

# Les Transformers

## principes et architecture

Une liste (partiellement ordonnée) de sources

Prérequis

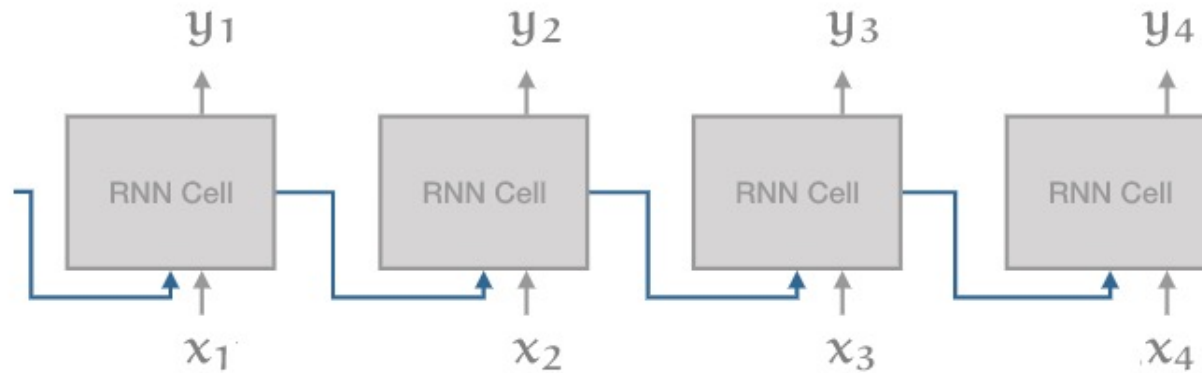
[Seq2Seq with Attention and Beam Search](#)

Articles choisis sur le(s) Transformer(s)

1. [The Illustrated Transformer](#) (introduction visuelle)
2. [Transformers from Scratch](#) (présentation concise + code)
3. [Attention is All You Need](#) (article original)
4. [The Annotated Transformer](#) (article original + code commenté)
5. [BERT – Pre-training of Deep Bidirectional Transformers for Language Understanding](#) (article original)
6. [Transformers are Graph Neural Networks](#) (lien avec les GNN)
7. [Efficient Transformers: A Survey](#) (variantes pour des textes longs)
8. [On the Opportunities and Risks of Foundation Models](#) (analyse critique pluridisciplinaire)



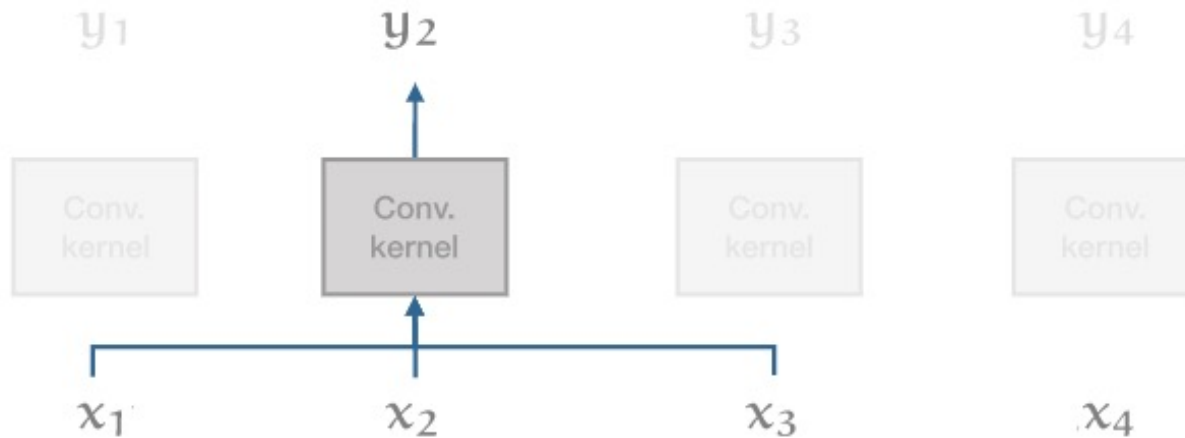
# Il était une fois les modèles seq2seq...



