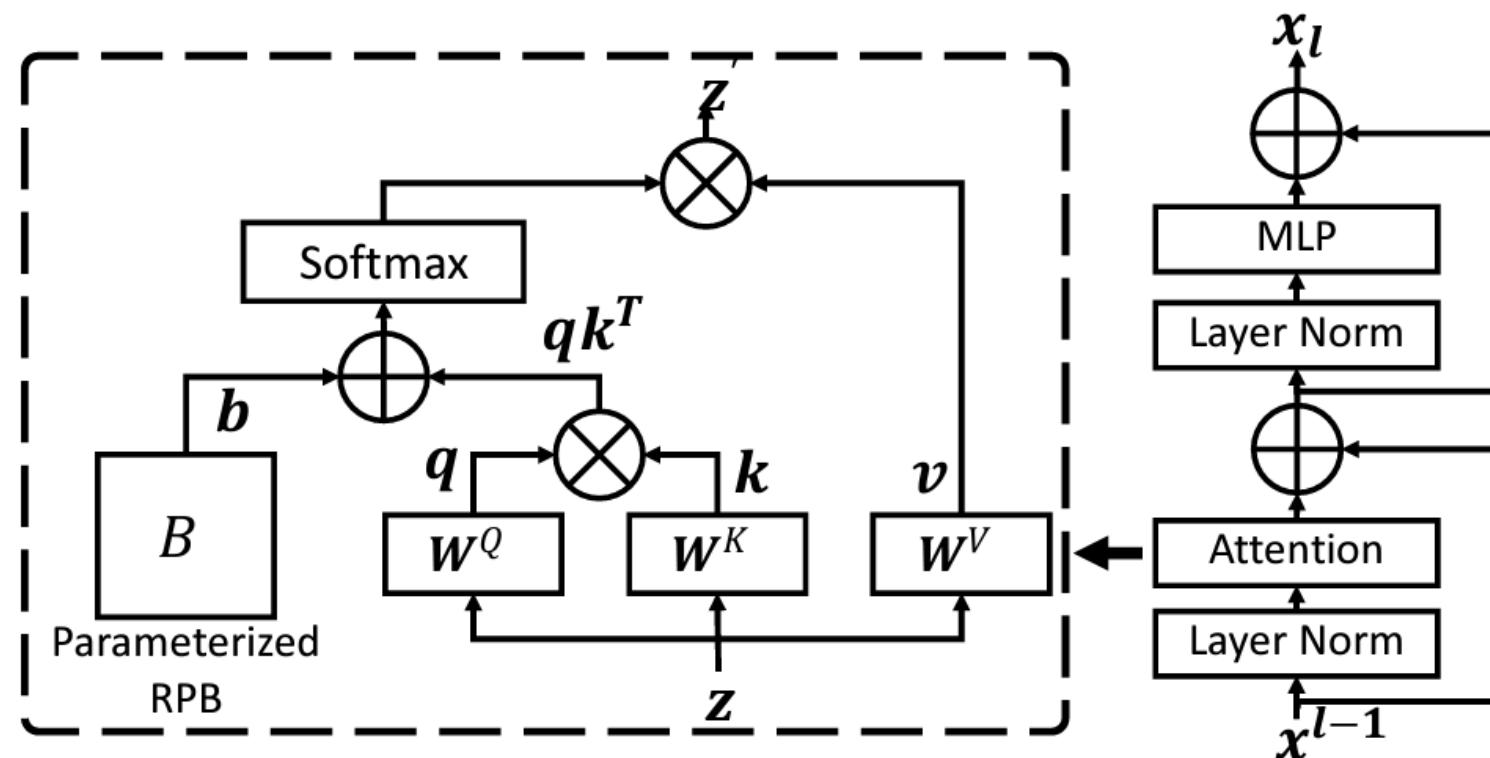


**TRANSFORMATEC**

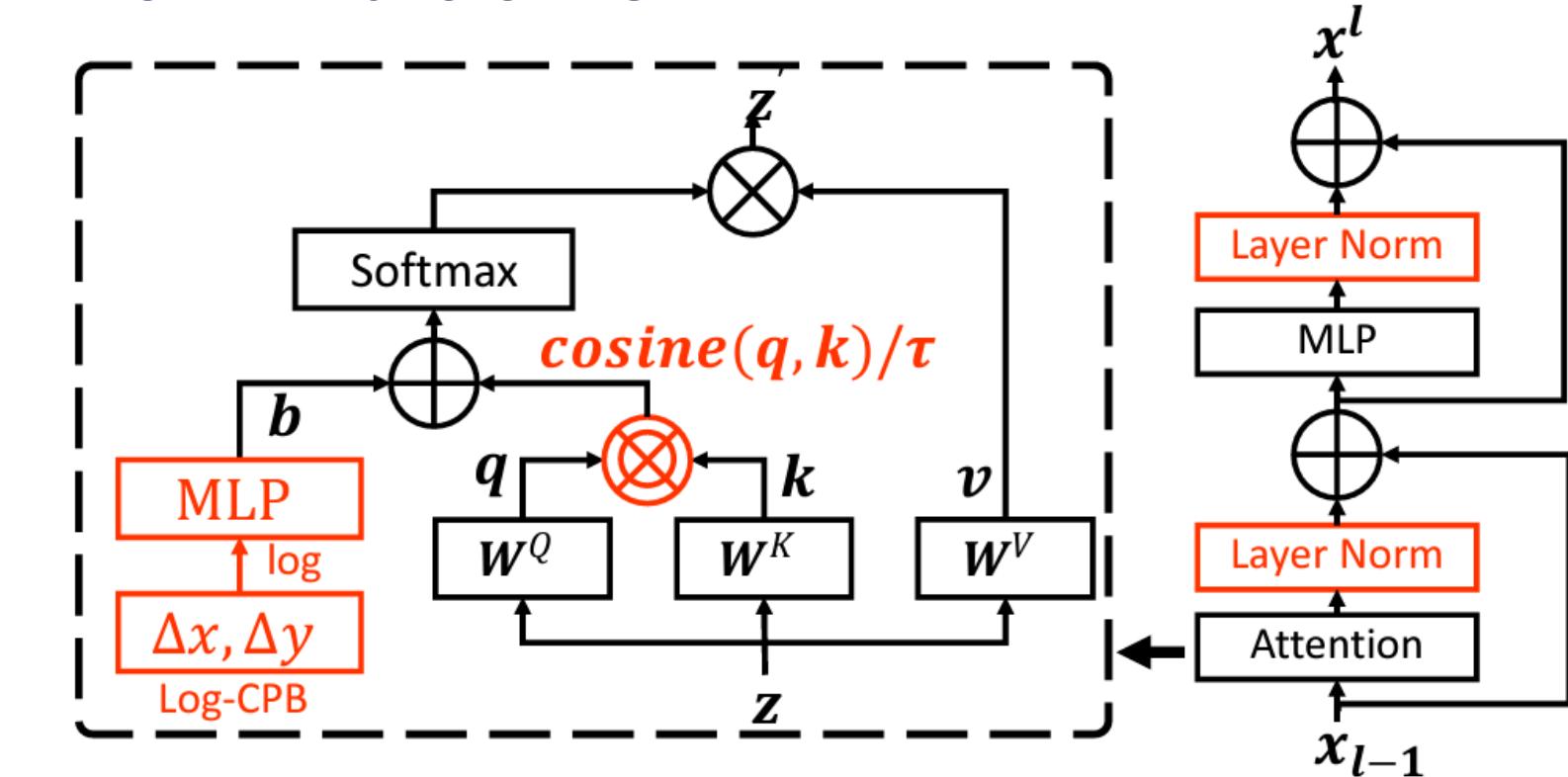
Ze Liu et al. (2022) "Swin Transformer V2: Scaling Up Capacity and Resolution".  
Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 12009-12019).

