

# Architecture de réseau de neurones

-

## Le Transformer

Guillaume Bourmaud

# PLAN

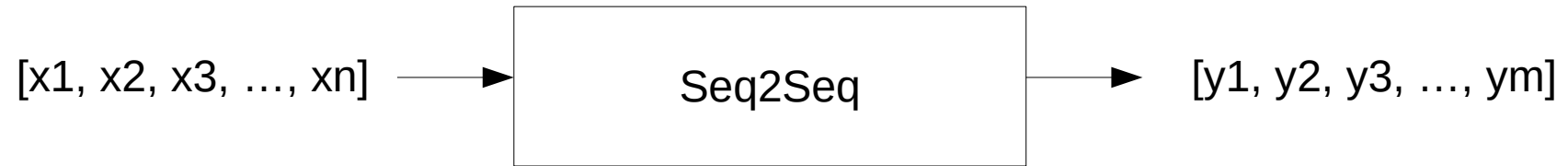
- I. Histoire du “Transformer”
- II. Couche d’attention à softmax
- III. Équivariance par permutation et encodage de la position
- IV. Application à des images
- V. Limites et tendances actuelles

# I) Histoire du “Transformer”

“Attention is All You Need”, NIPS 2017

1)

# Sequence-to-sequence

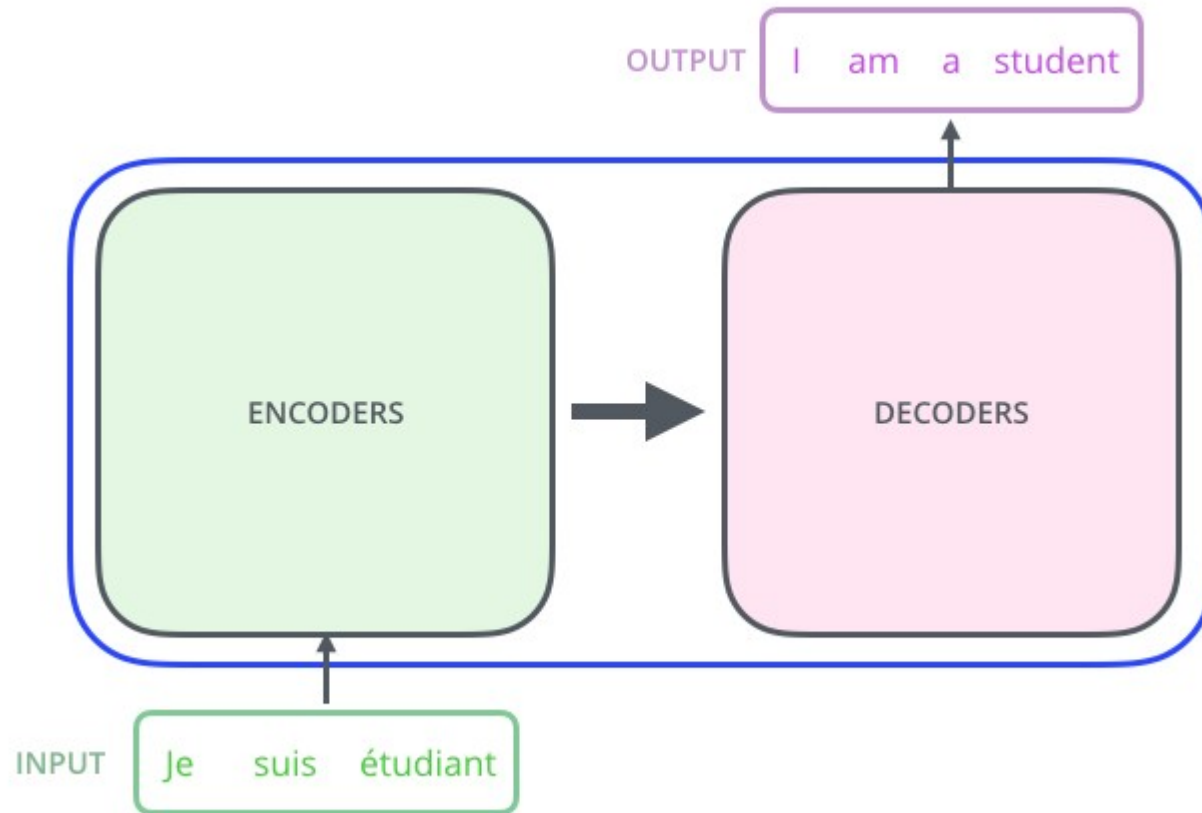


```
"the cat sat on the mat" -> [Seq2Seq model] -> "le chat etait assis sur le tapis"
```

1)

# Architecture encodeur-décodeur

Encodeur :  
- RNN  
- ou CNN  
- ou Transformer

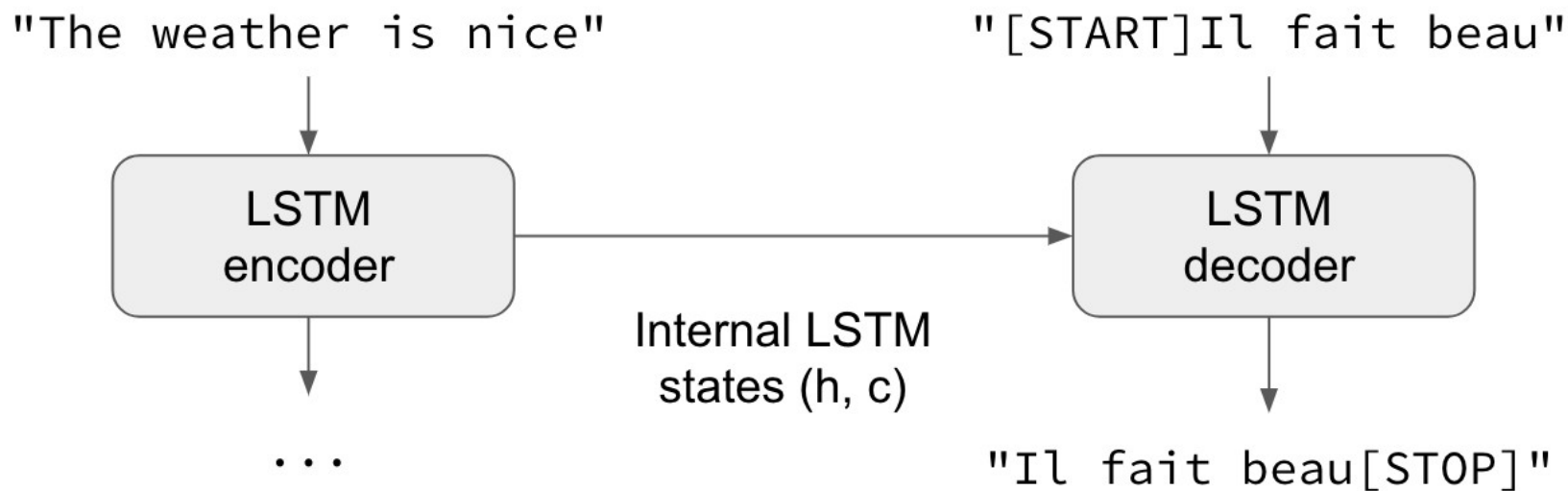


Décodeur :  
- RNN  
- ou CNN  
- ou Transformer

1)