RNN's

Atout : longueur arbitraire

Inconvénient : non parallélisable



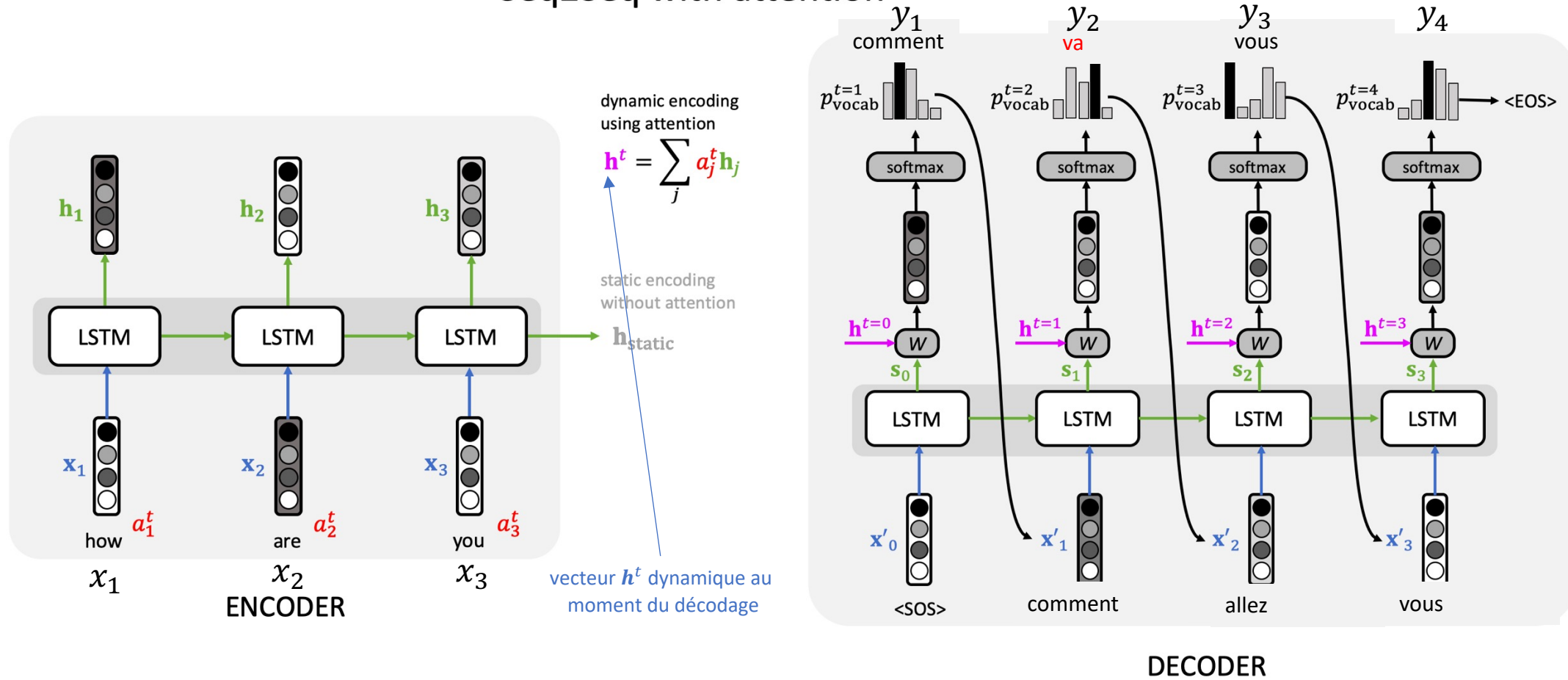
1D – CNN's

Atout : parallélisable

Inconvénient : longueur limitée

# Il était une fois le mécanisme d'attention...

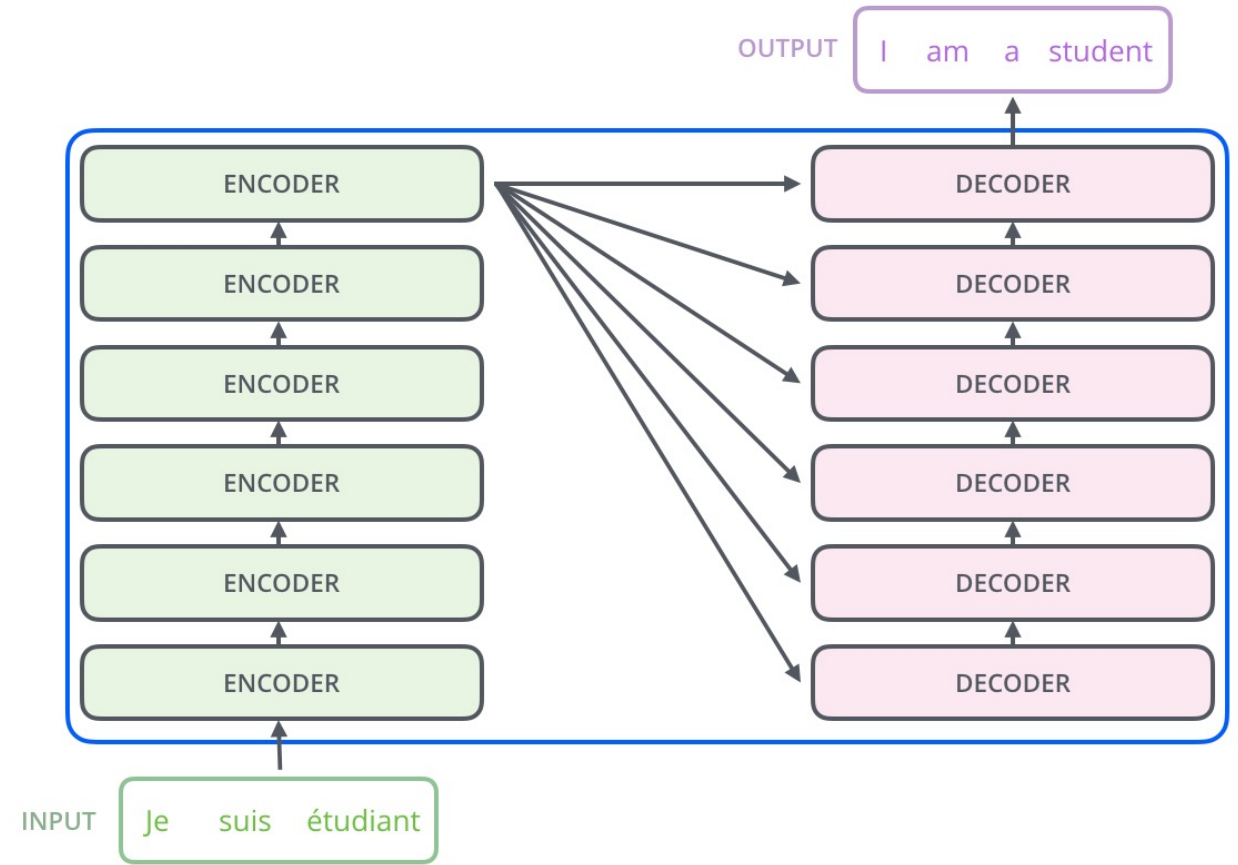
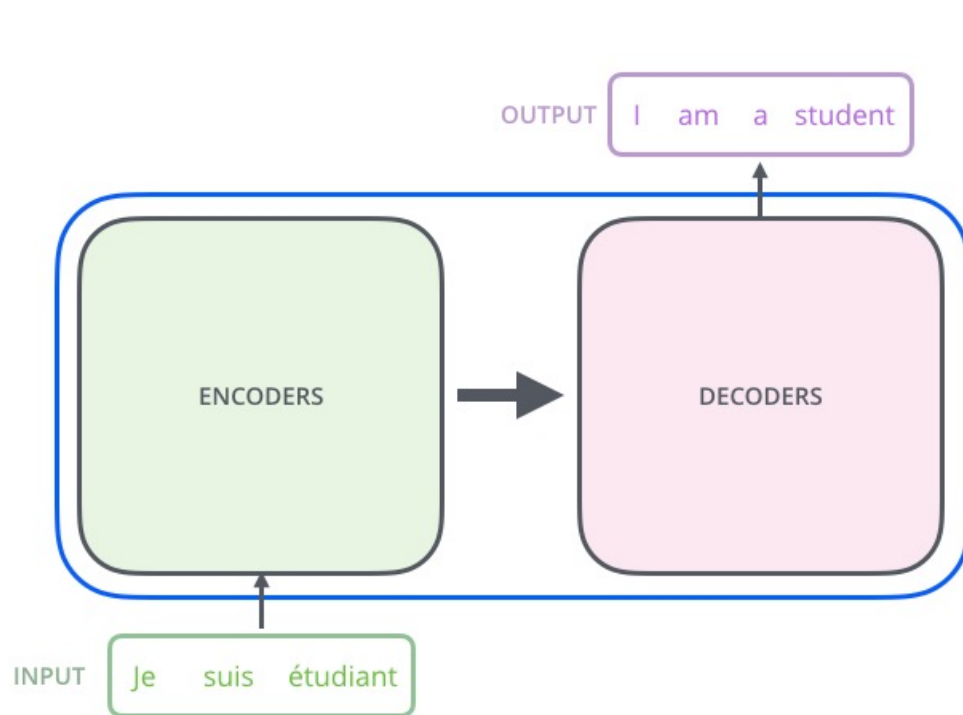
## Seq2Seq with attention



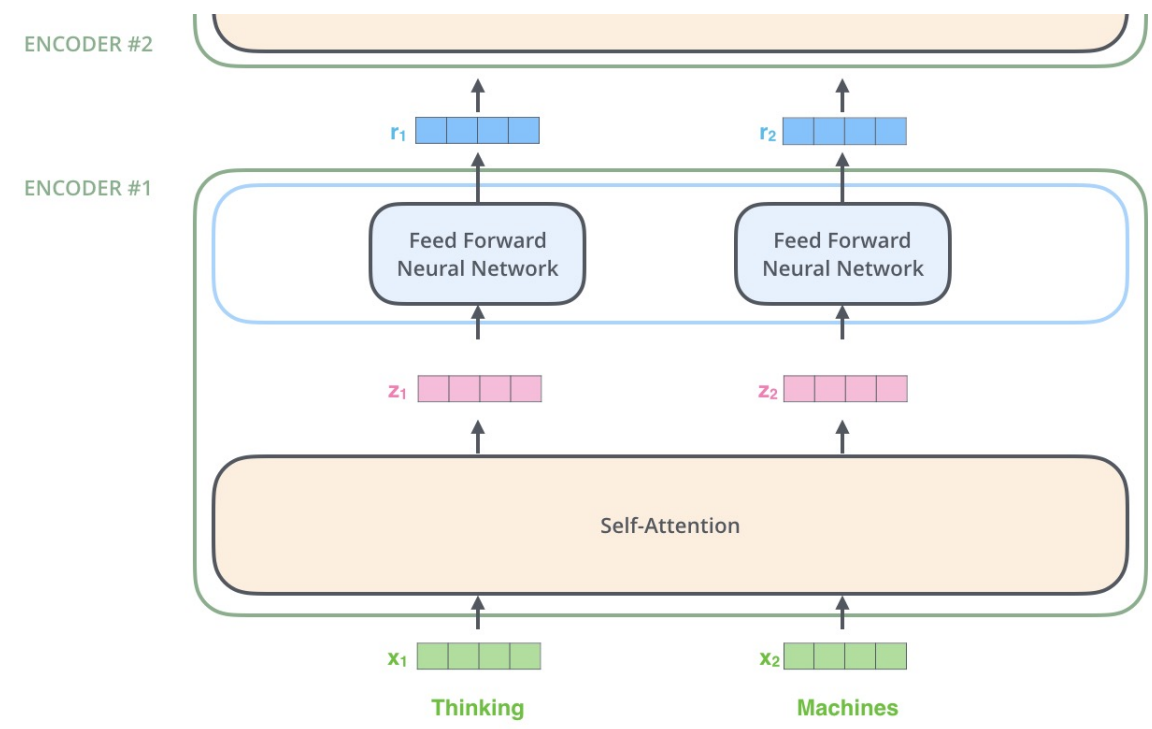
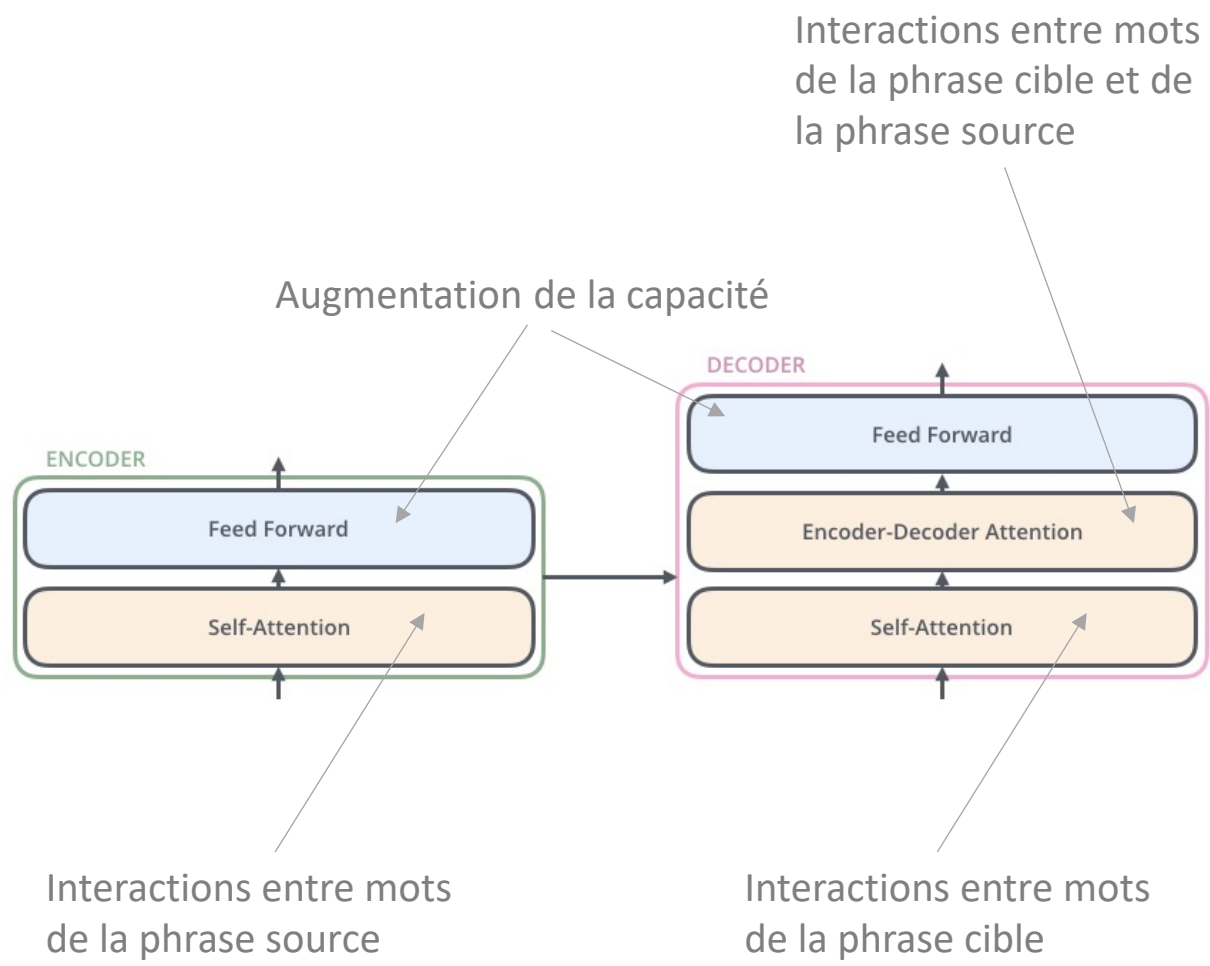
But : modéliser la distribution  $P(y_{t+1} | y_1, \dots, y_t; x_1, \dots, x_n)$

