

## TRANSFORMERS ET LLMs

Jérémie Cabessa

Laboratoire DAVID, UVSQ

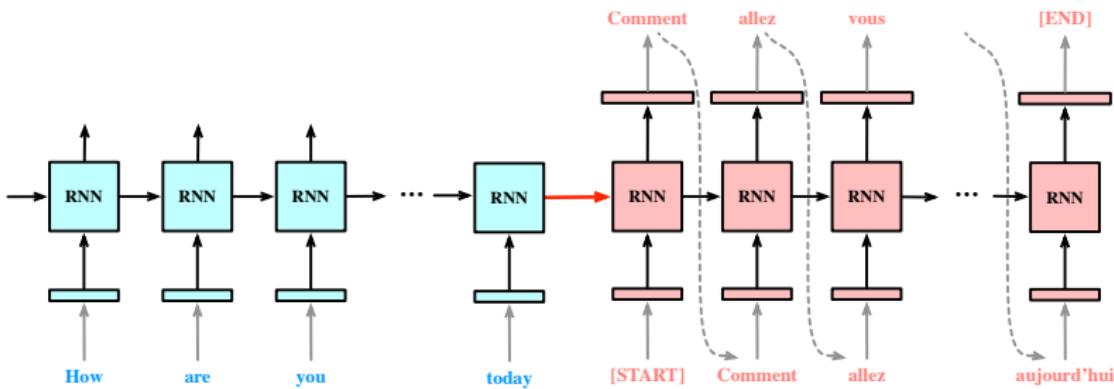
# ARCHITECTURE ENCODEUR-DÉCODEUR

- ▶ Les **transformers** sont les modèles révolutionnaires qui ont donné lieu à toute la famille des **large language models (LLMs)**.
  - ▶ Ils ont par la suite été généralisés à d'autres domaines que le NLP: vision, tabular, etc.

# ARCHITECTURE ENCODEUR-DÉCODEUR

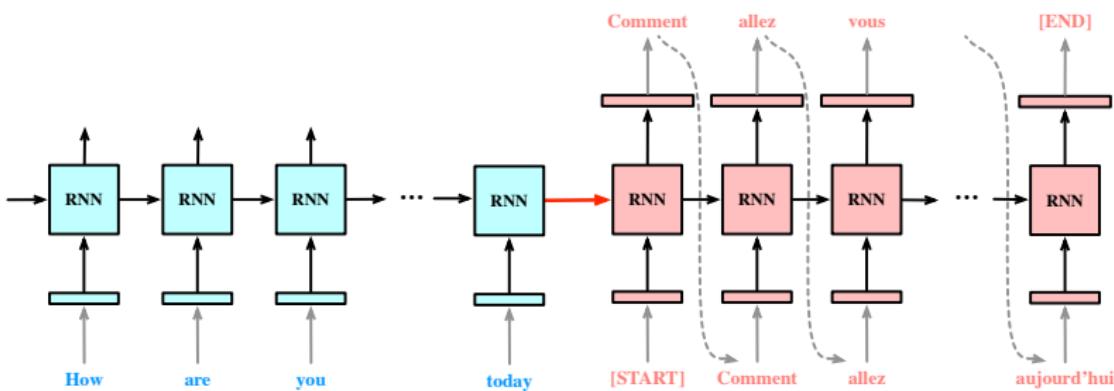
- ▶ Les **transformers** sont les modèles révolutionnaires qui ont donné lieu à toute la famille des **large language models (LLMs)**.
  - ▶ Ils ont par la suite été généralisés à d'autres domaines que le NLP: vision, tabular, etc.

# ARCHITECTURE ENCODEUR-DÉCODEUR



- ▶ Le dernier état de l'encodeur (flèche horizontale rouge) est le *context vector*.
- ▶ C'est un *embedding* qui encode tous les inputs et permet le décodage.

# ARCHITECTURE ENCODEUR-DÉCODEUR



- ▶ Le dernier état de l'encodeur (flèche horizontale rouge) est le *context vector*.
  - ▶ C'est un *embedding* qui encode tous les inputs et permet le décodage.

## ARCHITECTURE ENCODEUR-DÉCODEUR

- ▶ Les réseaux récurrents (RNNs) souffrent de la *disparition du gradient (vanishing gradient)*.
  - ▶ Les RNNs ne permettent pas la *parallélisation*: la forward pass doit être calculée de manière séquentielle, car les états successifs du réseaux sont donnés par l'historique des inputs.
  - ▶ Les RNNs peinent à capturer les dépendances de long termes entre les inputs (long-term dependencies).
  - ⇒ **Solution:** architecture *feedforward* (no vanishing/exploding gradients, parallelizable) couplée à un mécanisme d' *attention* (short and long term dependencies).

## ARCHITECTURE ENCODEUR-DÉCODEUR

- ▶ Les réseaux récurrents (RNNs) souffrent de la *disparition du gradient (vanishing gradient)*.
  - ▶ Les RNNs ne permettent pas la *parallélisation*: la forward pass doit être calculée de manière séquentielle, car les états successifs du réseaux sont donnés par l'historique des inputs.
  - ▶ Les RNNs peinent à capturer les dépendances de long termes entre les inputs (long-term dependencies).

⇒ **Solution:** architecture *feedforward* (no vanishing/exploding gradients, parallelizable) couplée à un mécanisme d'*attention* (short and long term dependencies).

## ARCHITECTURE ENCODEUR-DÉCODEUR

- ▶ Les réseaux récurrents (RNNs) souffrent de la *disparition du gradient (vanishing gradient)*.
  - ▶ Les RNNs ne permettent pas la *parallelisation*: la forward pass doit être calculée de manière séquentielle, car les états successifs du réseaux sont donnés par l'historique des inputs.
  - ▶ Les RNNs peinent à capturer les dépendances de long termes entre les inputs (*long-term dependencies*).

⇒ **Solution:** architecture *feedforward* (no vanishing/exploding gradients, parallelizable) couplée à un mécanisme d' *attention* (short and long term dependencies).

## ARCHITECTURE ENCODEUR-DÉCODEUR

- ▶ Les réseaux récurrents (RNNs) souffrent de la *disparition du gradient (vanishing gradient)*.
  - ▶ Les RNNs ne permettent pas la *parallélisation*: la forward pass doit être calculée de manière séquentielle, car les états successifs du réseaux sont donnés par l'historique des inputs.
  - ▶ Les RNNs peinent à capturer les dépendances de long termes entre les inputs (long-term dependencies).
  - ⇒ **Solution:** architecture *feedforward* (no vanishing/exploding gradients, parallelizable) couplée à un mécanisme d'*attention* (short and long term dependencies).

## TRANSFORMER

- ▶ Papier original: [\[Vaswani et al., 2017\]](#)
  - ▶ Architecture **encodeur-décodeur**.
  - ▶ Réseau de neurones **feedforward** (no vanishing/exploding gradients, parallelizable)
  - ▶ Mécanismes de **self-attention** et d'attention dans l'encodeur et le décodeur (short- and long-term dependencies).

# TRANSFORMER

- ▶ Papier original: [\[Vaswani et al., 2017\]](#)
- ▶ Architecture **encodeur-décodeur**.
- ▶ Réseau de neurones feedforward (no vanishing/exploding gradients, parallelizable)
- ▶ Mécanismes de **self-attention** et d'attention dans l'encodeur et le décodeur (short- and long-term dependencies).

# TRANSFORMER

- ▶ Papier original: [[Vaswani et al., 2017](#)]
- ▶ Architecture **encodeur-décodeur**.
- ▶ Réseau de neurones **feedforward** (no vanishing/exploding gradients, parallelizable)
- ▶ Mécanismes de **self-attention** et d'attention dans l'encodeur et le décodeur (short- and long-term dependencies).