# Exemple avec un RNN

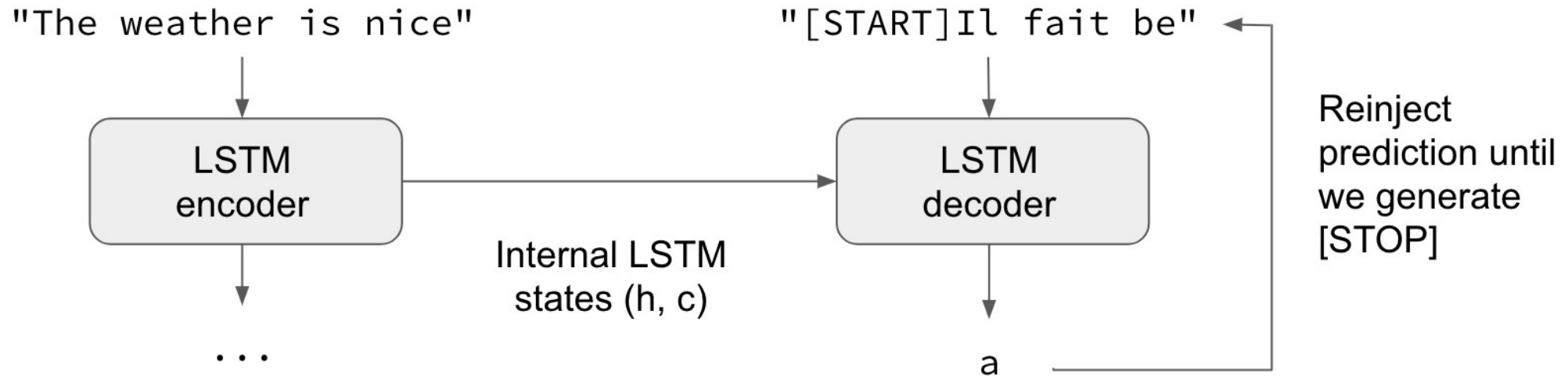
## Apprentissage (en "Teacher-Forcing")



l)

# Exemple avec un RNN (suite)

## Inférence

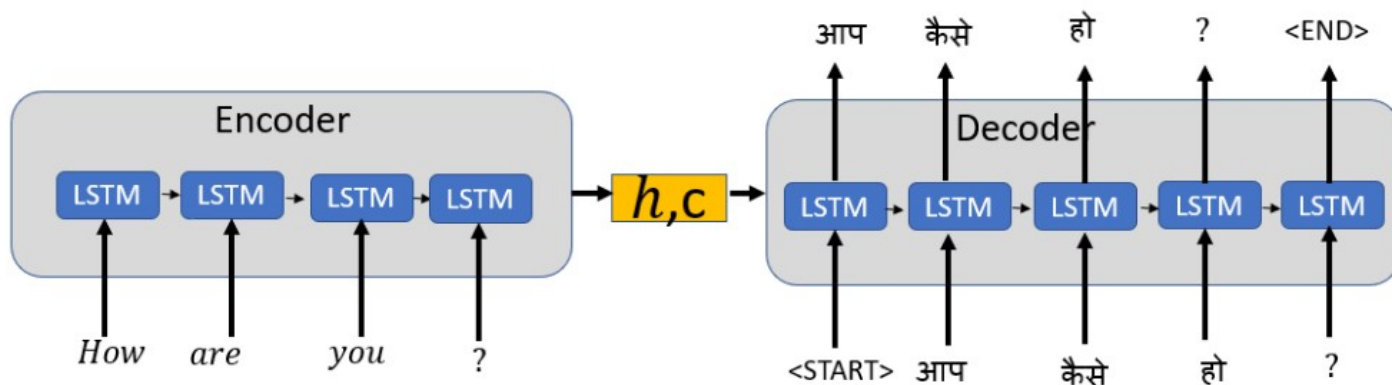


1)

# Exemple avec un RNN (suite)

## Limites

Calcul séquentiel  
↓  
Apprentissage lent  
↓  
Comment paralléliser ?



Difficile d'apprendre de  
longues dépendances  
avec un RNN.

Coronavirus pandemic is spread across 175 countries, it is a serious problem especially in Italy, Spain and US as of March 2020.

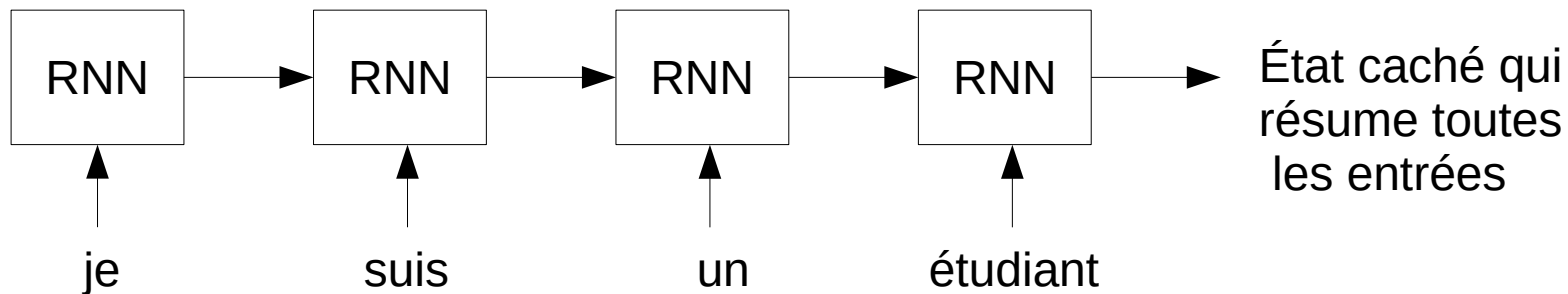


I)

# Transformer vs RNN

## Encodeur RNN

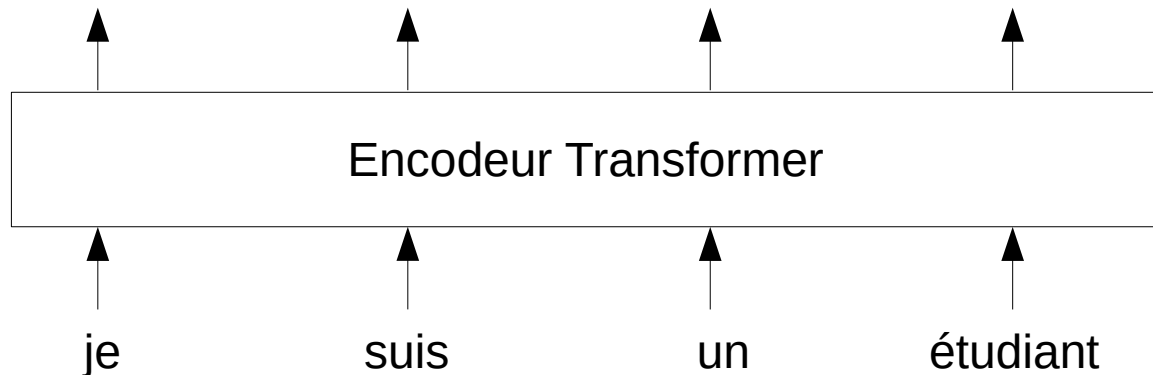
Calcul séquentiel  
+  
Chemin de la longueur  
de la phrase d'entrée



N entrées, N sorties. Chaque sortie dépend de **toutes** les entrées.