But : concilier les atouts des RNN et des CNN

Le Transformer vu comme un ensemble de poupées russes



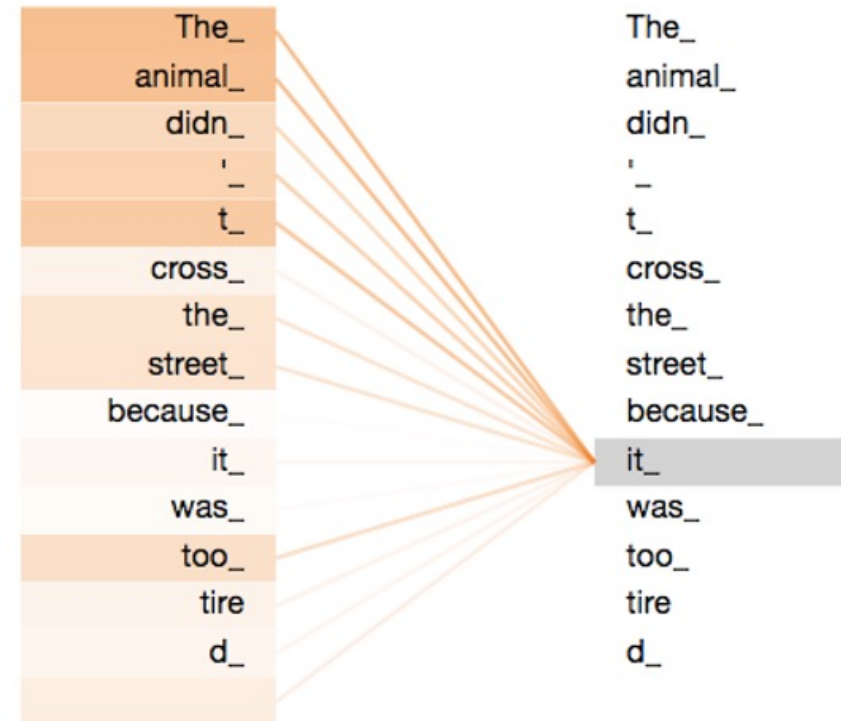
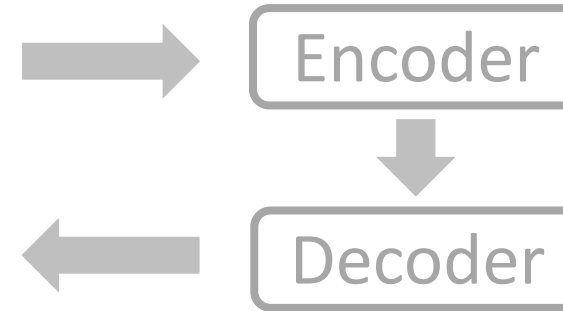
# Ouvrons le Transformer !



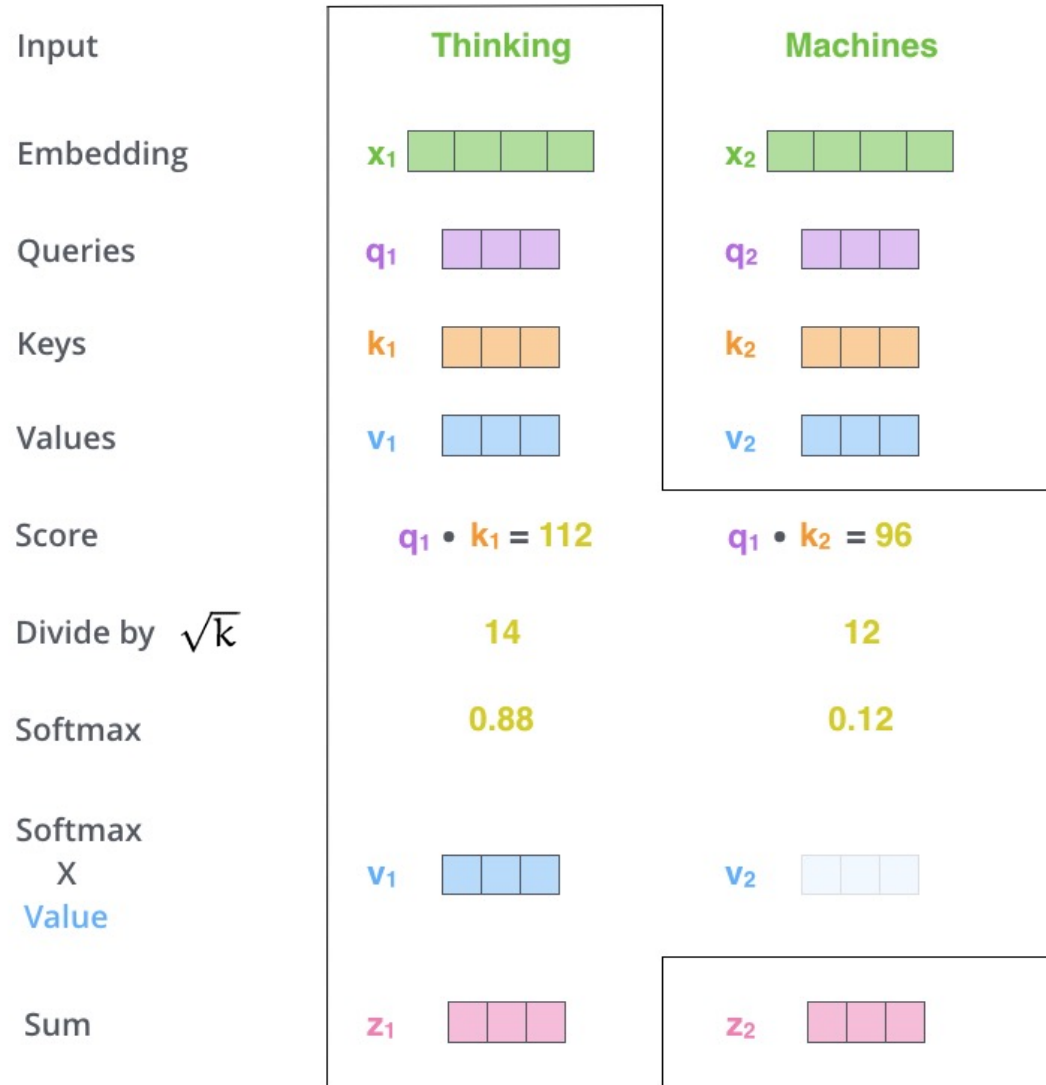
# Utilité d'un mécanisme d'attention

The **animal** didn't cross the **street** because **it** was too **wide**.  
 The **animal** didn't cross the **street** because **it** was too **tired**.

L'**animal** n'a pas traversé la **rue** car **elle** était trop **large**.  
 L'**animal** n'a pas traversé la **rue** car **il** était trop **fatigué**.