## TRANSFORMER

- ▶ Papier original: [\[Vaswani et al., 2017\]](#)
  - ▶ Architecture **encodeur-décodeur**.
  - ▶ Réseau de neurones **feedforward** (no vanishing/exploding gradients, parallelizable)
  - ▶ Mécanismes de **self-attention** et d'**attention** dans l'encodeur et le décodeur (short- and long-term dependencies).

# TRANSFORMER

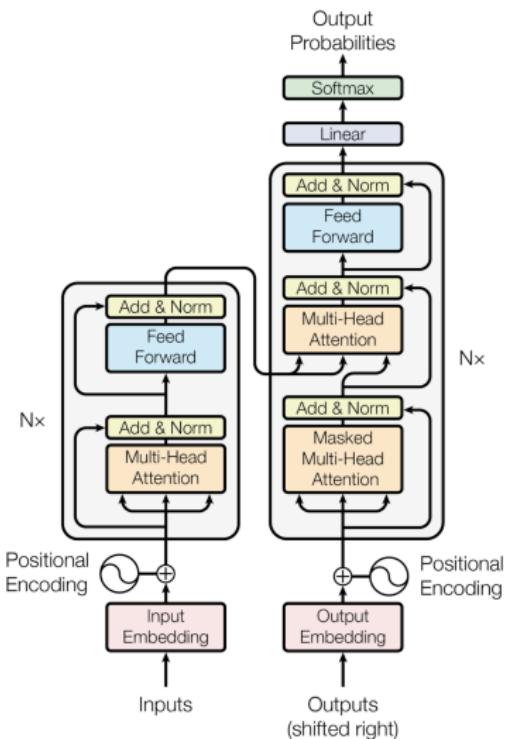


Figure taken from [Vaswani et al., 2017].

## TRANSFORMER

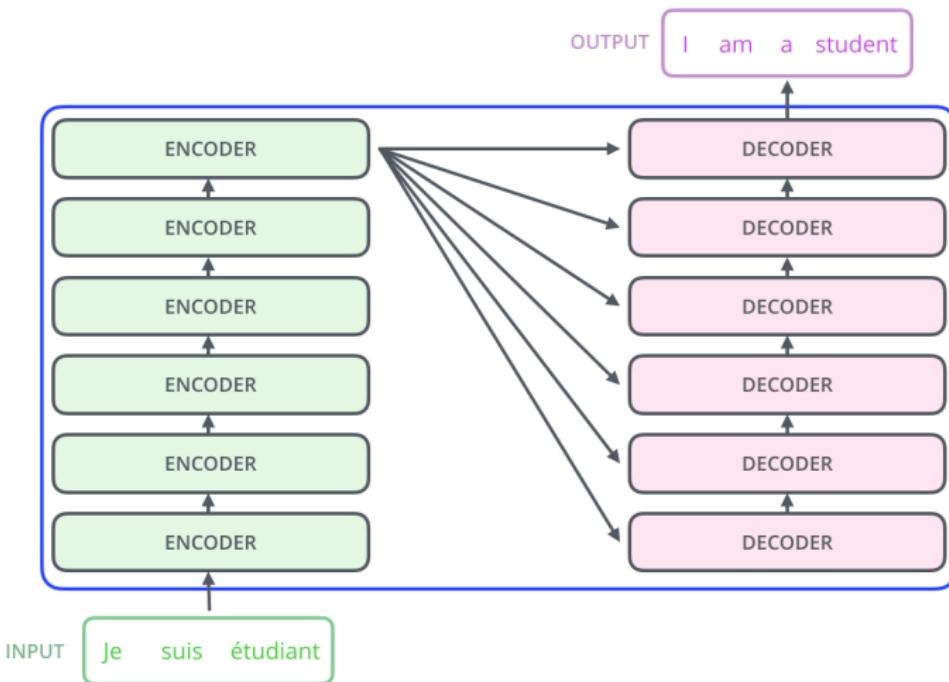


Figure taken from [\[Alammar, 2018\]](#).

## TRANSFORMER

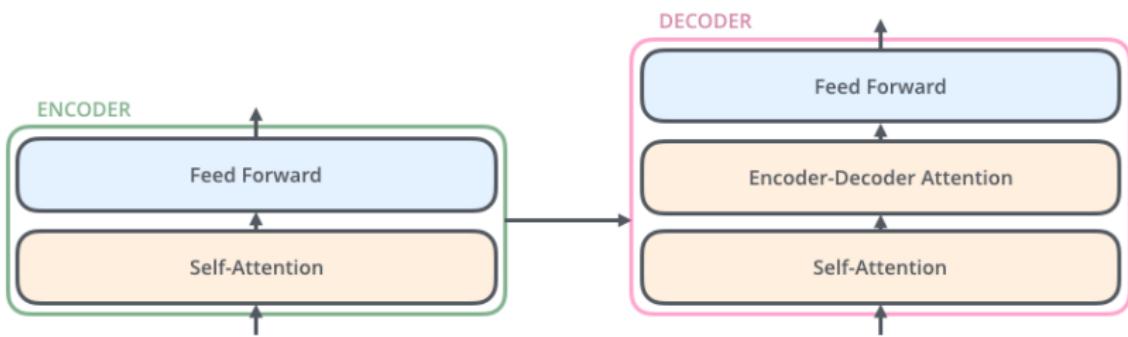


Figure taken from [\[Alammar, 2018\]](#).

## TRANSFORMER

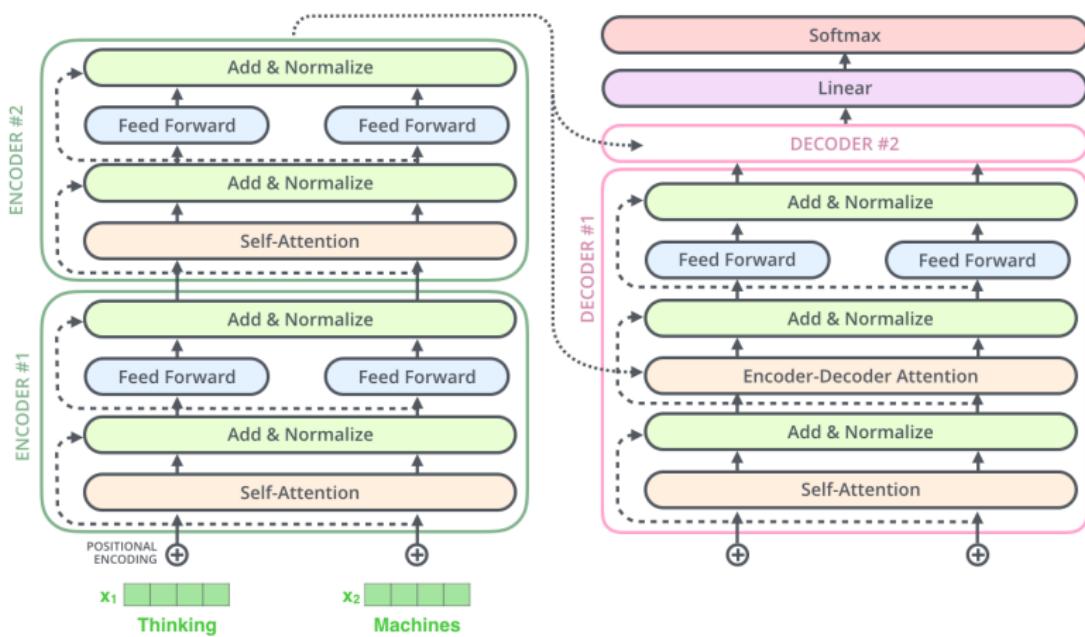


Figure taken from [Alammar, 2018](#).