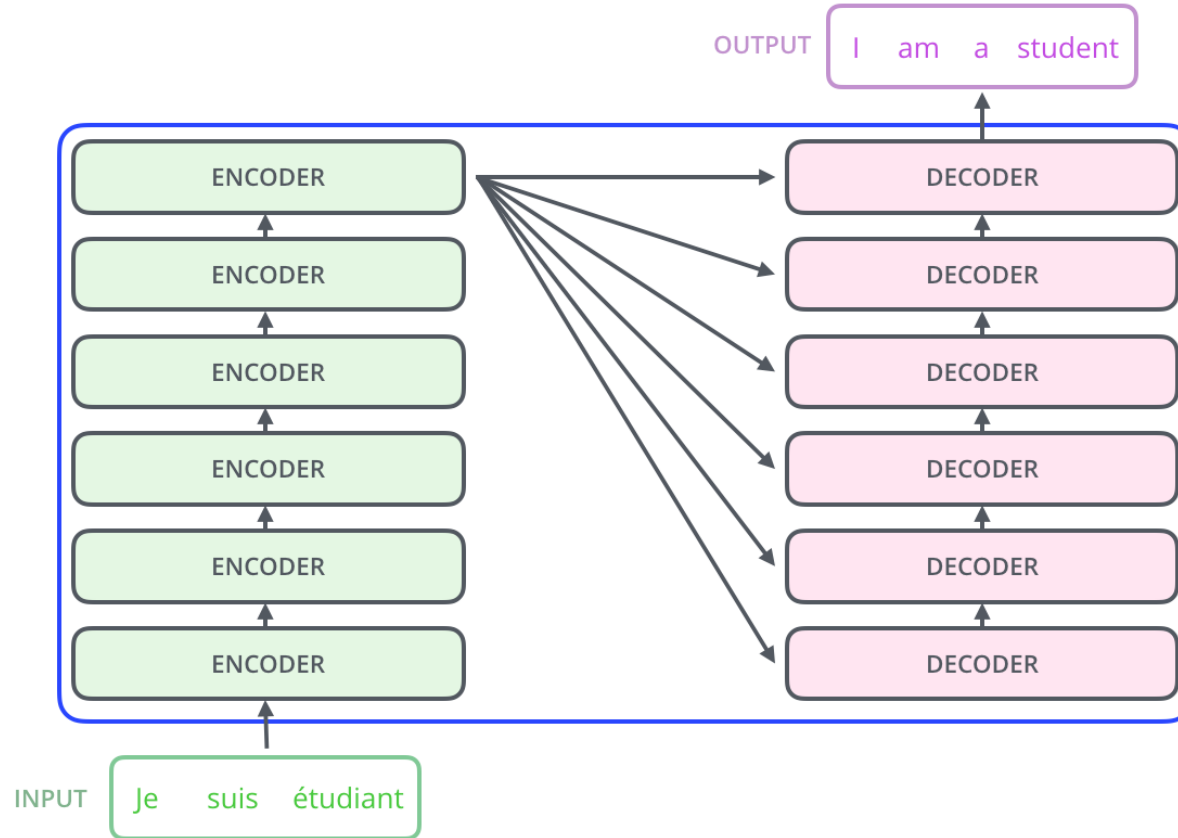
## Encodeur Transformer

Calcul en parallèle  
+  
Chemin de longueur  
constante



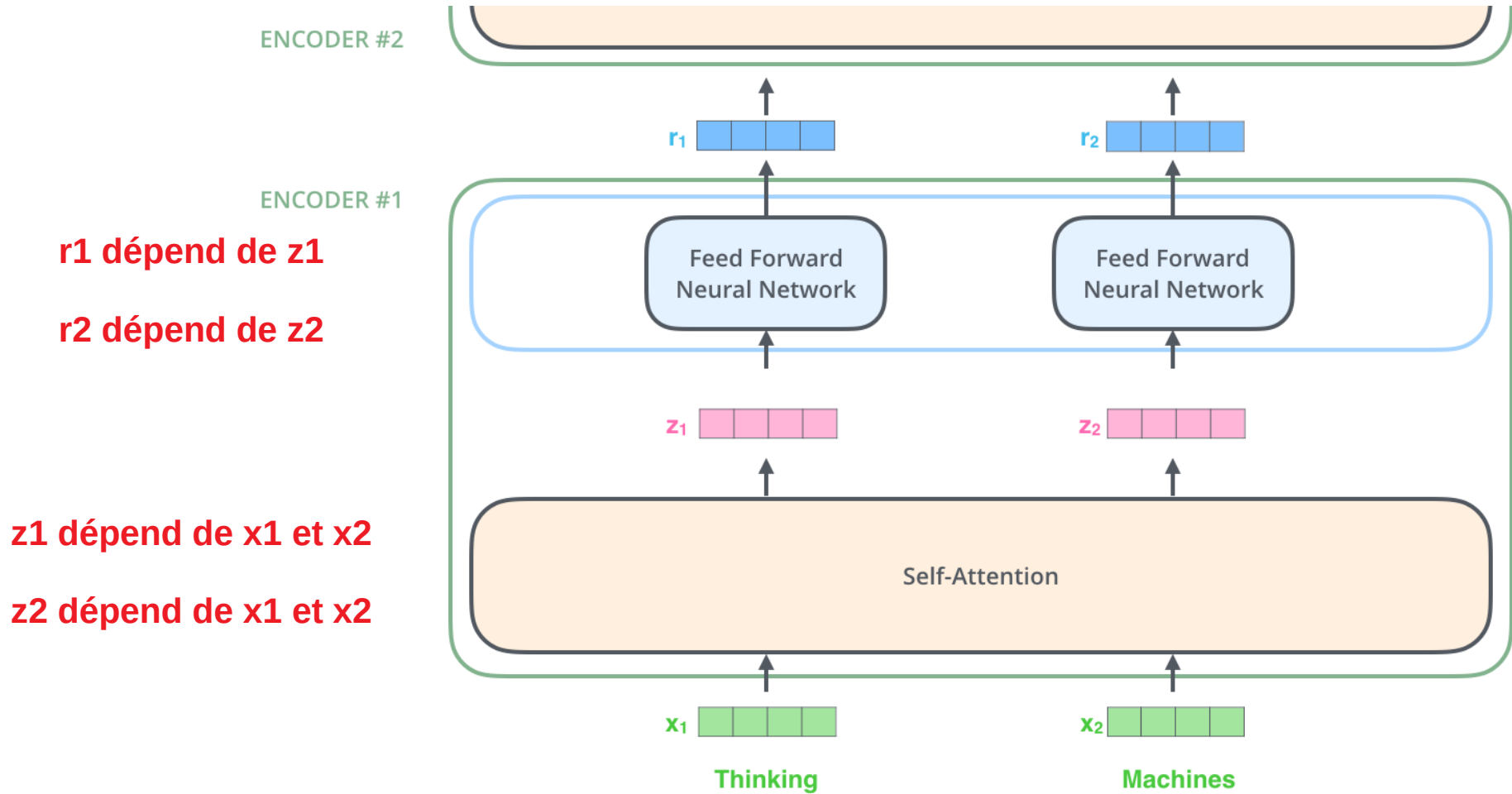
1)

# Transformer : vue globale



l)

# Transformer : encodeur



## II) Couche d'attention à softmax

# Attention utilisant la fonction softmax

Entrées	$\mathbf{x}$	: vecteur de dimension	$1 \times D_x$
	$\{\mathbf{y}_i\}_{i=1 \dots N_y}$	: ensemble de vecteurs de dimension	$1 \times D_y$
Sortie	$\mathbf{x}'$	: vecteur de dimension	$1 \times D_x$
Paramètres	$Q$	: matrice “query” de taille	$D_x \times L$
	$K$	: matrice “key” de taille	$D_y \times L$
	$V$	: matrice “value” de taille	$D_y \times D_x$
Fonction	$\mathbf{x}' = \mathbf{x} + \sum_{j=1}^{N_y} \underbrace{\frac{\exp(\mathbf{x}Q(\mathbf{y}_jK)^\top)}{\sum_{k=1}^{N_y} \exp(\mathbf{x}Q(\mathbf{y}_kK)^\top)}}_{\text{Softmax}} \mathbf{y}_jV$		

scalaire

II)

## Attention utilisant la fonction softmax (suite)

$$\mathbf{x}' = \mathbf{x} + \left( \sum_{j=1}^{N_y} s_j \mathbf{y}_j \right) \mathbf{V} \quad \text{où} \quad s_j = \frac{\exp(\mathbf{xQ}(\mathbf{y}_j\mathbf{K})^\top)}{\sum_{k=1}^{N_y} \exp(\mathbf{xQ}(\mathbf{y}_k\mathbf{K})^\top)}$$

Transfert d'information depuis  $\{\mathbf{y}_i\}_{i=1\dots N_y}$  vers  $\mathbf{x}$  le tout stocké dans  $\mathbf{x}'$

Q

K

V

permettent d'apprendre à transférer l'information pertinente pour la tâche concernée

Terminologie : “Dot-product attention”, plus rarement “softmax attention”

II)