# Implémentation du MA : key, query, value



On associe 3 vecteurs dans  $\mathbb{R}^k$  à chaque vecteur  $x_i$  :

- Un vecteur « value »  $v_i$  : contribution
- Un vecteur « key »  $k_i$  : identification
- Un vecteur « query »  $q_i$  : recherche

3 matrices à apprendre

$$q_i = W_q x_i \quad k_i = W_k x_i \quad v_i = W_v x_i$$

$$w'_{ij} = \frac{q_i^T k_j}{\sqrt{k}} \quad i, j = 1, \dots, L$$

$$w_{ij} = \text{softmax}(w'_{ij})$$

$$z_i = \sum_j w_{ij} v_j$$

Transformers



$$O(k L^2)$$

$$\Rightarrow L < k$$

RNN

$$O(k^2 L)$$



# Formulation matricielle du calcul de l'attention

Chaque ligne correspond à un mot en entrée

$$\begin{matrix} \mathbf{X} & & \mathbf{W}^Q & = & \mathbf{Q} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} & \times & \begin{matrix} \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \end{matrix} & = & \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

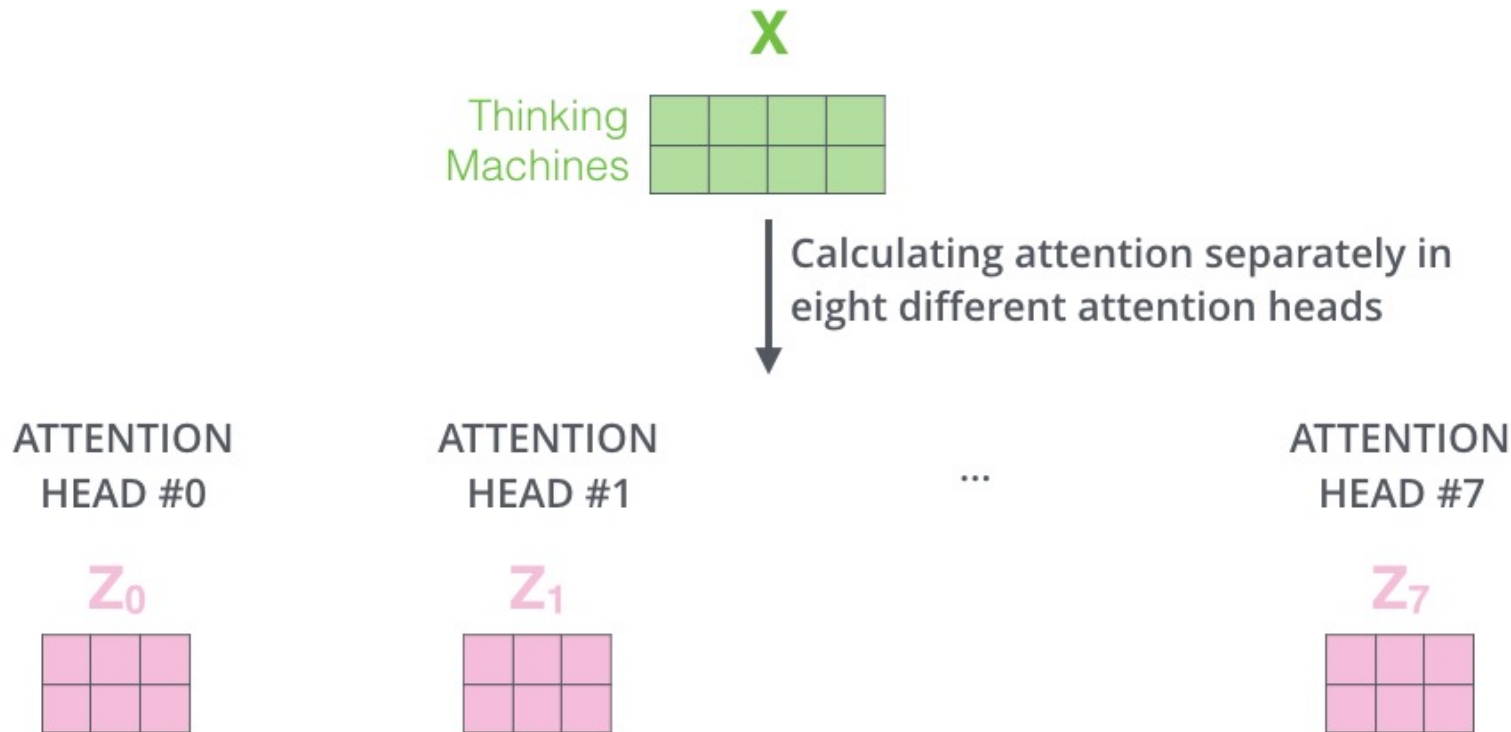
$$\begin{matrix} \mathbf{X} & & \mathbf{W}^K & = & \mathbf{K} \\ \begin{matrix} \square & \square & \square & \square \\ \square & \square & \square & \square \end{matrix} & \times & \begin{matrix} \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \\ \square & \square & \square & \square \end{matrix} & = & \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

$$\begin{matrix} \mathbf{X} & & \mathbf{W}^V & = & \mathbf{V} \\ \begin{matrix} \square & \square & \square & \square \\ \square & \square & \square & \square \end{matrix} & \times & \begin{matrix} \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \end{matrix} & = & \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

$$\begin{aligned} & \text{softmax} \left( \frac{\begin{matrix} \mathbf{Q} & & \mathbf{K}^T \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} & \times & \begin{matrix} \square & \square \\ \square & \square \end{matrix} \end{matrix}}{\sqrt{d_k}} \right) \begin{matrix} \mathbf{V} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix} \\ & = \begin{matrix} \mathbf{Z} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix} \end{aligned}$$



# Plusieurs têtes !



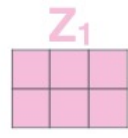
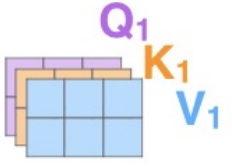
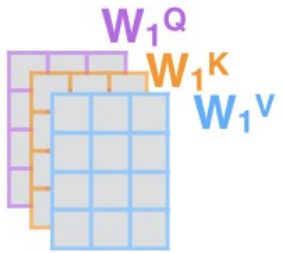
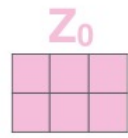
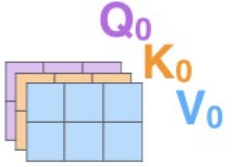
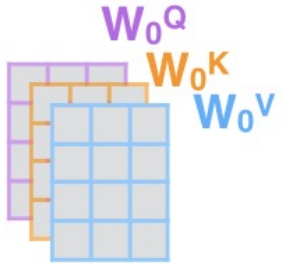
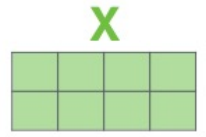
Deux interprétations possibles du « Multi-head Attention Mechanism » :

1. Identifier plusieurs types de relations.
2. Régulariser le mécanisme d'attention.

# On réconcilie tout le monde !

- 1) This is our input sentence\*
- 2) We embed each word\*
- 3) Split into 8 heads. We multiply  $X$  or  $R$  with weight matrices
- 4) Calculate attention using the resulting  $Q/K/V$  matrices
- 5) Concatenate the resulting  $Z$  matrices, then multiply with weight matrix  $W^o$  to produce the output of the layer

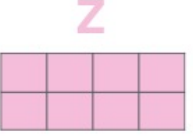
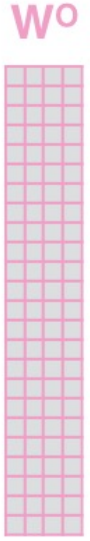
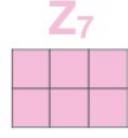
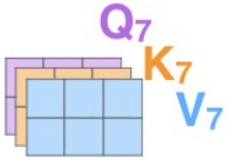
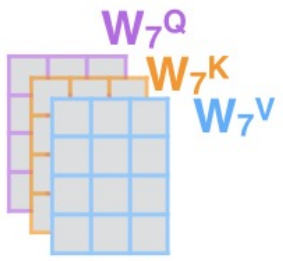
Thinking Machines



...

...

...



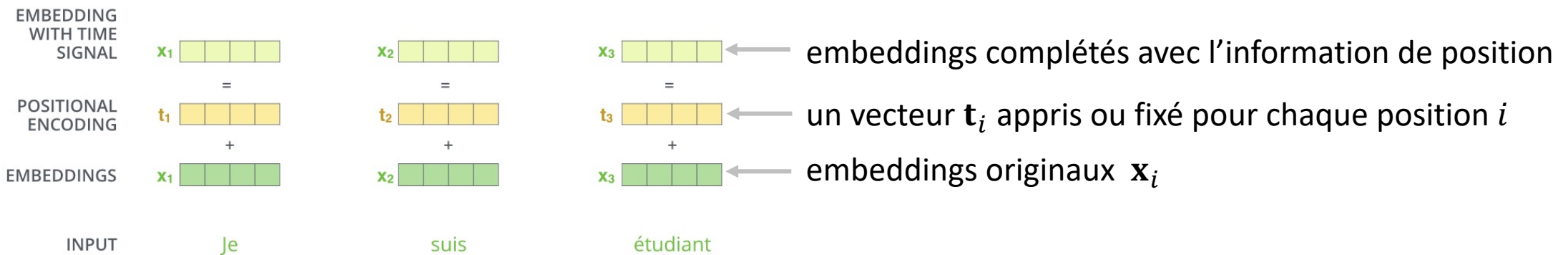
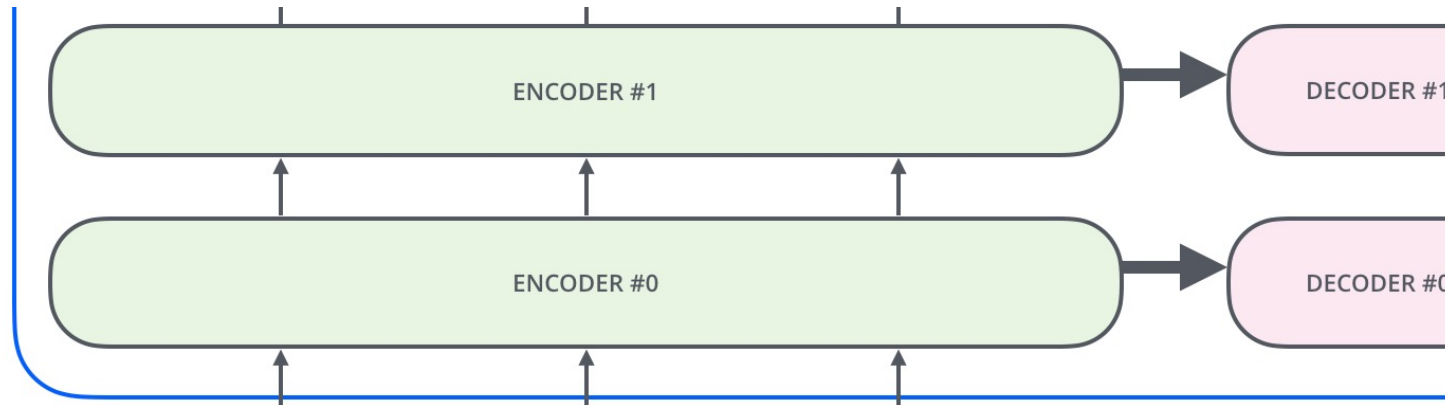
\* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