## TOKENIZATION ET EMBEDDING

- ▶ Le texte passé en `input` est "tokenisé", i.e., convertit en une séquence de "input tokens".
  - ▶ Chaque input token est associé à un entier appelé "token id", qui est à un indice dans un vocabulaire d'environ 40K mots.
  - ▶ Chaque token id est ensuite "embeddé" en un vecteur de taille 512 grâce à un embedding statique sous forme de "lookup matrice" de dimension  $40K \times 512$ .
  - ▶ À chaque "input embedding" est ajouté un vecteur de "positional encoding" qui capture la position de cet input dans le texte de départ.

## TOKENIZATION ET EMBEDDING

- ▶ Le text passé en input est “tokenisé”, i.e., convertit en une séquence de “input tokens”.
  - ▶ Chaque input token est associé à un entier appelé “token id”, qui est à un indice dans un vocabulaire d'environ 40K mots.
  - ▶ Chaque token id est ensuite “embeddé” en un vecteur de taille 512 grâce à un embedding statique sous forme de “lookup matrice” de dimension  $40K \times 512$ .
  - ▶ À chaque “input embedding” est ajouté un vecteur de “positional encoding” qui capture la position de cet input dans le texte de départ.

## TOKENIZATION ET EMBEDDING

- ▶ Le texte passé en input est “tokenisé”, i.e., convertit en une séquence de “input tokens”.
  - ▶ Chaque input token est associé à un entier appelé “token id”, qui est à un indice dans un vocabulaire d'environ 40K mots.
  - ▶ Chaque token id est ensuite “embeddé” en un vecteur de taille 512 grâce à un embedding statique sous forme de “lookup matrice” de dimension  $40K \times 512$ .
  - ▶ À chaque “input embedding” est ajouté un vecteur de “positional encoding” qui capture la position de cet input dans le texte de départ.

## TOKENIZATION ET EMBEDDING

- ▶ Le texte passé en input est “tokenisé”, i.e., convertit en une séquence de “input tokens”.
  - ▶ Chaque input token est associé à un entier appelé “token id”, qui est à un indice dans un vocabulaire d'environ 40K mots.
  - ▶ Chaque token id est ensuite “embeddé” en un vecteur de taille 512 grâce à un embedding statique sous forme de “lookup matrice” de dimension  $40K \times 512$ .
  - ▶ À chaque “input embedding” est ajouté un vecteur de “positional encoding” qui capture la position de cet input dans le texte de départ.

## TOKENIZATION

### # Original Sentence

## TOKENIZATION

### # Original Sentence

Let's learn deep learning!

## TOKENIZATION

### # Original Sentence

Let's learn deep learning!

## # Tokenized Sentence

## TOKENIZATION

### # Original Sentence

Let's learn deep learning!

## # Tokenized Sentence

```
["Let", "", "s", "learn", "deep", "learning", "!"]
```

## TOKENIZATION

### # Original Sentence

Let's learn deep learning!

## # Tokenized Sentence

```
["Let", "", "s", "learn", "deep", "learning", "!"]
```

# Adding [CLS] and [SEP] Tokens

## TOKENIZATION

### # Original Sentence

Let's learn deep learning!

## # Tokenized Sentence

```
["Let", "", "s", "learn", "deep", "learning", "!"]
```

# Adding [CLS] and [SEP] Tokens

["[CLS]", "Let", "", "s", "learn", "deep", "learning", "!", "[SEP]"]

## TOKENIZATION

### # Original Sentence

Let's learn deep learning!

## # Tokenized Sentence

```
["Let", "", "s", "learn", "deep", "learning", "!"]
```

# Adding [CLS] and [SEP] Tokens

["[CLS]", "Let", "'", "s", "learn", "deep", "learning", "!", "[SEP]"]

# Padding

## TOKENIZATION

```
# Original Sentence
Let's learn deep learning!

# Tokenized Sentence
["Let", "'", "s", "learn", "deep", "learning", "!"]

# Adding [CLS] and [SEP] Tokens
["[CLS]", "Let", "'", "s", "learn", "deep", "learning", "!", "[SEP]"]

# Padding
["[CLS]", "Let", "'", "s", "learn", "deep", "learning", "!", "[SEP]",
 "[PAD]", "[PAD]", "[PAD]", "[PAD]", "[PAD]", "[PAD]", "[PAD]"]

# Converting to token ids
[101, 2421, 112, 188, 3858, 1996, 3776, 106, 102, 0, 0, 0, 0, 0, 0, 0, 0]
```

## TOKENIZATION

## TOKENIZATION

```
# Original Sentence
Let's learn deep learning!

# Tokenized Sentence
["Let", "'", "s", "learn", "deep", "learning", "!"]

# Adding [CLS] and [SEP] Tokens
["[CLS]", "Let", "'", "s", "learn", "deep", "learning", "!", "[SEP]"]

# Padding
["[CLS]", "Let", "'", "s", "learn", "deep", "learning", "!", "[SEP]",
 "[PAD]", "[PAD]", "[PAD]", "[PAD]", "[PAD]", "[PAD]", "[PAD]"]

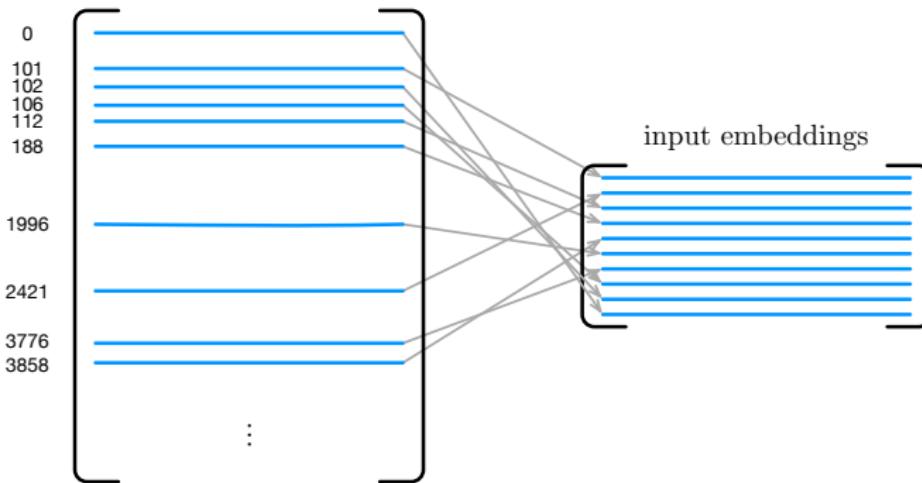
# Converting to token ids
[101, 2421, 112, 188, 3858, 1996, 3776, 106, 102, 0, 0, 0, 0, 0, 0, 0]
```

## EMBEDDING

input token ids

[101, 2421, 112, 188, 3858, 1996, 3776, 106, 102, 0]

embedding matrix

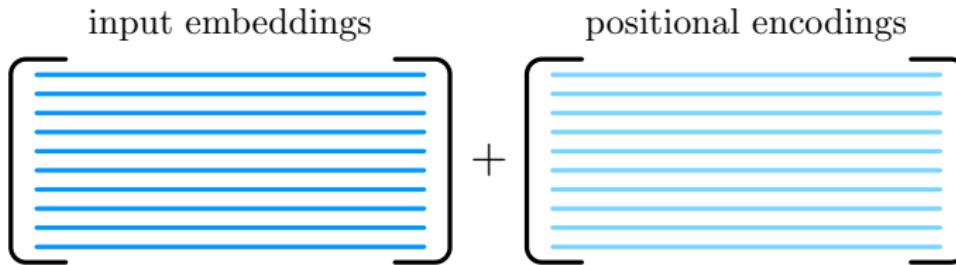


## POSITIONAL ENCODING

$$PE(pos, 2i) = \sin\left(pos/10000^{\frac{2i}{d}}\right)$$

$$PE(pos, 2i + 1) = \cos\left(pos/10000^{\frac{2i}{d}}\right)$$

où  $pos$  est la position du input token ( $0, 1, 2, \dots$ ),  $i$  se réfère à la dimension de l'encoding ( $0 \leq i \leq 255$ ), et  $d = 512$ .