## Inter-attention (“Cross-attention”)

**Entrées**  $\{\mathbf{x}_i\}_{i=1\dots N_x}$  : vecteurs de dimension  $D_x$   $\longrightarrow$   $\mathbf{X}$  : matrice de taille  $N_x \times D_x$   
 $\{\mathbf{y}_i\}_{i=1\dots N_y}$  : vecteurs de dimension  $D_y$   $\longrightarrow$   $\mathbf{Y}$  : matrice de taille  $N_y \times D_y$

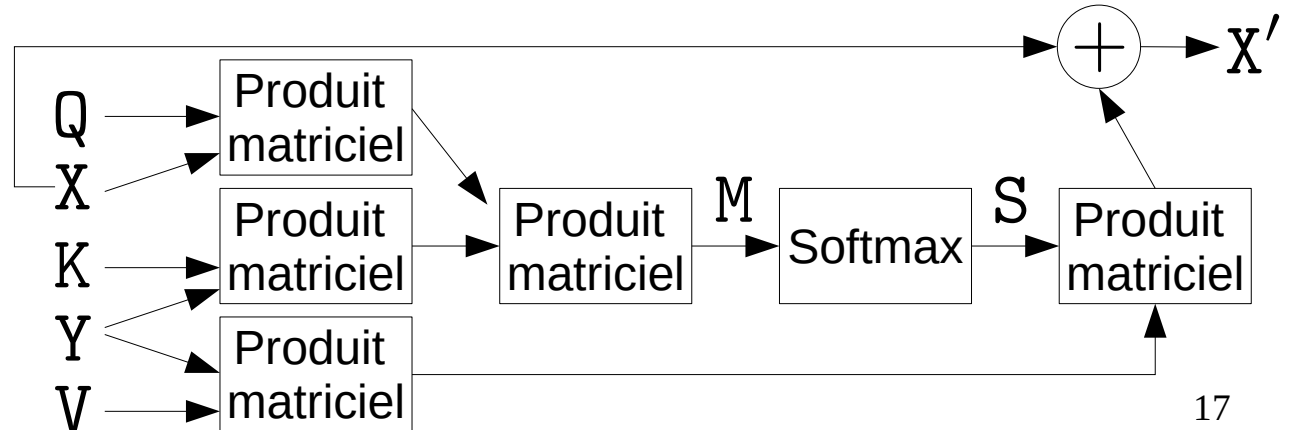
**Sorties**  $\{\mathbf{x}'_i\}_{i=1\dots N_x}$  : vecteurs de dimension  $D_x$   $\longrightarrow$   $\mathbf{X}'$  : matrice de taille  $N_x \times D_x$

$$\mathbf{M} = \mathbf{XQ}(\mathbf{YK})^\top$$

$$\mathbf{S} = \text{softmax}(\mathbf{M}, \text{dim}=1)$$

$$\mathbf{X}' = \mathbf{X} + \mathbf{SYV}$$

Transfert d'information depuis  $\mathbf{Y}$   
 vers  $\mathbf{X}$  le tout stocké dans  $\mathbf{X}'$



II)

## Inter-attention (“Cross-attention”) (suite)

$$\mathbf{X} : N_x \times D_x \quad \mathbf{Y} : N_y \times D_y \quad \mathbf{X}' : N_x \times D_x$$

---

$$\mathbf{M} = \mathbf{XQ}(\mathbf{YK})^\top : N_x \times N_y$$

$$\mathbf{S} = \text{softmax}(\mathbf{M}, \text{dim}=1) : N_x \times N_y$$

### Calcul

- Nombre d'opérations potentiellement très élevé

- + Parallélisable

### Stockage

- Mémoire requise pour stocker  $\mathbf{M}$  et  $\mathbf{S}$  potentiellement très élevée



II)

## Cas particulier : Auto-attention (“Self-attention”)

**Entrées**  $\{\mathbf{x}_i\}_{i=1\dots N_x}$  : vecteurs de dimension  $D_x$   $\longrightarrow$   $\mathbf{X}$  : matrice de taille  $N_x \times D_x$

---

**Sorties**  $\{\mathbf{x}'_i\}_{i=1\dots N_x}$  : vecteurs de dimension  $D_x$   $\longrightarrow$   $\mathbf{X}'$  : matrice de taille  $N_x \times D_x$

---

$$\mathbf{M} = \mathbf{XQ}(\mathbf{XK})^\top : N_x \times N_x$$

$$\mathbf{S} = \text{softmax}(\mathbf{M}, \text{dim}=1) : N_x \times N_x$$

$$\mathbf{X}' = \mathbf{X} + \mathbf{SXV} : N_x \times D_x$$

Transfert d'information depuis  $\mathbf{X}$  vers lui-même le tout stocké dans  $\mathbf{X}'$

II)

## Couche d'attention softmax à têtes multiples "Multi-head dot-product attention"

$$\mathbf{X} : N_x \times D_x \quad \mathbf{Y} : N_y \times D_y \quad \mathbf{X}' : N_x \times D_x$$

---

Pour  $k = 1$  à  $H$

1)  $H$  couches d'attention indépendantes (têtes)

$$\mathbf{M}_k = \mathbf{XQ}_k(\mathbf{YK}_k)^\top : N_x \times N_y$$

$$\mathbf{S}_k = \text{softmax}(\mathbf{M}_k, \text{dim}=1) : N_x \times N_y$$

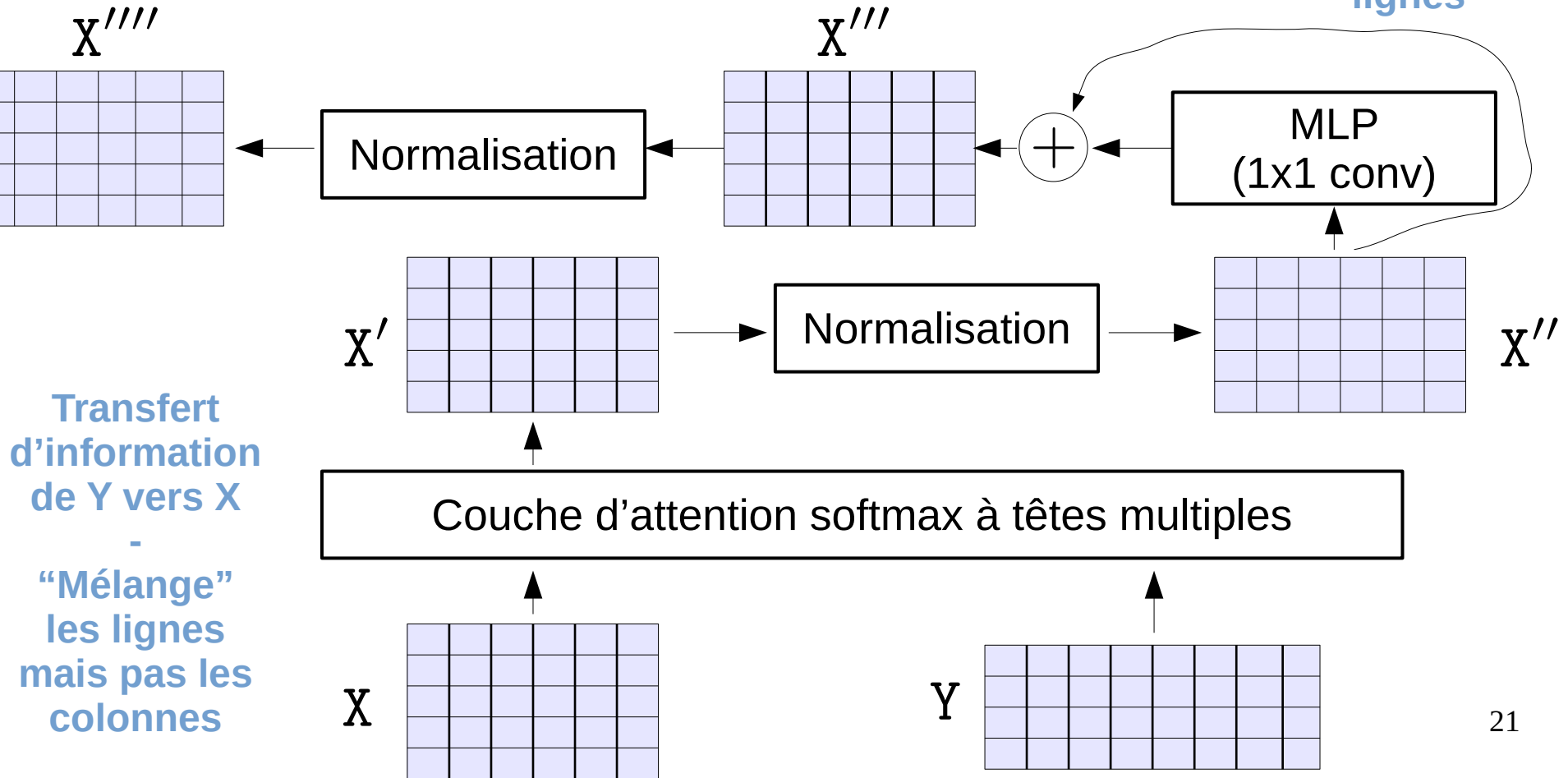
$$\mathbf{Z}_k = \mathbf{S}_k \mathbf{YV}_k : N_x \times D_v$$