# Swin Transformer v2

Swin Transformer



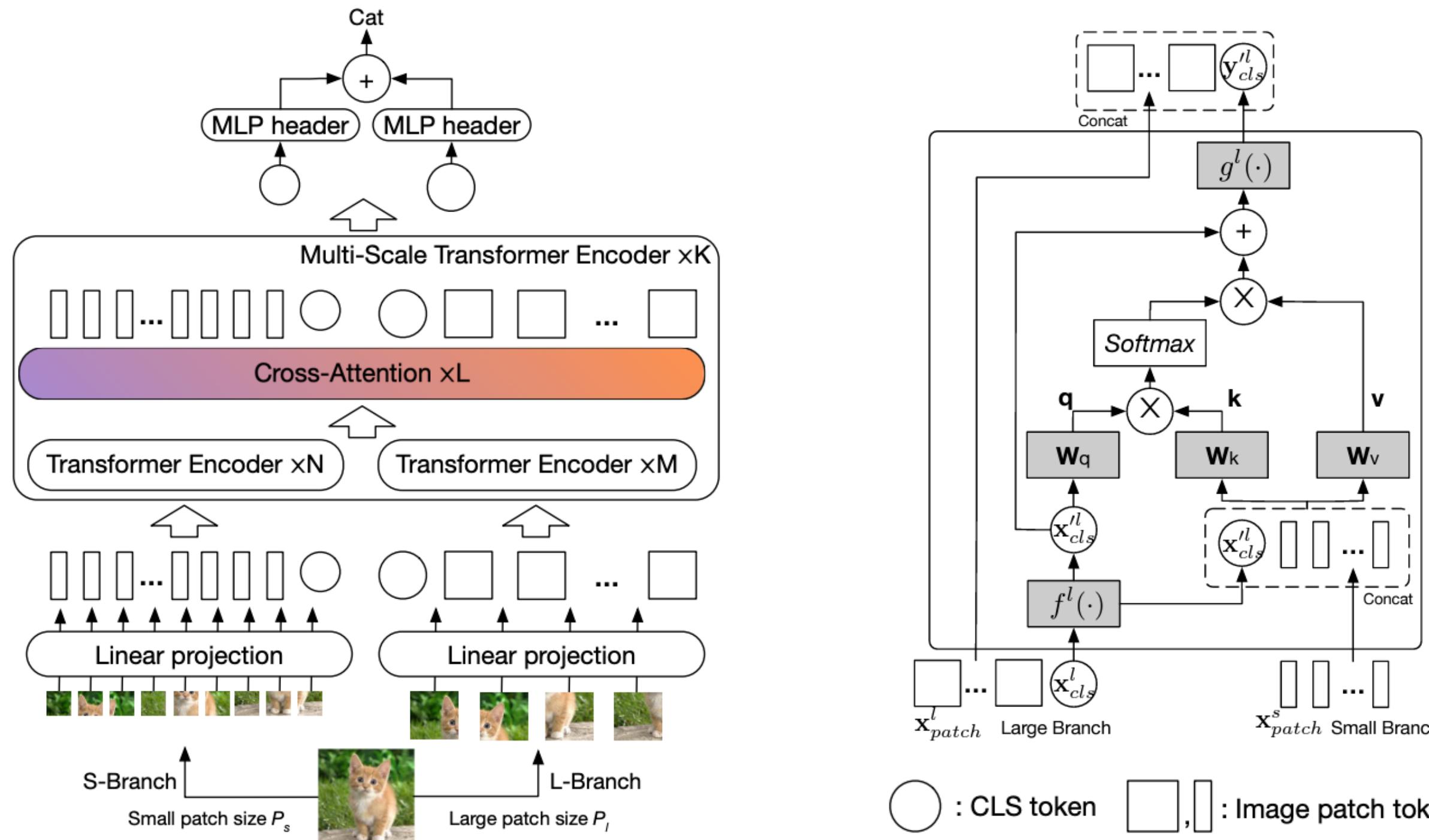
Swin Transformer V2



**TRANSFORMATEC**

Ze Liu et al. (2022) "Swin Transformer V2: Scaling Up Capacity and Resolution".  
Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 12009-12109).

# Cross ViT



# GRACIAS

*Victor Flores Benites*