## POSITIONAL ENCODING

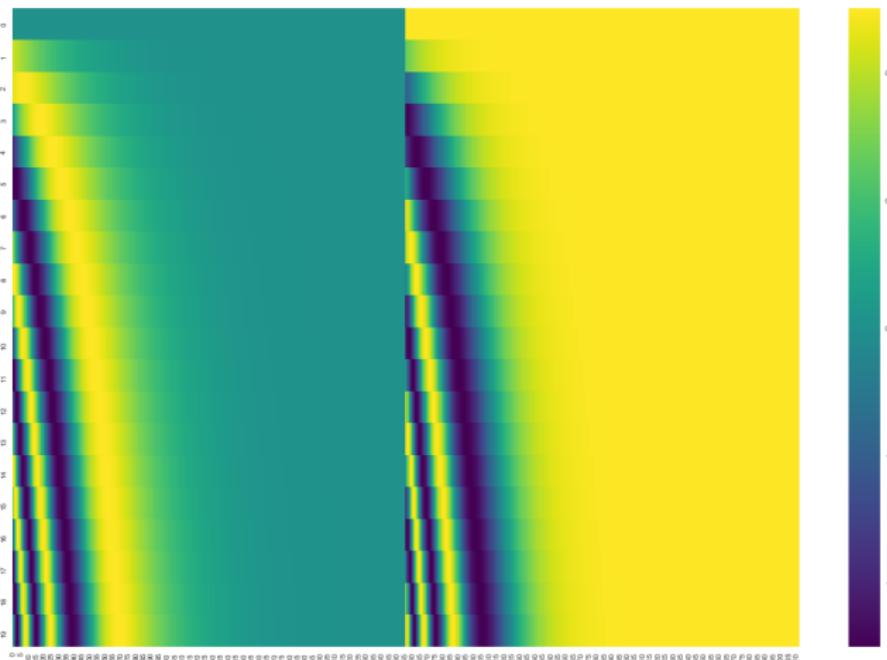


Figure taken from [Alammar, 2018]. Exemple de "positional encoding" légèrement différent de celui présenté au slide précédent: les fonctions  $\cos$  et  $\sin$  ne s'alternent pas en fonction de l'indice  $i$ ...

## ATTENTION

- ▶ On distingue 2 types d'attention:
    - ▶ l'attention classique
    - ▶ la self-attention

## ATTENTION

- ▶ On distingue 2 types d'attention:
    - ▶ **l'attention classique**
    - ▶ **la self-attention**

## ATTENTION

- ▶ On distingue 2 types d'attention:
    - ▶ **l'attention classique**
    - ▶ **la self-attention**

## ATTENTION

- ▶ Le **mécanisme d'attention** crée des embeddings de mots (tokens) qui prennent en compte d'autres mots de la phrase.
  - ▶ L'attention joue le rôle de *mémoire*.
  - ▶ Exemple: il serait bon que l'embedding du mot "bateau" dans la phrase "le bateau bleu est amarré dans le port d'Amsterdam" mette de l'attention sur les mots "bleu" et "amarré", puisque ces derniers apportent des précisions sur le bateau.
  - ▶ En pratique, l'embedding d'un mot correspondra (en gros) à une combinaison pondérée des embeddings des autres mots.

## ATTENTION

- ▶ Le **mécanisme d'attention** crée des embeddings de mots (tokens) qui prennent en compte d'autres mots de la phrase.
  - ▶ L'attention joue le rôle de *mémoire*.
  - ▶ Exemple: il serait bon que l'embedding du mot "bateau" dans la phrase "le bateau bleu est amarré dans le port d'Amsterdam" mette de l'attention sur les mots "bleu" et "amarré", puisque ces derniers apportent des précisions sur le bateau.
  - ▶ En pratique, l'embedding d'un mot correspondra (en gros) à une combinaison pondérée des embeddings des autres mots.

## ATTENTION

- ▶ Le **mécanisme d'attention** crée des embeddings de mots (tokens) qui prennent en compte d'autres mots de la phrase.
  - ▶ L'attention joue le rôle de *mémoire*.
  - ▶ Exemple: il serait bon que l'embedding du mot "bateau" dans la phrase "le bateau bleu est amarré dans le port d'Amsterdam" mette de l'attention sur les mots "bleu" et "amarré", puisque ces derniers apportent des précisions sur le bateau.
  - ▶ En pratique, l'embedding d'un mot correspondra (en gros) à une combinaison pondérée des embeddings des autres mots.

## ATTENTION

- ▶ Le **mécanisme d'attention** crée des embeddings de mots (tokens) qui prennent en compte d'autres mots de la phrase.
  - ▶ L'attention joue le rôle de *mémoire*.
  - ▶ Exemple: il serait bon que l'embedding du mot "bateau" dans la phrase "le bateau bleu est amarré dans le port d'Amsterdam" mette de l'attention sur les mots "bleu" et "amarré", puisque ces derniers apportent des précisions sur le bateau.
  - ▶ En pratique, l'embedding d'un mot correspondra (en gros) à une combinaison pondérée des embeddings des autres mots.

## ATTENTION

- ▶ Analogie avec un *système de recherche (retrieval system)* basé sur des *requêtes*, des *clés* et des *valeurs (queries, keys, and values)*.
  - ▶ On a une base de données d'éléments représentés par des couples clé-valeur. Pour une requête donnée, on cherche les clés qui possèdent le plus de similarités avec cette dernière et on extrait les valeurs correspondantes.
  - ▶ A “query” matches different “keys” and retrieves the corresponding “values”.

## ATTENTION

- ▶ Analogie avec un *système de recherche (retrieval system)* basé sur des *requêtes*, des *clés* et des *valeurs (queries, keys, and values)*.
  - ▶ On a une base de données d'éléments représentés par des couples clé-valeur. Pour une requête donnée, on cherche les clés qui possèdent le plus de similarités avec cette dernière et on extrait les valeurs correspondantes.
  - ▶ *A “query” matches different “keys” and retrieves the corresponding “values”.*

## ATTENTION

- ▶ Analogie avec un *système de recherche (retrieval system)* basé sur des *requêtes*, des *clés* et des *valeurs (queries, keys, and values)*.
  - ▶ On a une base de données d'éléments représentés par des couples clé-valeur. Pour une requête donnée, on cherche les clés qui possèdent le plus de similarités avec cette dernière et on extrait les valeurs correspondantes.
  - ▶ A “query” matches different “keys” and retrieves the corresponding “values”.