$$\mathbf{Z} = [\mathbf{Z}_1 \mathbf{Z}_1 \dots \mathbf{Z}_H] : N_x \times (H \times D_v) \quad 2) \text{ "mélange" des résultats des } H \text{ têtes}$$

$$\mathbf{X}' = \mathbf{X} + \mathbf{ZW} : N_x \times D_x$$

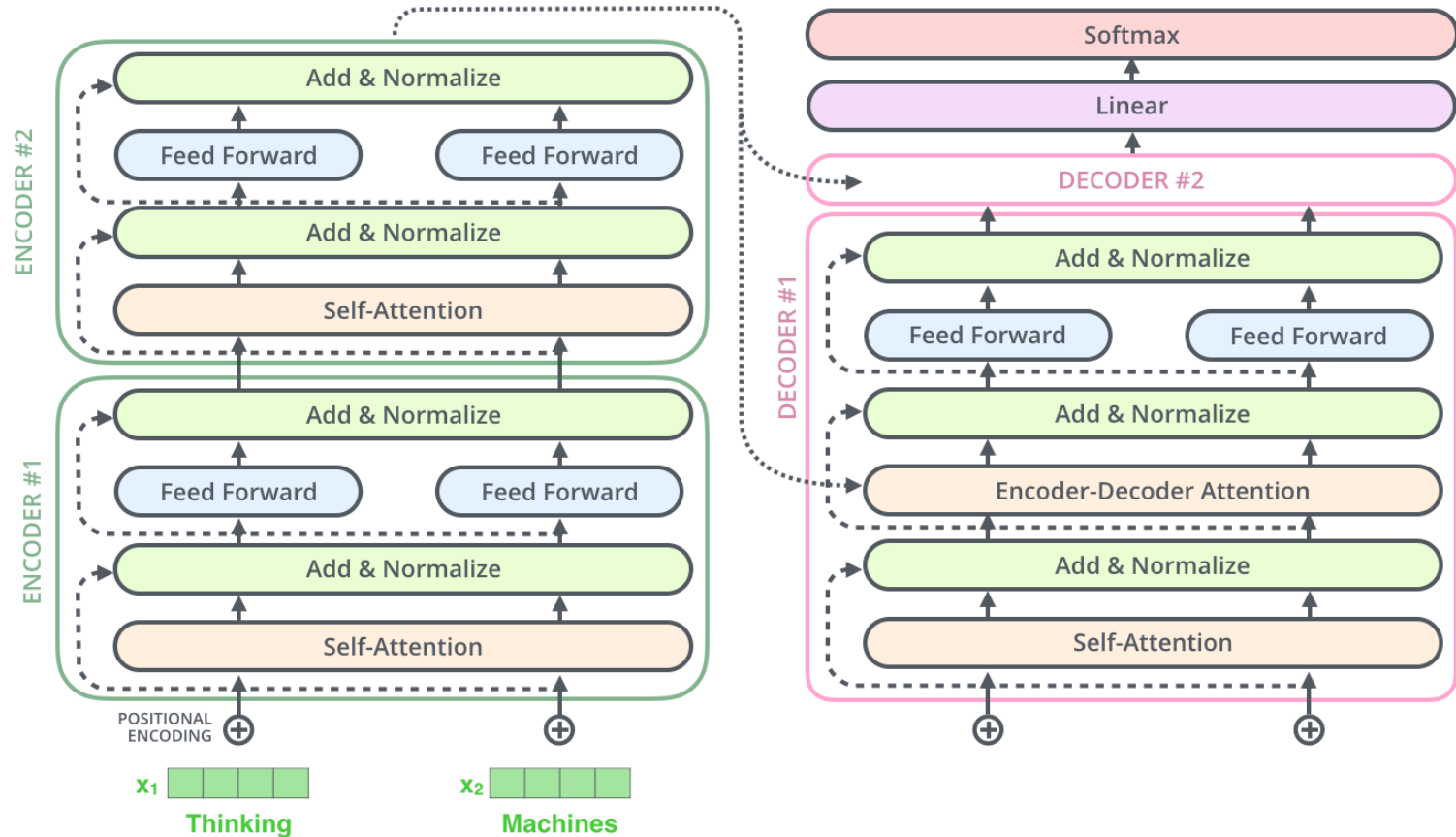
II)

# Bloc d'attention "classique"



II)

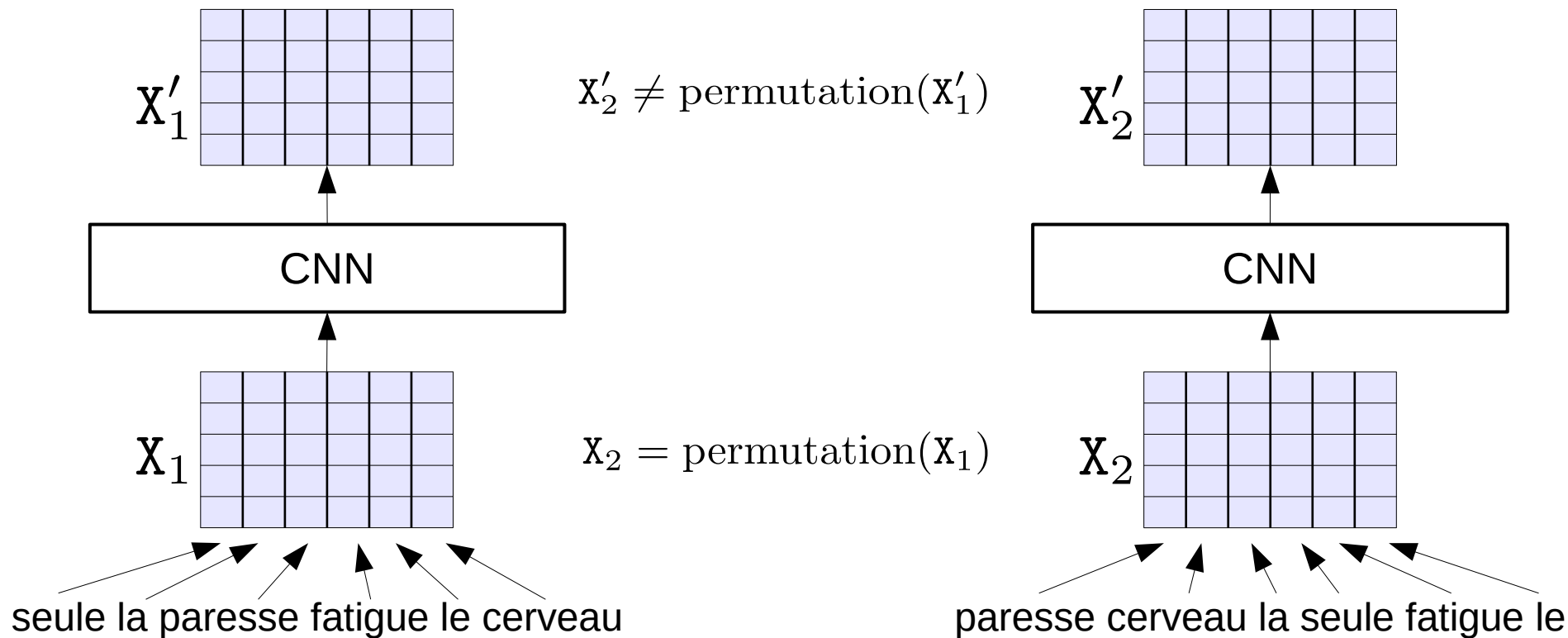
# Vue détaillée du Transformer



### III) Équivariance par permutation et encodage de la position

III)