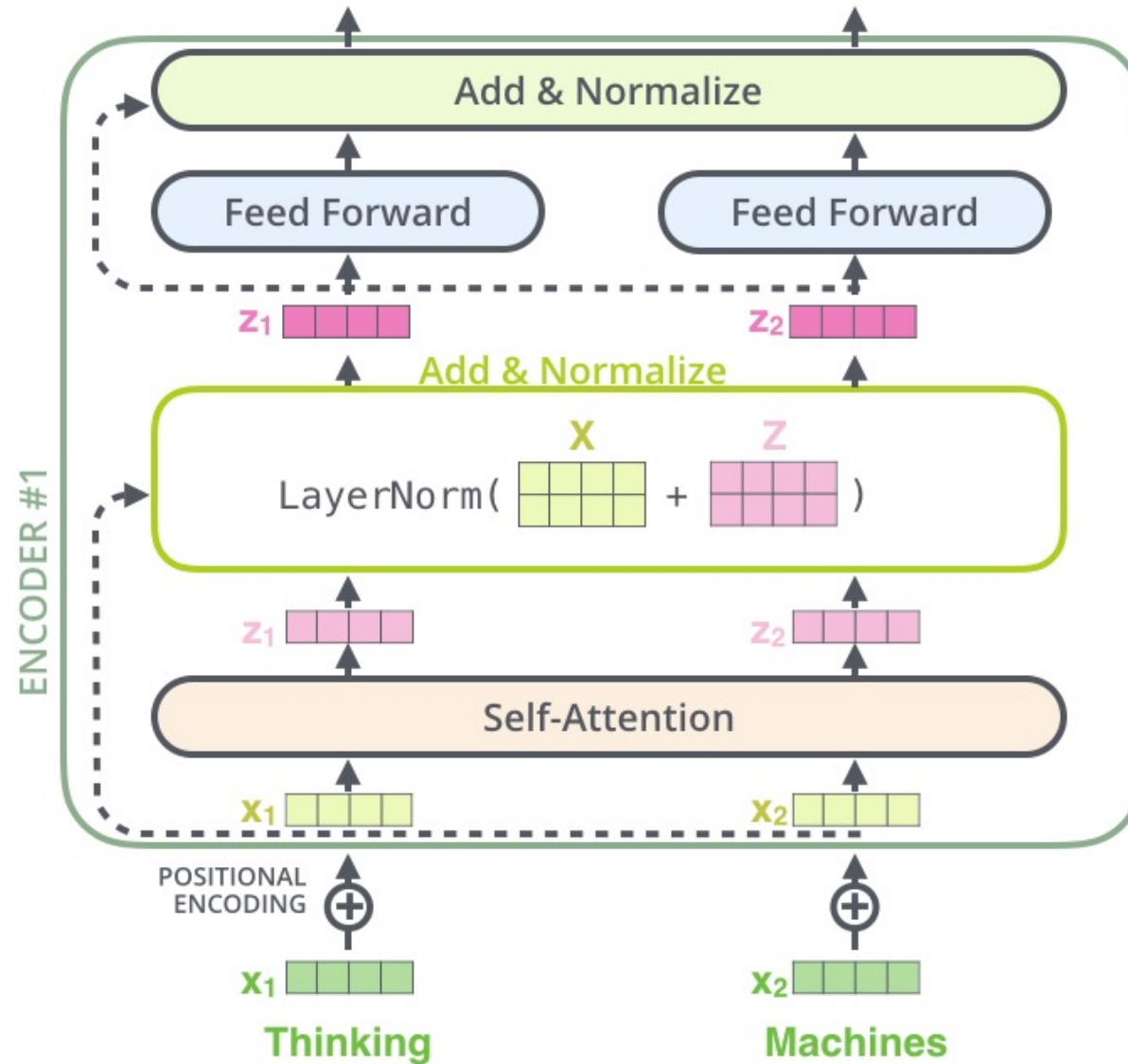
# N'oublions pas l'ordre des mots ! ~~bag of words~~

Jusque-là : le Transformer convertit un **ensemble** de vecteurs en un autre **ensemble** de vecteurs.

Ce qui est souhaité : convertir une **liste ordonnée** de vecteurs en une autre **liste ordonnée** de vecteurs.



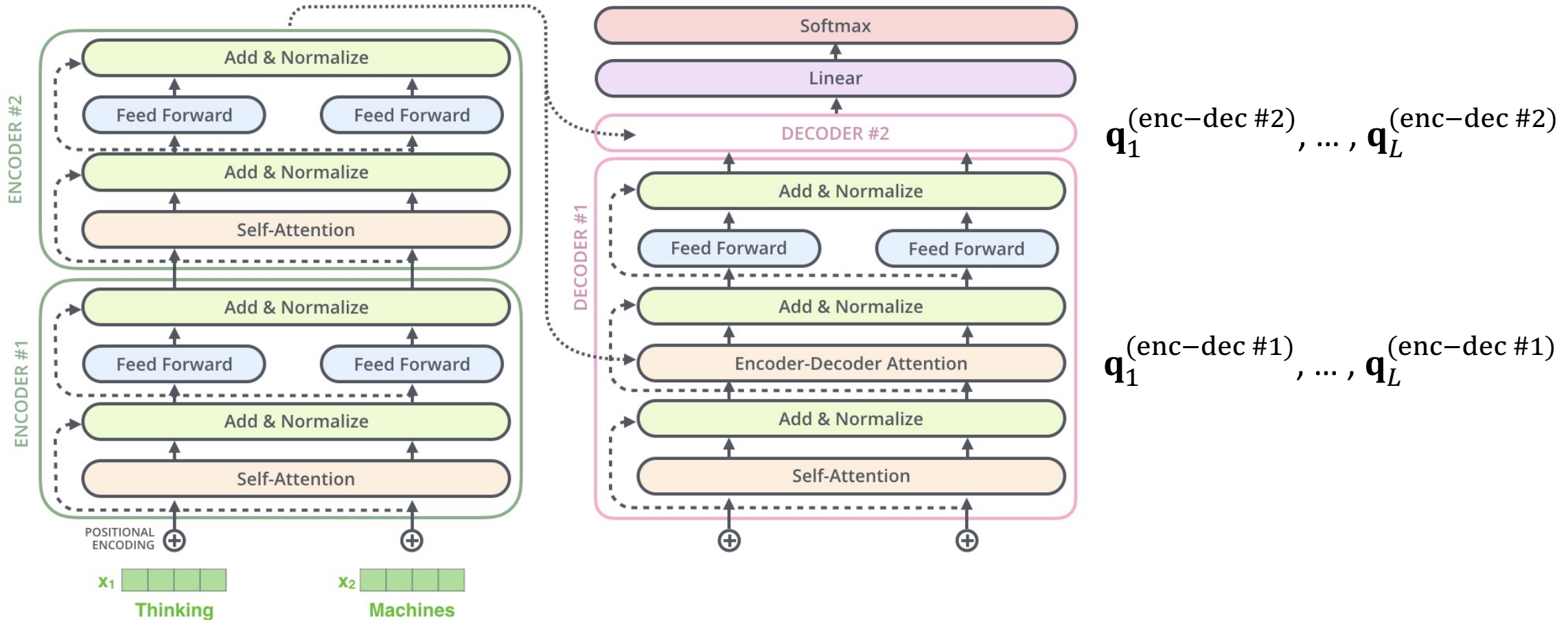
# Skip-connections, layer normalization et dropout



# Décodeur (1) : un mécanisme d'attention croisé

$\mathbf{k}_1^{(\text{enc-dec})}, \dots, \mathbf{k}_L^{(\text{enc-dec})}$

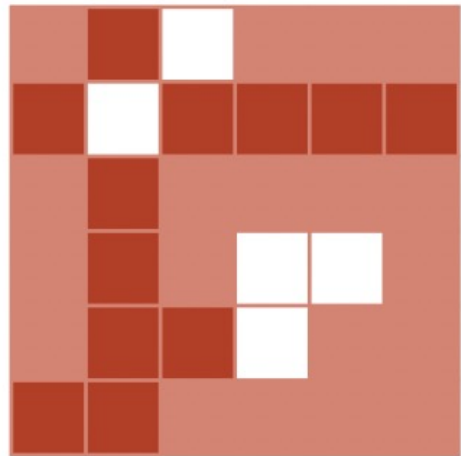
$\mathbf{v}_1^{(\text{enc-dec})}, \dots, \mathbf{v}_L^{(\text{enc-dec})}$