## ATTENTION

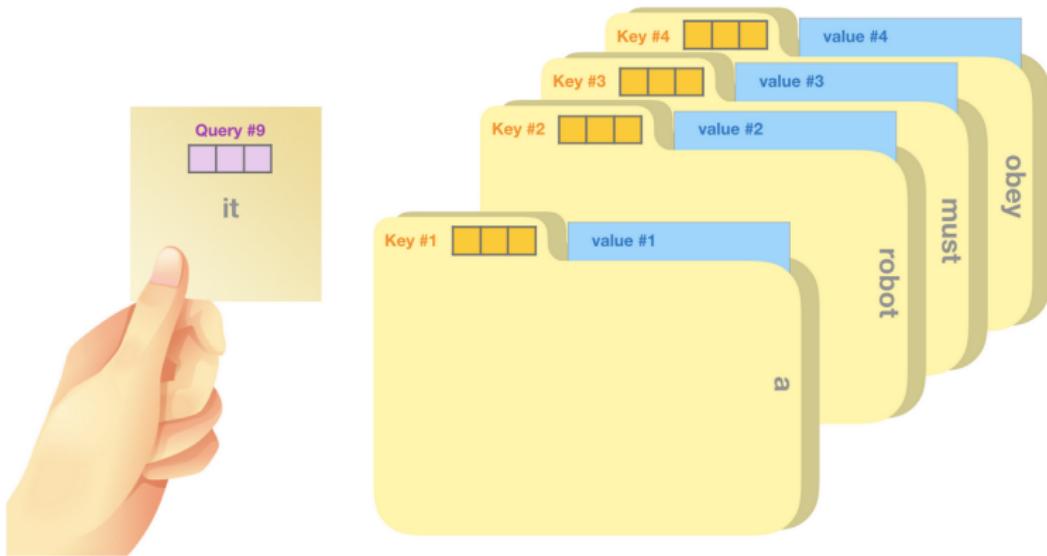


Figure taken from [\[Alammar, 2018\]](#).

## ATTENTION



Figure taken from [\[Alammar, 2018\]](#).

## SCALED DOT-PRODUCT ATTENTION

- ▶ Pour la self-attention qui se déroule dans l'encodeur, les requêtes, les clés et les valeurs proviennent de la séquence des inputs.
  - ▶ Pour la self-attention qui se déroule dans le décodeur, les requêtes, les clés et les valeurs proviennent de la séquence des mots décodés jusqu'à maintenant.
  - ▶ Pour l'attention qui se déroule dans le décodeur, les requêtes proviennent des outputs de l'encodeur, et les clés et les valeurs proviennent de la séquence des mots décodés jusqu'à maintenant.

## SCALED DOT-PRODUCT ATTENTION

- ▶ Pour la self-attention qui se déroule dans l'encodeur, les requêtes, les clés et les valeurs proviennent de la séquence des inputs.
  - ▶ Pour la self-attention qui se déroule dans le décodeur, les requêtes, les clés et les valeurs proviennent de la séquence des mots décodés jusqu'à maintenant.
  - ▶ Pour l'attention qui se déroule dans le décodeur, les requêtes proviennent des outputs de l'encodeur, et les clés et les valeurs proviennent de la séquence des mots décodés jusqu'à maintenant.

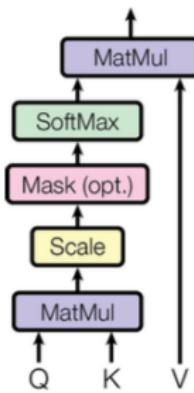
## SCALED DOT-PRODUCT ATTENTION

- ▶ Pour la self-attention qui se déroule dans l'encodeur, les requêtes, les clés et les valeurs proviennent de la séquence des inputs.
  - ▶ Pour la self-attention qui se déroule dans le décodeur, les requêtes, les clés et les valeurs proviennent de la séquence des mots décodés jusqu'à maintenant.
  - ▶ Pour l'attention qui se déroule dans le décodeur, les requêtes proviennent des outputs de l'encodeur, et les clés et les valeurs proviennent de la séquence des mots décodés jusqu'à maintenant.

## SCALED DOT-PRODUCT ATTENTION

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

## Scaled Dot-Product Attention



## Multi-Head Attention

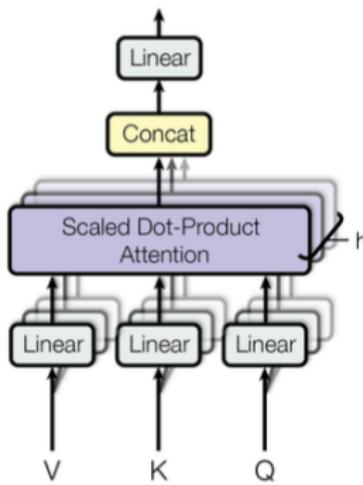
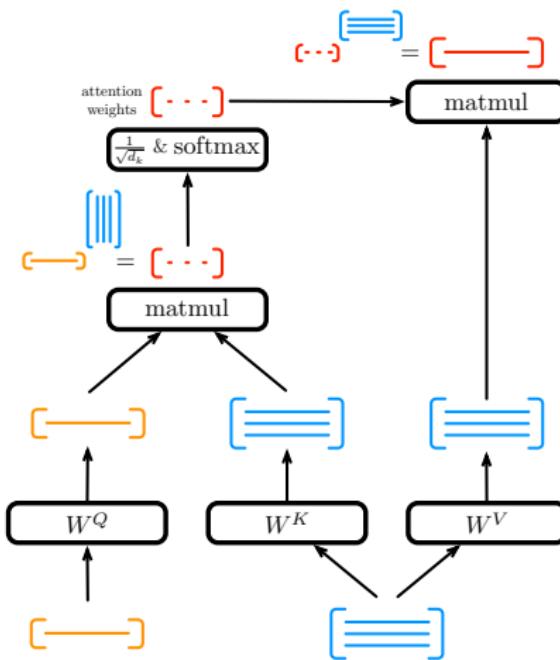


Figure taken from [\[Vaswani et al., 2017\]](#).

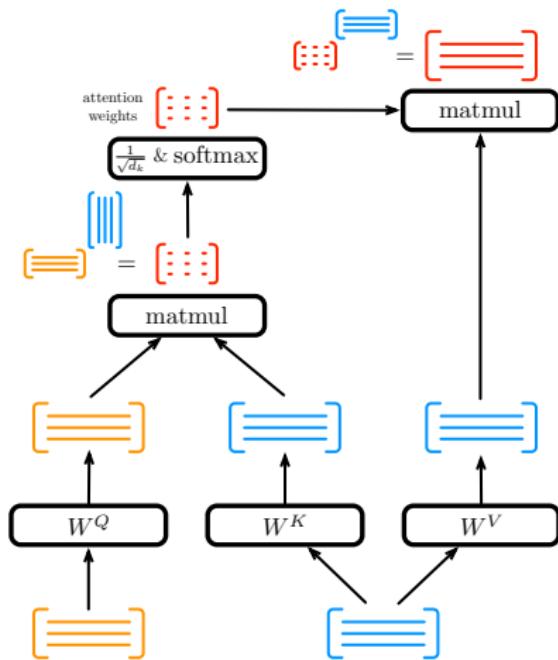
## SCALED DOT-PRODUCT ATTENTION

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$



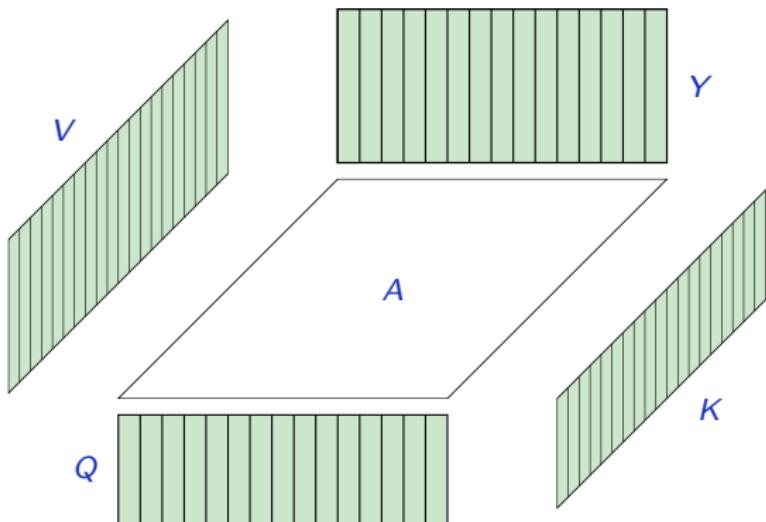
## SCALED DOT-PRODUCT ATTENTION

- On peut appliquer le processus à plusieurs *queries* en parallèle.