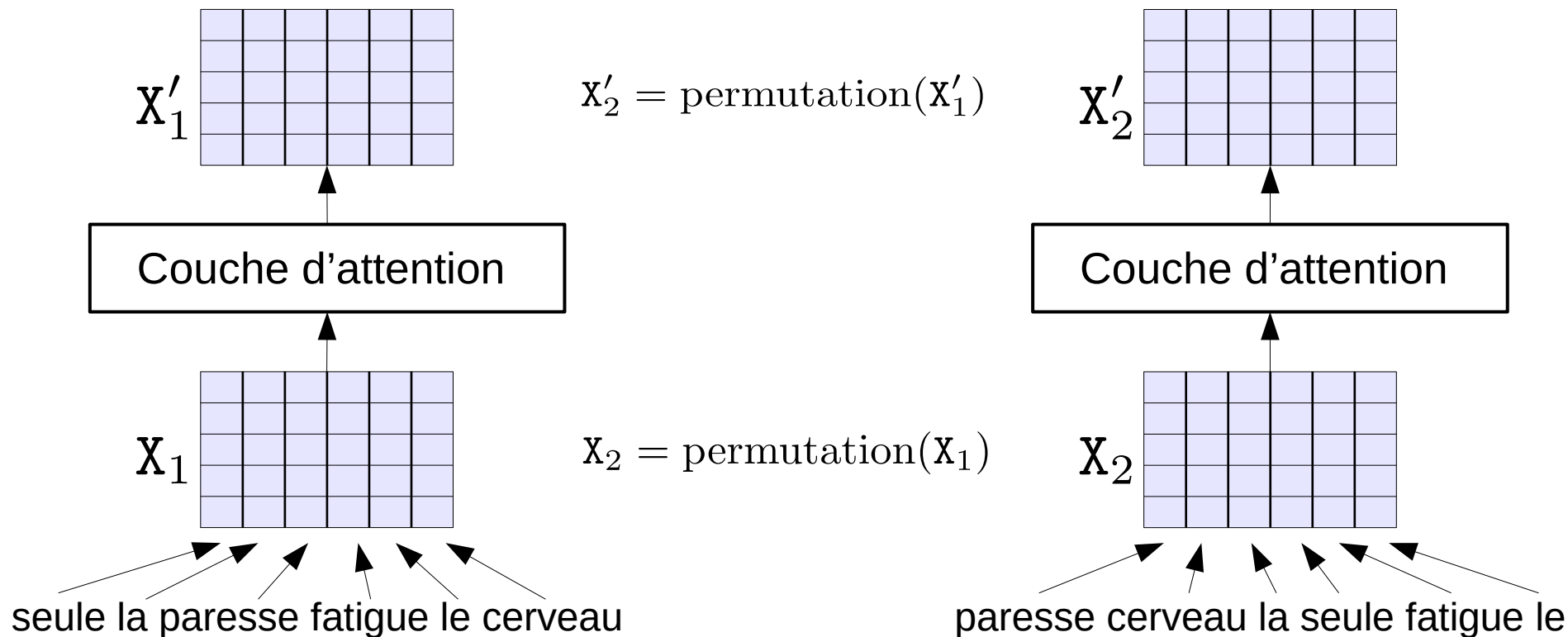
# Équivariance par permutation



**CNN adapté aux ensembles ordonnés (phrase, signal, image, etc.)**  
→ ses filtres s'appliquent sur un voisinage

III)

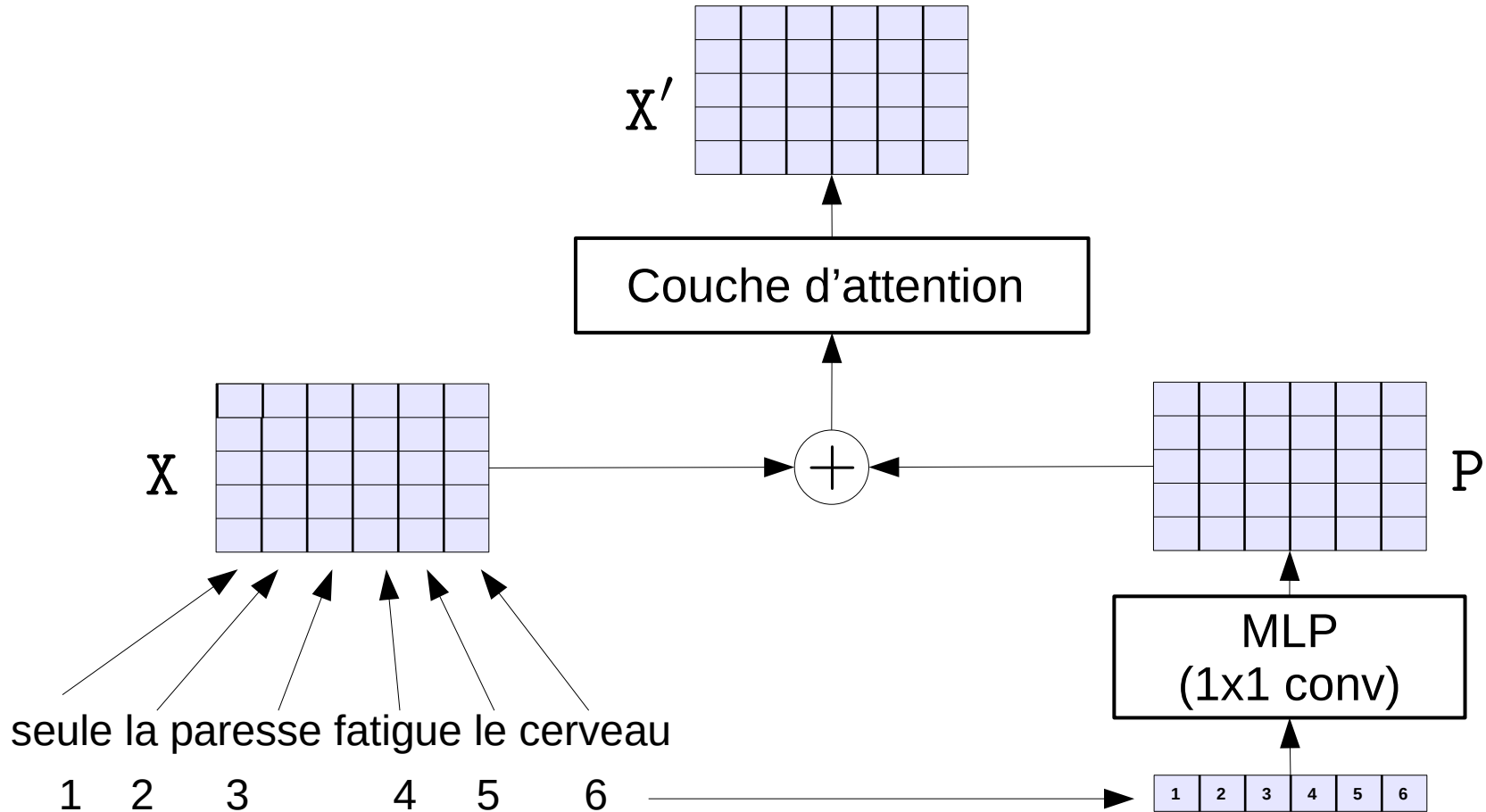
## Équivariance par permutation (suite)