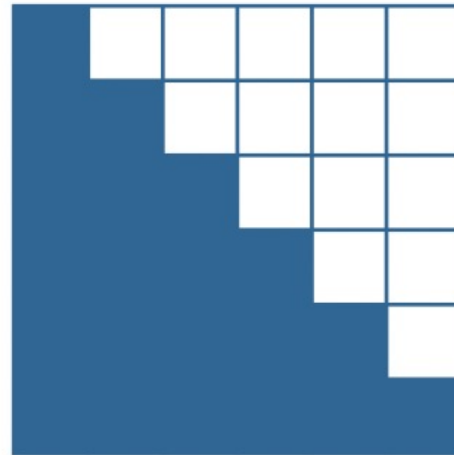


# Décodeur (2) : interdiction de regarder le futur !

$$\left( \mathbf{q}_i^{(\text{dec})} \cdot \mathbf{k}_j^{(\text{dec})} \right)_{i=1, \dots, L}^{j=1, \dots, L}$$

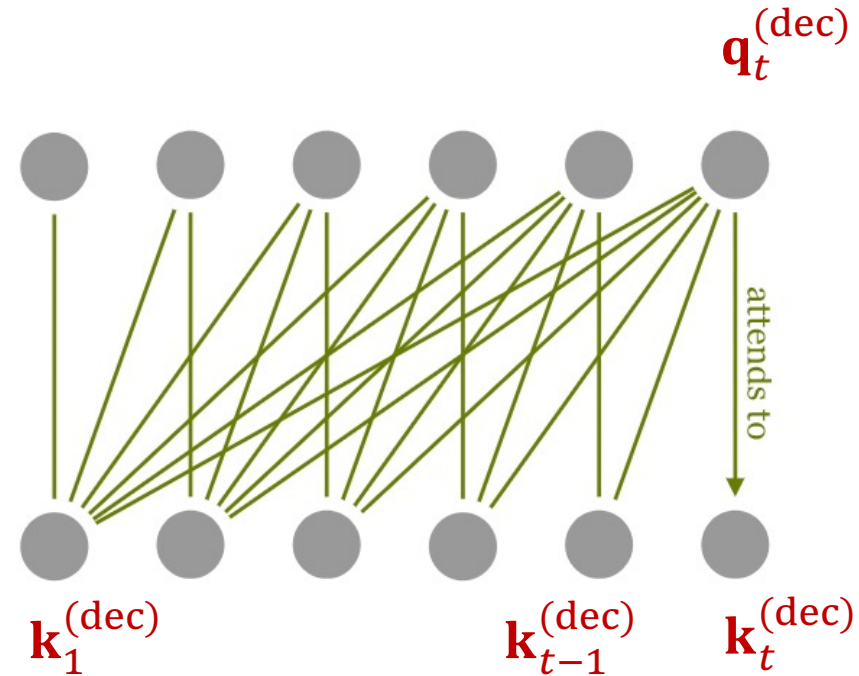


raw attention weights



mask

Un **masque** est appliqué avant le softmax du mécanisme de **self-attention** pour interdire au modèle de regarder le futur

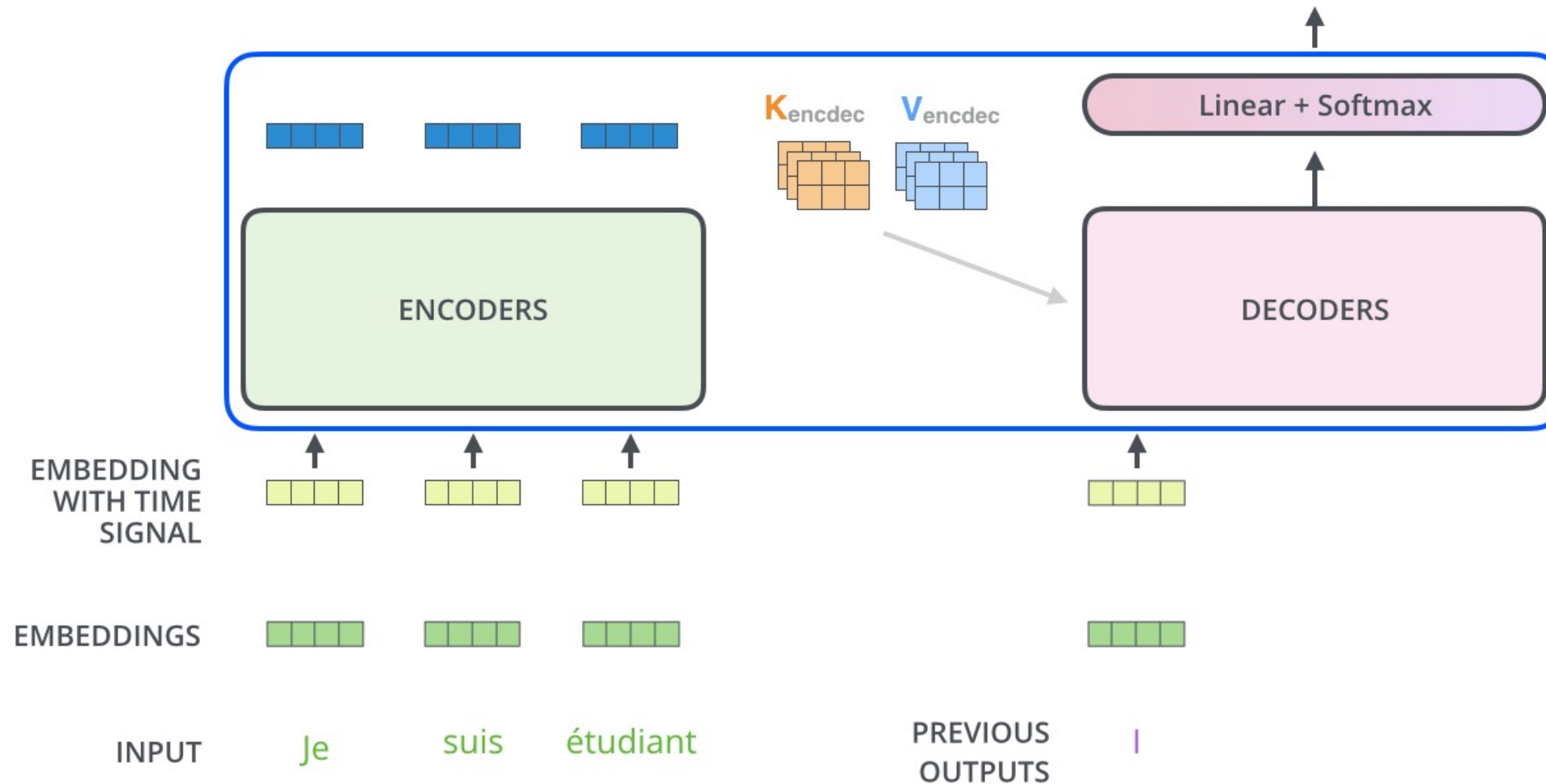


Lorsque le décodeur fait une prédiction à l'instant  $t$  il n'a accès qu'aux informations antérieures à cet instant. Lors de l'entraînement il faut donc lui **cacher le futur** de  $t$  bien qu'il soit connu dans le train set !