## SCALED DOT-PRODUCT ATTENTION

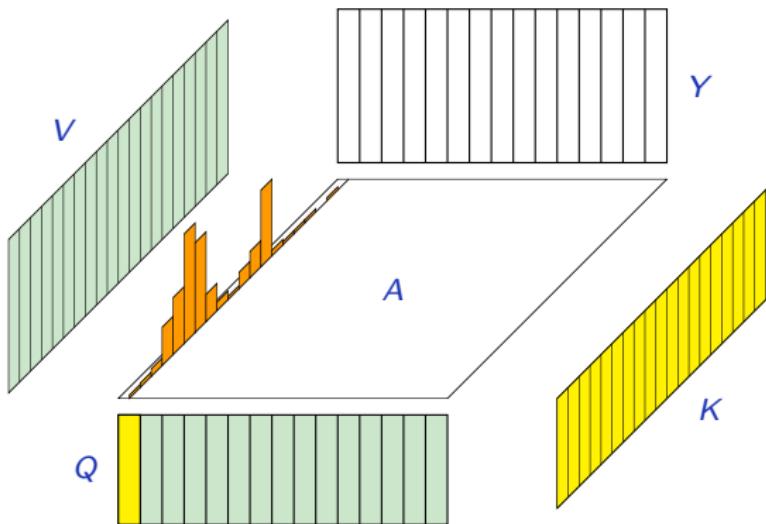
$$A_i = \text{softmax} \left( \frac{KQ_i}{\sqrt{D}} \right) \quad Y_i = V^\top A_i$$



Figures taken from [\[Fleuret, 2022\]](#).

## SCALED DOT-PRODUCT ATTENTION

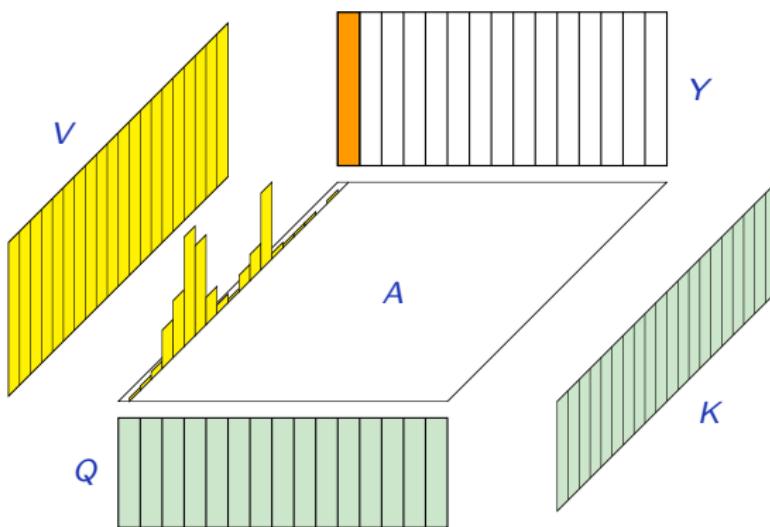
$$A_i = \text{softmax} \left( \frac{KQ_i}{\sqrt{D}} \right)$$



Figures taken from [Fleuret, 2022].

## SCALED DOT-PRODUCT ATTENTION

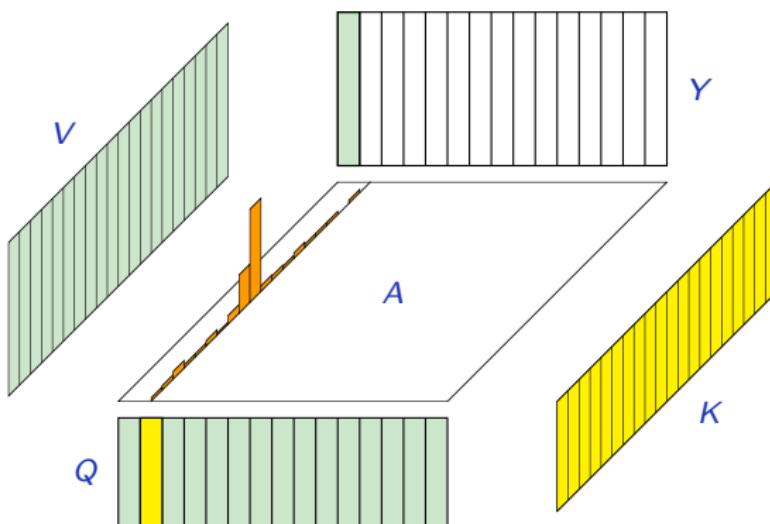
$$Y_i = V^\top A_i$$



Figures taken from [Fleuret, 2022].

## SCALED DOT-PRODUCT ATTENTION

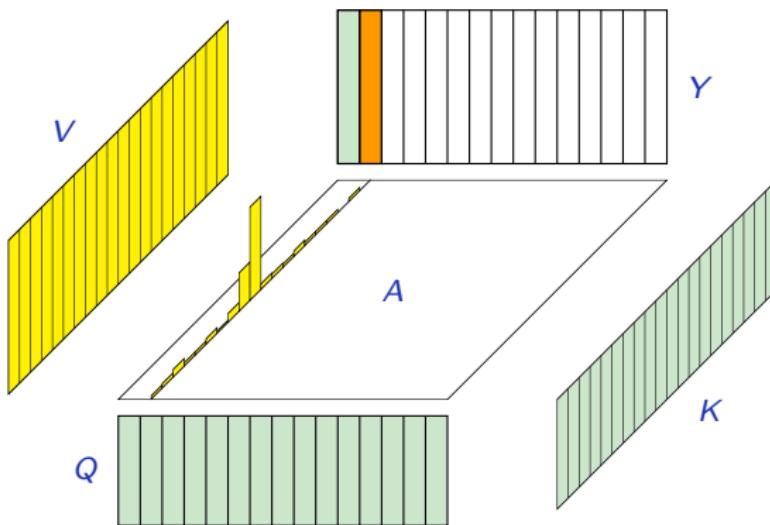
$$A_i = \text{softmax} \left( \frac{KQ_i}{\sqrt{D}} \right)$$



Figures taken from [Fleuret, 2022].

## SCALED DOT-PRODUCT ATTENTION

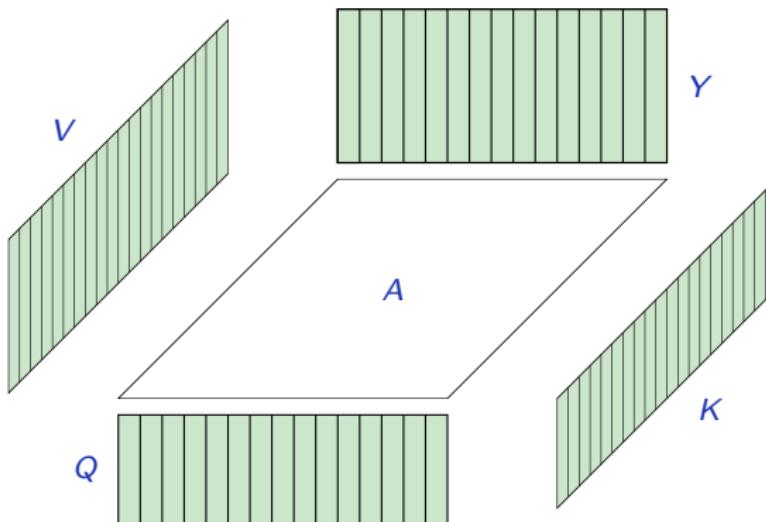
$$Y_i = V^\top A_i$$



Figures taken from [Fleuret, 2022].