**Pour un ensemble ordonné (phrase, signal, image, etc.)  
→ nécessité d'encoder la position**

III)

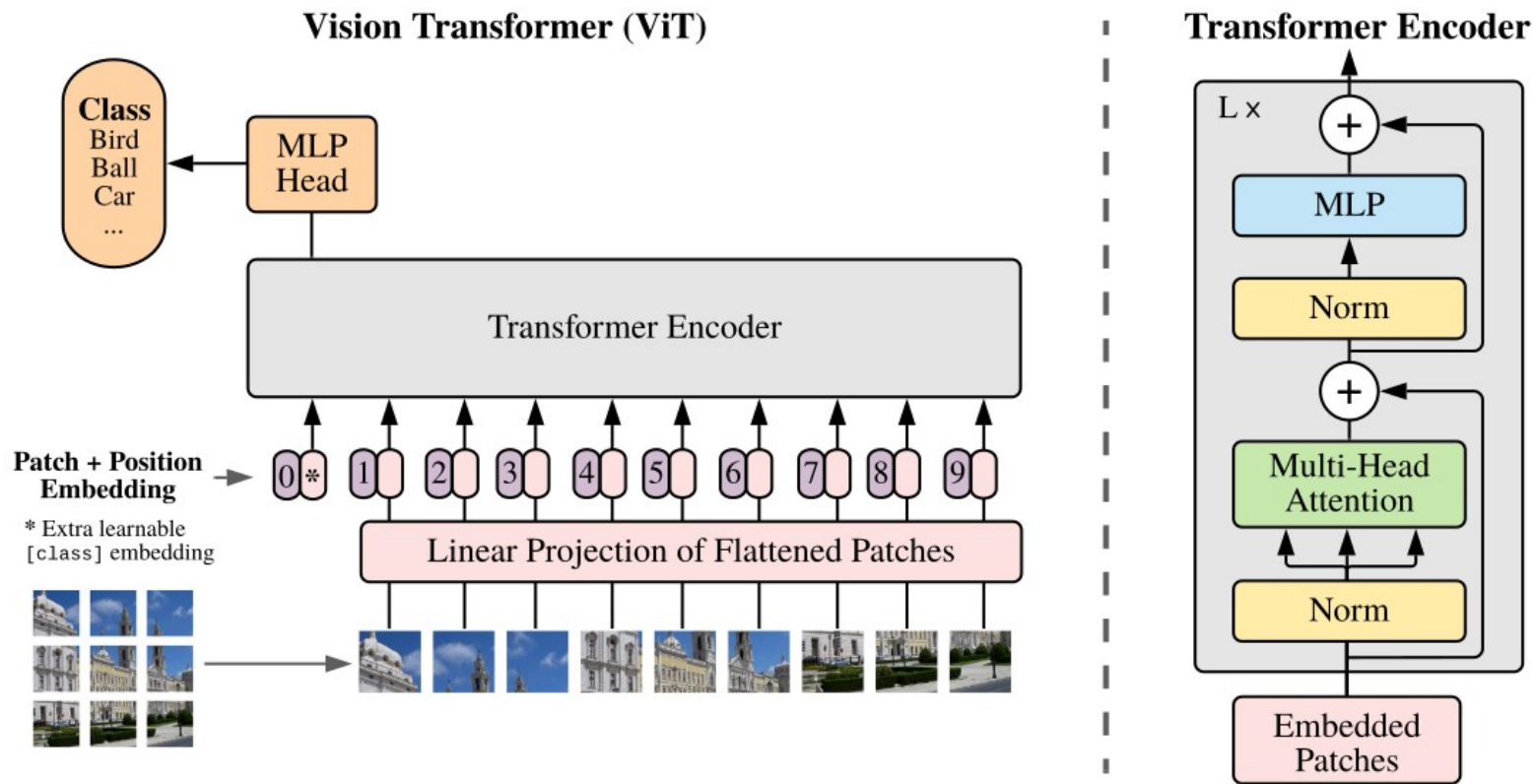
## Encodage de la position





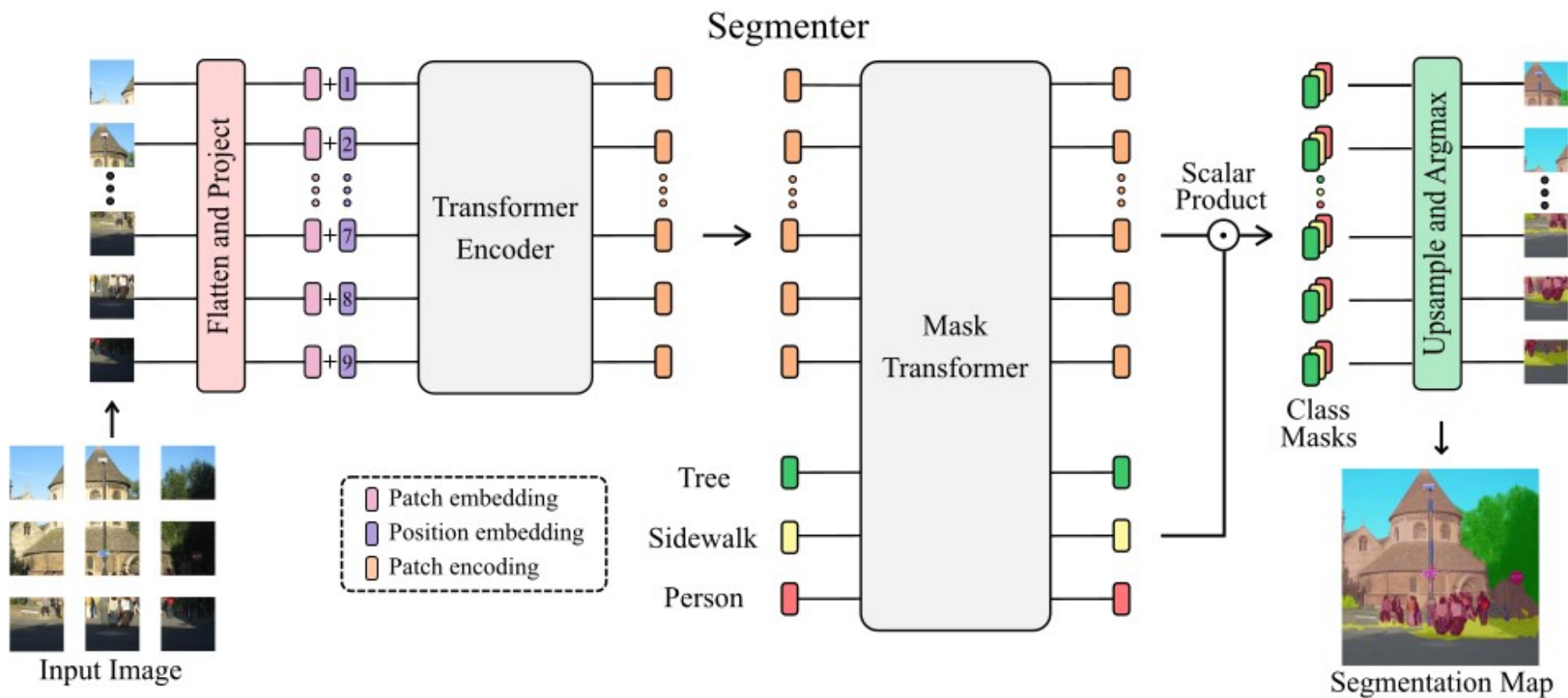
## IV) Application à des images

“An image is worth 16x16 words”, ICLR 2021



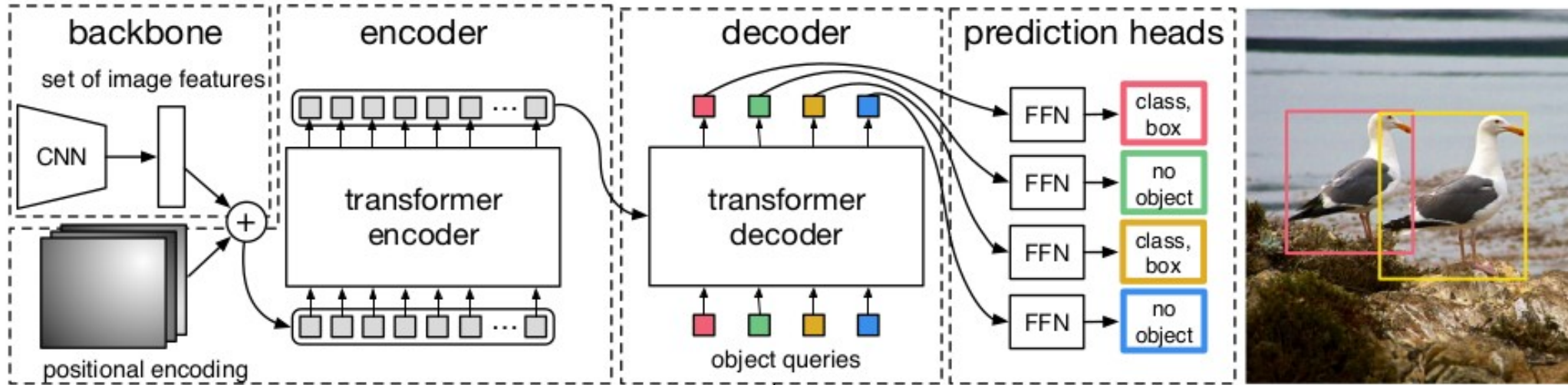
IV)

# “Segmenter: Transformer for Semantic Segmentation”, ICCV 2021



IV)

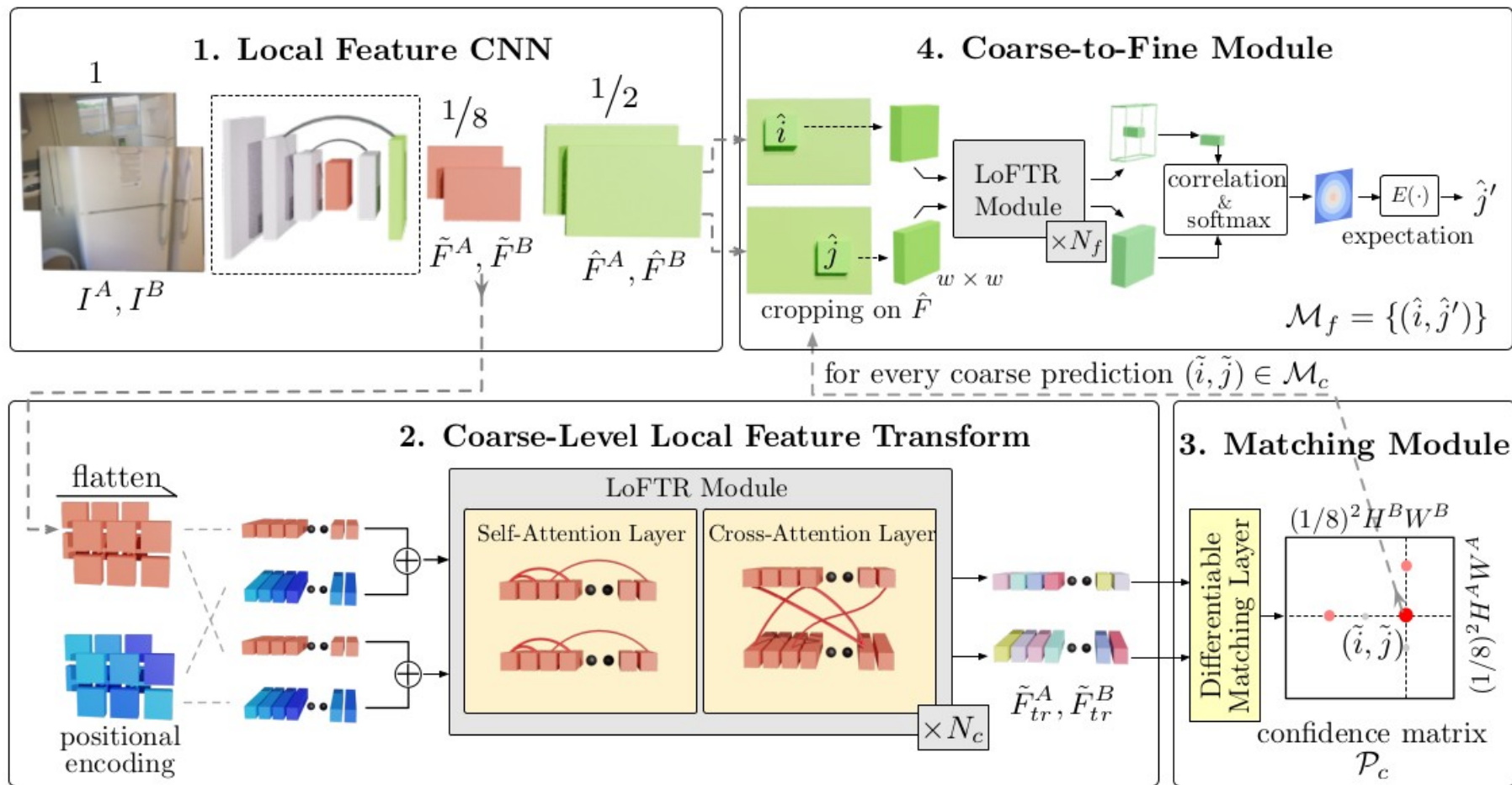
## “DETR: End-to-End Object Detection With Transformers”, ECCV 2020



En pratique, il y a 100 “object queries”, donc 100 boîtes englobantes prédites.

IV)

# “LoFTR: Detector-Free Local Feature Matching with Transformers”, CVPR 2021



## V) Limites et tendances actuelles

v)

## Limites des couches d'attention à softmax

$$M = XQ(YK)^{\top} : N_x \times N_y \quad S = \text{softmax}(M, \text{dim}=1) : N_x \times N_y$$

Problème : Inapplicable pour des ensembles de grandes tailles.

Solutions :

- Appliquer des couches d'attention à softmax en cherchant à réduire  $N_x$  et/ou  $N_y$