# Décodeur (3) : prédiction des mots

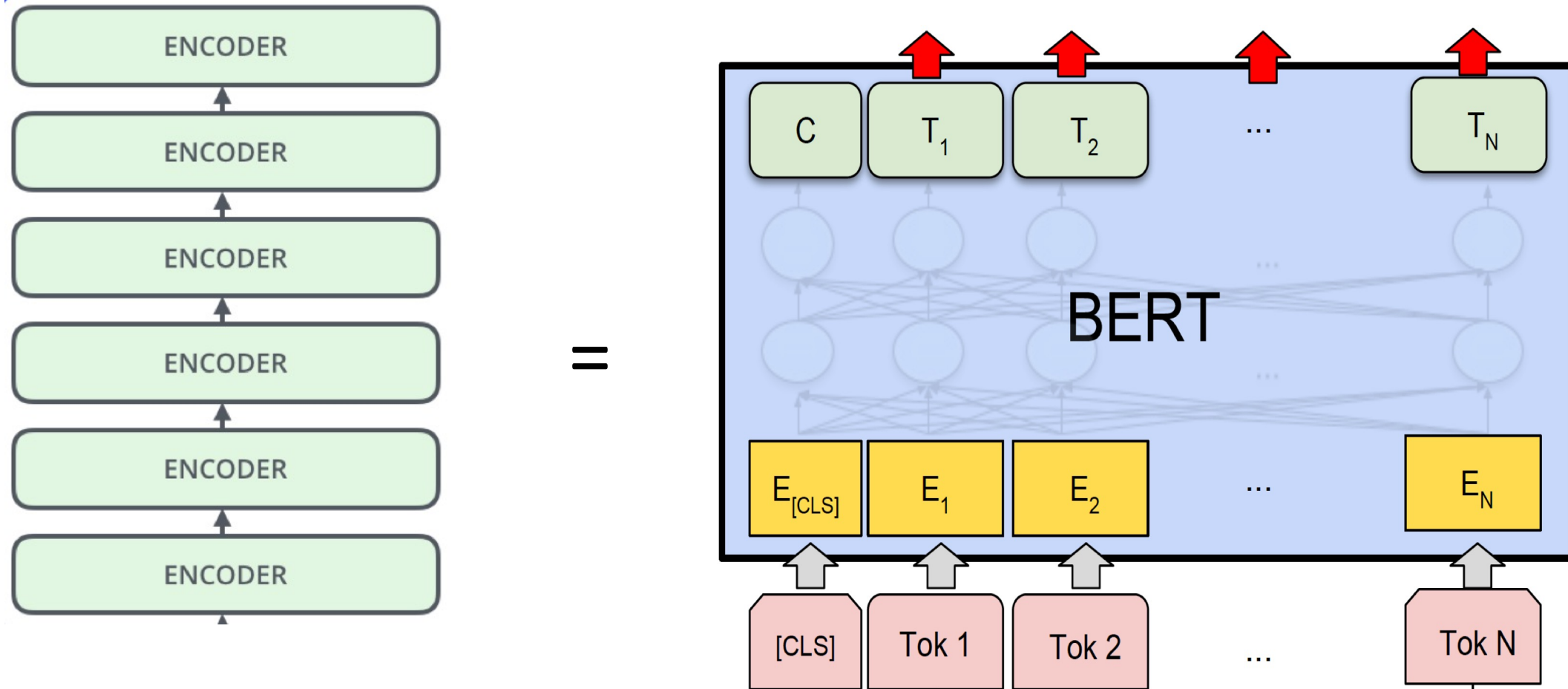
Decoding time step: 1 2 3 4 5 6

OUTPUT |

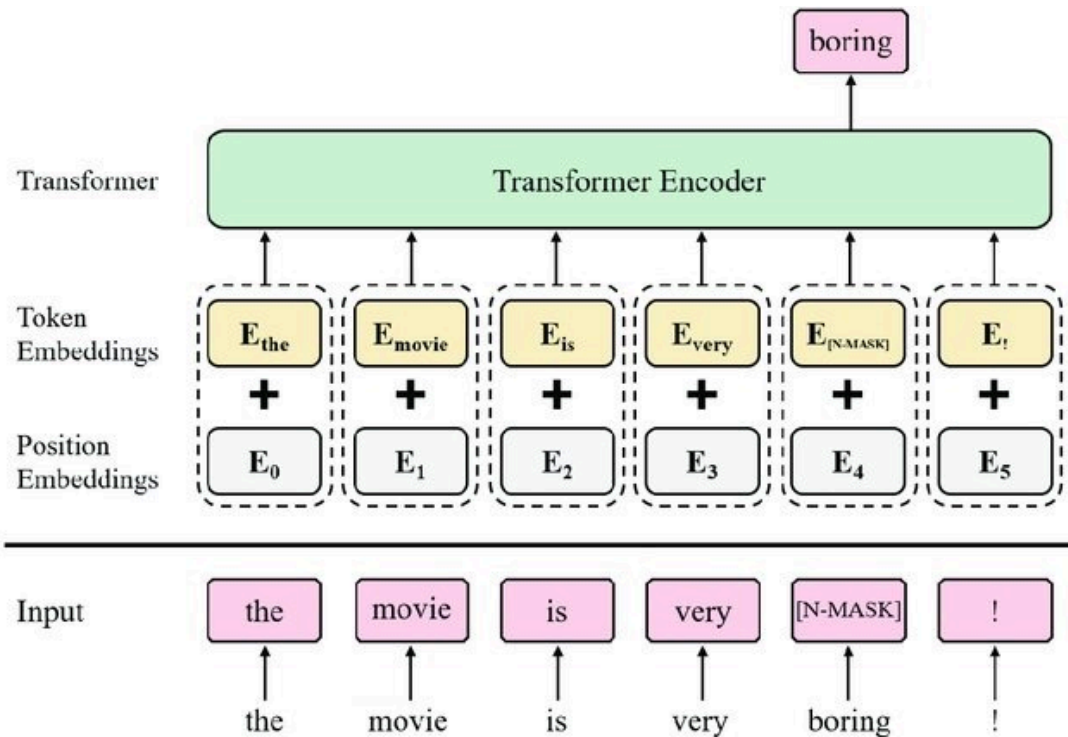




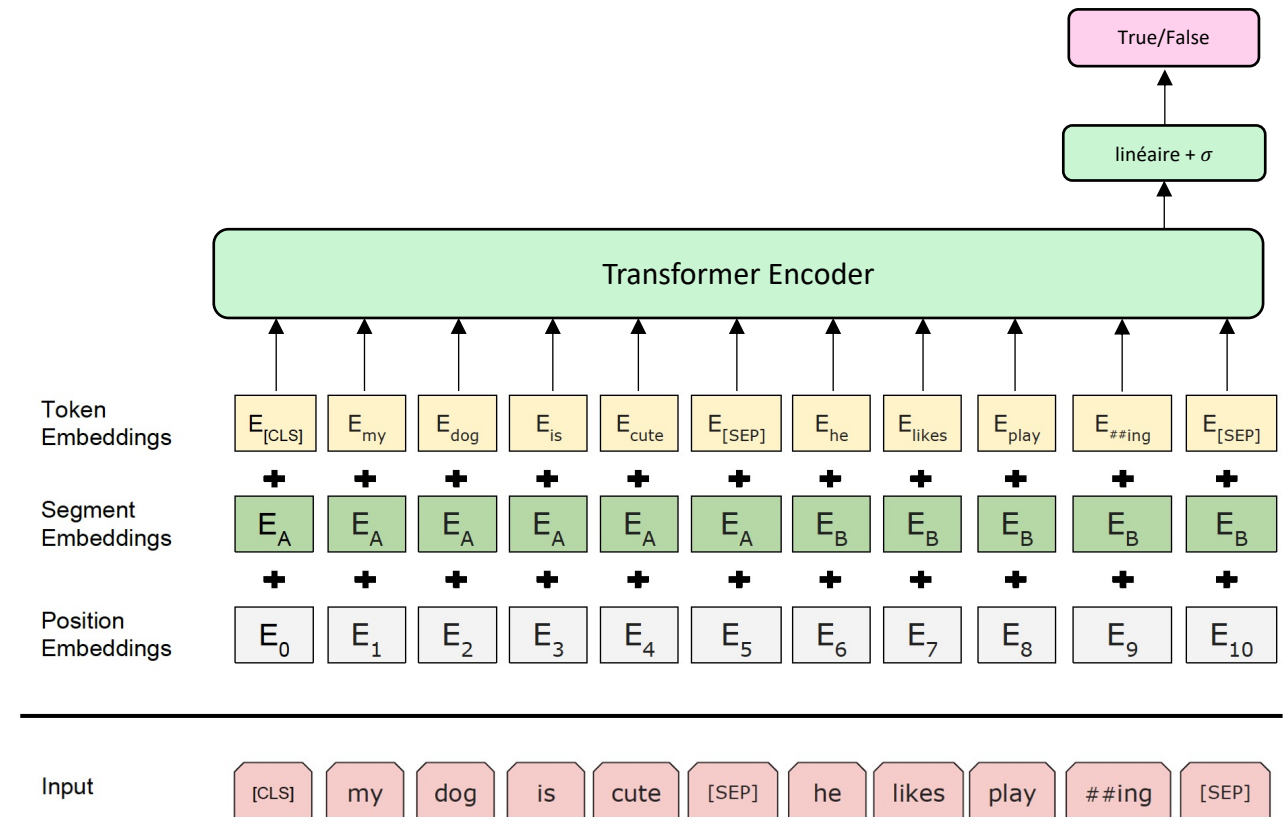
# BERT = Encoder Only Transformer



# BERT : 2 tâches de pré-entraînement

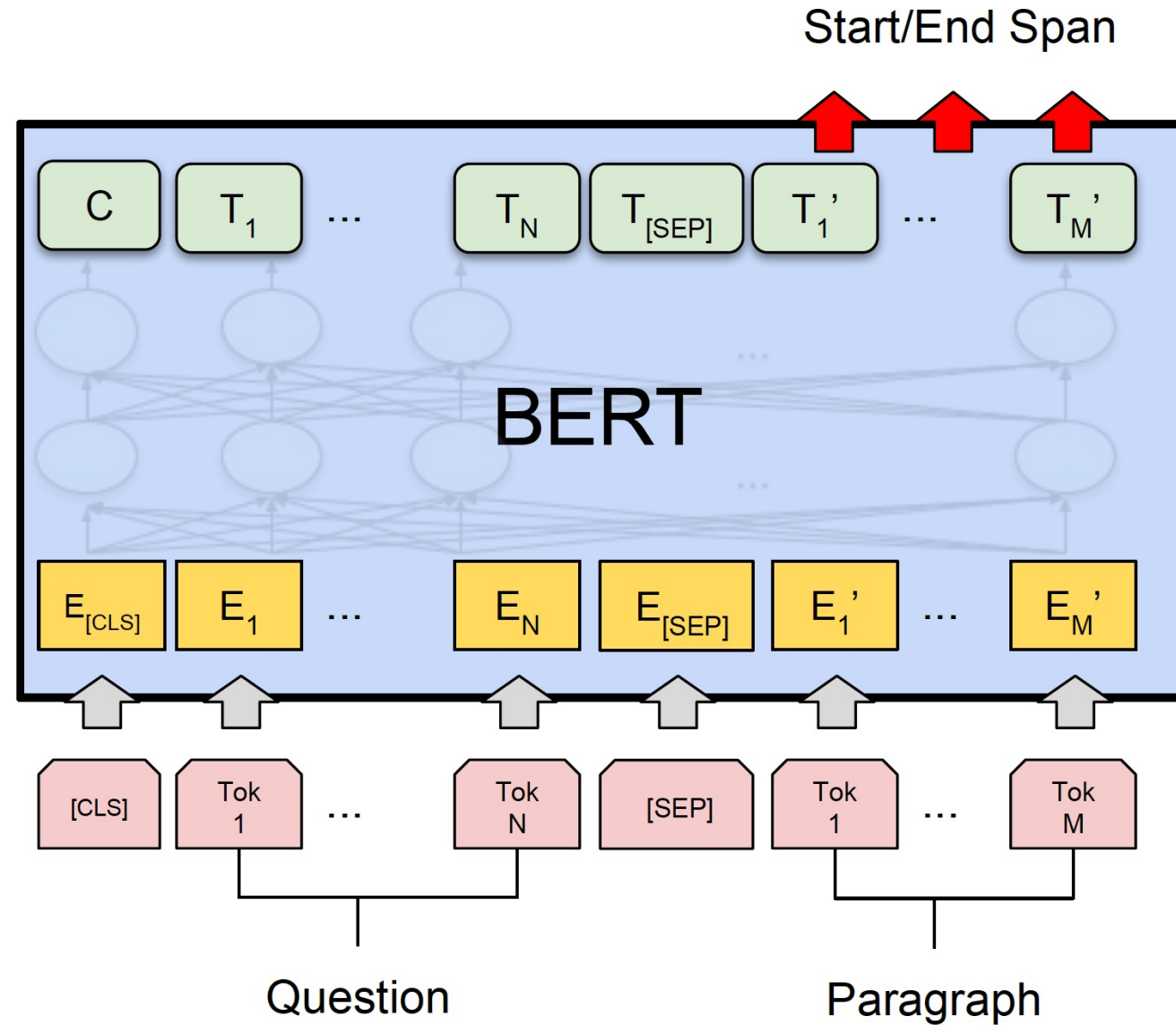


(1) Masked Language Model



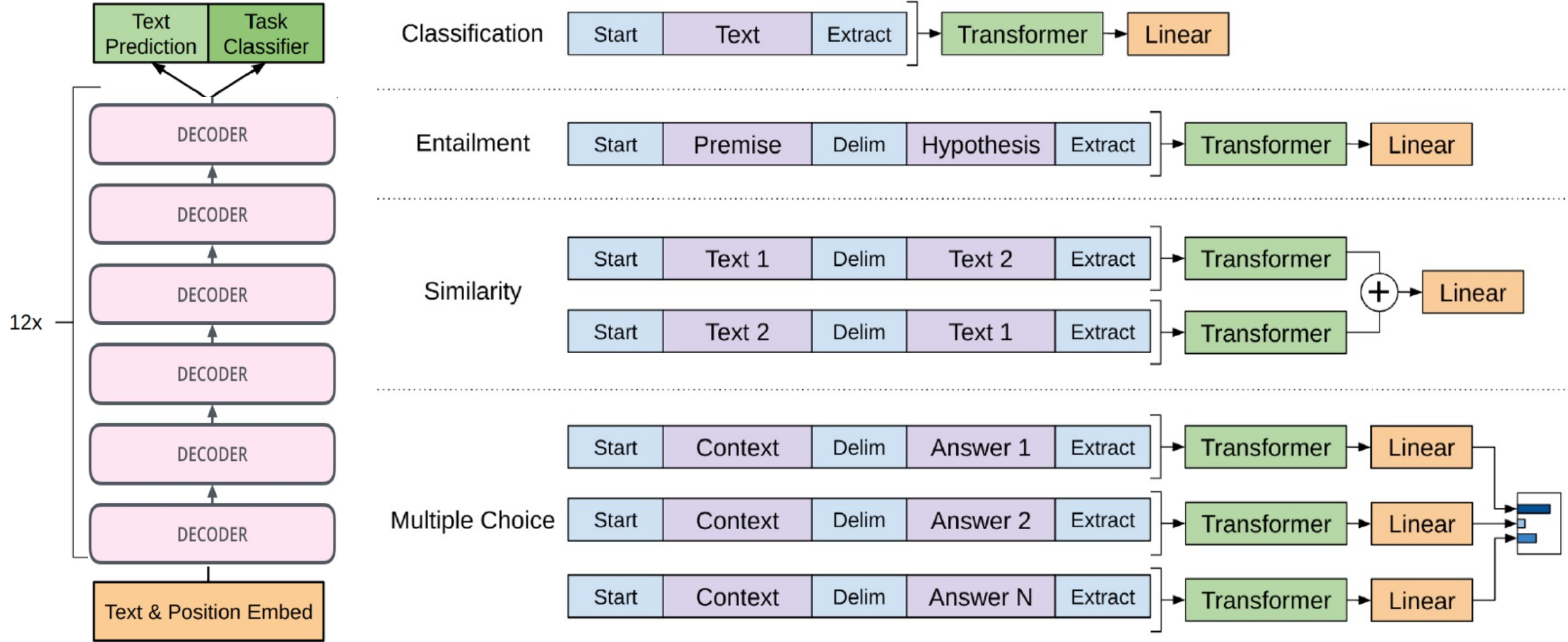
(2) Next Sentence Prediction