## SCALED DOT-PRODUCT ATTENTION

$$A_i = \text{softmax} \left( \frac{KQ_i}{\sqrt{D}} \right) \quad Y_i = V^T A_i$$



Figures taken from [\[Fleuret, 2022\]](#).

# SCALED DOT-PRODUCT ATTENTION

$$Q = W^Q(X) \quad K = W^K(X') \quad V = W^V(X')$$
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

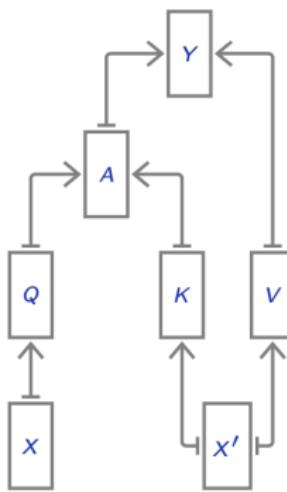
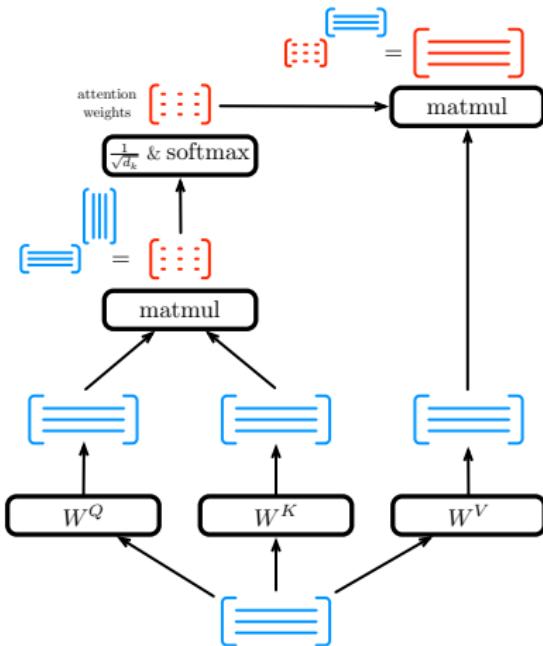


Figure taken from [Fleuret, 2022].

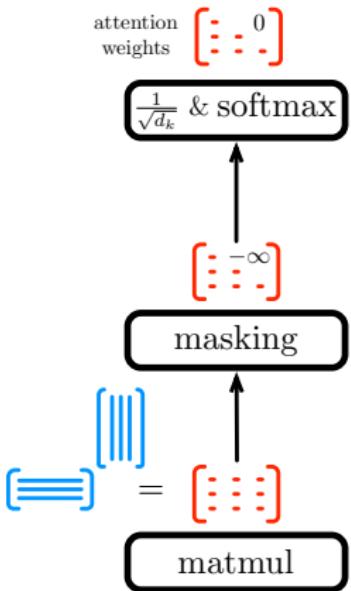
# SELF-ATTENTION

- Lorsque les *queries* proviennent de la même sequence que les *keys* et les *values*, on parle de **self-attention**.



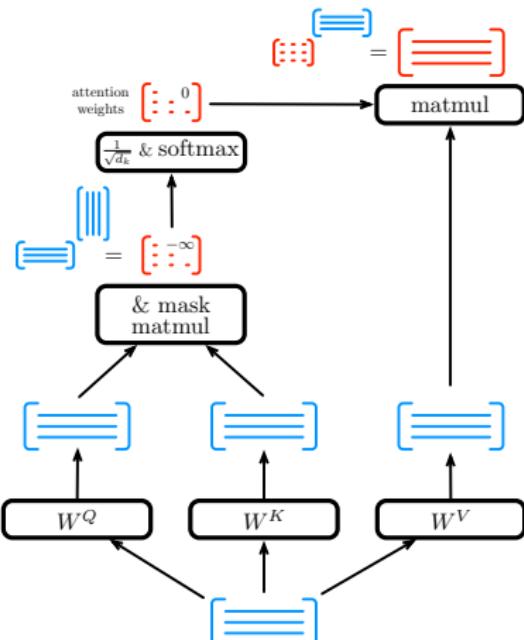
## MASKED SELF-ATTENTION

- ▶ Lorsqu'on interdit aux *queries* de porter de l'attention sur les éléments suivants de la séquence, on parle de **masked self-attention**.



# MASKED SELF-ATTENTION

- ▶ Lorsqu'on interdit aux *queries* de porter de l'attention sur les éléments suivants de la séquence, on parle de **masked self-attention**.



## SELF-ATTENTION AND MASKED SELF-ATTENTION

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

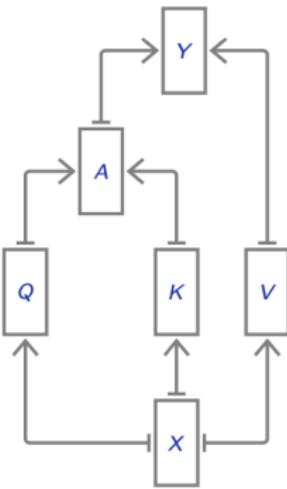


Figure taken from [Fleuret, 2022].

## MULTIPLE ATTENTION HEADS

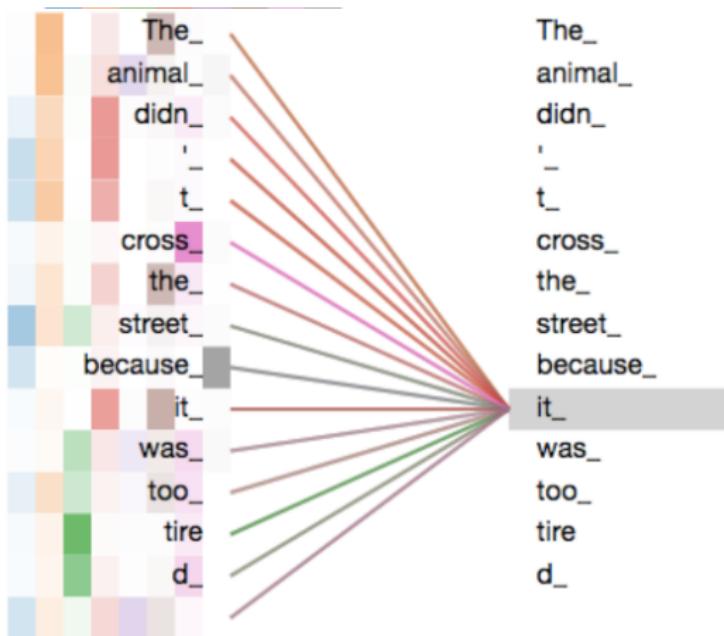


Figure taken from [Alammar, 2018](#).

## EXEMPLE

## ► Tâche: Détection de motif

# EXEMPLE

- ▶ **Tâche: Détection de motif**
- ▶ On donne une suite de 20 lettres parmi A, B, C et D:  
Par exemple [B, A, C, A, D, D, A, B, C, B, C, A, B, A, D, D, A, B, C, B].
- ▶ On veut prédire le symbole juste après le dernier A (B ici).
- ▶ Un *Recurrent Neural Network (RNN)* ou *Multi Layer Perceptron (MLP)* doit parcourir la séquence entière et mémoriser où se trouve le dernier A.
- ▶ Un *mini transformer avec attention* peut apprendre très vite à pointer le dernier A, puis regarder le symbole qui suit.