- Modifier la couche d'attention softmax

v)

# Exemple de modification de la couche d'attention softmax : couche d'attention linéaire

“Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention”, ICML 2020

$$\mathbf{x}' = \mathbf{x} + \left( \sum_{j=1}^{N_y} \frac{\phi(\mathbf{xQ})\phi(\mathbf{y}_j\mathbf{K})^\top}{\sum_{k=1}^{N_y} \phi(\mathbf{xQ})\phi(\mathbf{y}_k\mathbf{K})^\top} \mathbf{y}_j \right) \mathbf{V}$$

Remplacement du  
noyau exponentiel par  
un noyau linéaire

Se simplifie en



$$\mathbf{x}' = \mathbf{x} + \phi(\mathbf{xQ}) \frac{\sum_{j=1}^{N_y} \phi(\mathbf{y}_j\mathbf{K})^\top \mathbf{y}_j \mathbf{V}}{\phi(\mathbf{xQ}) \sum_{k=1}^{N_y} \phi(\mathbf{y}_k\mathbf{K})^\top}$$

Indépendant de  $\mathbf{x}$ , plus  
besoin de calculer ni  
stocker explicitement  
les matrices  $\mathbf{M}$  et  $\mathbf{S}$



# Exemple de réduction de $N_x$ et/ou $N_y$ : PerceiverIO

“Perceiver IO: A General Architecture for Structured Inputs & Outputs.” arXiv, 2021