# BERT : architecture « universelle » pour le NLP



# GPT 1 = Decoder Only Transformer + ...

pré-entraînement sur LM + adaptation et fine tuning spécifique par tâche

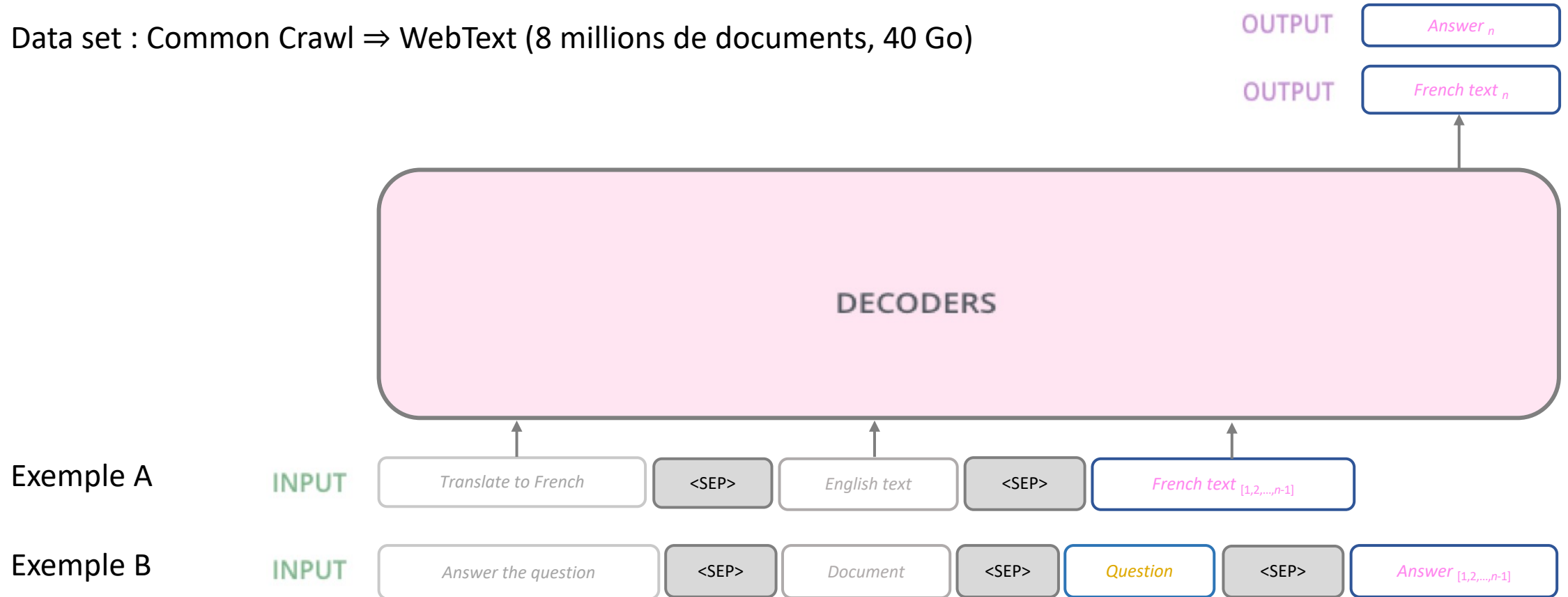


# GPT 2 & 3 = Decoder Only Transformer + ...

pré-entraînement sur très gros LM + ~~adaptation et fine tuning spécifique par tâche~~

Idée : spécification d'une tâche en langage naturel plutôt qu'en utilisant une architecture dédiée

Data set : Common Crawl  $\Rightarrow$  WebText (8 millions de documents, 40 Go)



...