# EXEMPLE

- ▶ **Tâche: Détection de motif**
- ▶ On donne une suite de 20 lettres parmi A, B, C et D:  
Par exemple [B, A, C, A, D, D, A, B, C, B, C, A, B, A, D, D, A, B, C, B].
- ▶ On veut prédire le symbole juste après le dernier A (B ici).
- ▶ Un *Recurrent Neural Network (RNN)* ou *Multi Layer Perceptron (MLP)* doit parcourir la séquence entière et mémoriser où se trouve le dernier A.
- ▶ Un *mini transformer avec attention* peut apprendre très vite à pointer le dernier A, puis regarder le symbole qui suit.

# EXEMPLE

- ▶ **Tâche: Détection de motif**
- ▶ On donne une suite de 20 lettres parmi A, B, C et D:  
Par exemple [B, A, C, A, D, D, A, B, C, B, C, A, B, A, D, D, A, B, C, B].
- ▶ On veut prédire le symbole juste après le dernier A (B ici).
- ▶ Un *Recurrent Neural Network (RNN)* ou *Multi Layer Perceptron (MLP)* doit parcourir la séquence entière et mémoriser où se trouve le dernier A.
- ▶ Un *mini transformer avec attention* peut apprendre très vite à pointer le dernier A, puis regarder le symbole qui suit.

## EXEMPLE

- ▶ **Tâche: Détection de motif**
  - ▶ On donne une suite de 20 lettres parmi A, B, C et D:  
Par exemple [B, A, C, A, D, D, A, B, C, B, C, A, B, A, D, D, A, B, C, B].
  - ▶ On veut prédire le symbole juste après le dernier A (B ici).
  - ▶ Un *Recurrent Neural Network (RNN)* ou *Multi Layer Perceptron (MLP)* doit parcourir la séquence entière et mémoriser où se trouve le dernier A.
  - ▶ Un *mini transformer* avec *attention* peut apprendre très vite à pointer le dernier A, puis regarder le symbole qui suit.

# EXEMPLE

## ▶ Modèle 1: MLP

- ▶ Embedding dim = 32 et output dim = 4 (nb de lettres)
- ▶ 2 fully connected layers:  $20 \times 32 \rightarrow 64$  et  $64 \rightarrow 4$
- ▶ Nb de paramètres = 41'509
- ▶ Entraînement: 300 époques sur le train set (1500 samples)
- ▶ Résultats: Précision de 56.4% sur le test set (500 samples)

## EXEMPLE

## ► Modèle 1: MLP

- ▶ Embedding dim = 32 et output dim = 4 (nb de lettres)
  - ▶ 2 fully connected layers:  $20 \times 32 \rightarrow 64$  et  $64 \rightarrow 4$
  - ▶ Nb de paramètres = 41'509
  - ▶ Entraînement: 300 époques sur le train set (1500 samples)
  - ▶ Résultats: Précision de 56.4% sur le test set (500 samples)

## EXEMPLE

## ► Modèle 1: MLP

- ▶ Embedding dim = 32 et output dim = 4 (nb de lettres)
  - ▶ 2 fully connected layers:  $20 \times 32 \rightarrow 64$  et  $64 \rightarrow 4$
  - ▶ Nb de paramètres = 41'509
  - ▶ Entraînement: 300 époques sur le train set (1500 samples)
  - ▶ Résultats: Précision de 56.4% sur le test set (500 samples)

## EXEMPLE

## ► Modèle 1: MLP

- ▶ Embedding dim = 32 et output dim = 4 (nb de lettres)
  - ▶ 2 fully connected layers:  $20 \times 32 \rightarrow 64$  et  $64 \rightarrow 4$
  - ▶ Nb de paramètres = 41'509
  - ▶ Entraînement: 300 époques sur le train set (1500 samples)
  - ▶ Résultats: Précision de 56.4% sur le test set (500 samples)

## EXEMPLE

- ▶ Modèle 1: MLP
  - ▶ Embedding dim = 32 et output dim = 4 (nb de lettres)
  - ▶ 2 fully connected layers:  $20 \times 32 \rightarrow 64$  et  $64 \rightarrow 4$
  - ▶ Nb de paramètres = 41'509
  - ▶ Entraînement: 300 époques sur le train set (1500 samples)
  - ▶ Résultats: Précision de 56.4% sur le test set (500 samples)

## EXEMPLE

- ▶ Modèle 1: MLP
  - ▶ Embedding dim = 32 et output dim = 4 (nb de lettres)
  - ▶ 2 fully connected layers:  $20 \times 32 \rightarrow 64$  et  $64 \rightarrow 4$
  - ▶ Nb de paramètres = 41'509
  - ▶ Entraînement: 300 époques sur le train set (1500 samples)
  - ▶ Résultats: Précision de 56.4% sur le test set (500 samples)

## EXEMPLE

```
class SimpleMLP(nn.Module):
    def __init__(self):
        super().__init__()
        self.embed = nn.Embedding(VOCAB_SIZE, EMBED_DIM)
        self.fc1 = nn.Linear(SEQ_LEN * EMBED_DIM, HIDDEN_DIM)
        self.fc2 = nn.Linear(HIDDEN_DIM, VOCAB_SIZE)

    def forward(self, x):
        emb = self.embed(x)  # (batch, seq, emb)
        flat = emb.view(x.size(0), -1)
        h = F.relu(self.fc1(flat))
        return self.fc2(h)
```

## EXEMPLE

## ► Modèle 2: Mini Transformer

## EXEMPLE

- ▶ Modèle 2: Mini Transformer
  - ▶ Embedding dim = 32, nb têtes d'attention = 1, nb de transformer blocks = 3
  - ▶ Nb de paramètres = 39'077 (à peu près comme le MLP)
  - ▶ Entraînement: 300 époques sur le train set (1500 samples)
  - ▶ Résultats: Précision de 95.6% sur le test set (500 samples)
  - ▶ Résultats drastiquement meilleurs que le MLP, mais entraînement plus lent.

## EXEMPLE

- ▶ Modèle 2: Mini Transformer
  - ▶ Embedding dim = 32, nb têtes d'attention = 1, nb de transformer blocks = 3
  - ▶ Nb de paramètres = 39'077 (à peu près comme le MLP)
  - ▶ Entraînement: 300 époques sur le train set (1500 samples)
  - ▶ Résultats: Précision de 95.6% sur le test set (500 samples)
  - ▶ Résultats drastiquement meilleurs que le MLP, mais entraînement plus lent.

## EXEMPLE

- ▶ Modèle 2: Mini Transformer
  - ▶ Embedding dim = 32, nb têtes d'attention = 1, nb de transformer blocks = 3
  - ▶ Nb de paramètres = 39'077 (à peu près comme le MLP)
  - ▶ Entraînement: 300 époques sur le train set (1500 samples)
  - ▶ Résultats: Précision de 95.6% sur le test set (500 samples)
  - ▶ Résultats drastiquement meilleurs que le MLP, mais entraînement plus lent.

## EXEMPLE

- ▶ Modèle 2: Mini Transformer
  - ▶ Embedding dim = 32, nb têtes d'attention = 1, nb de transformer blocks = 3
  - ▶ Nb de paramètres = 39'077 (à peu près comme le MLP)
  - ▶ Entraînement: 300 époques sur le train set (1500 samples)
  - ▶ Résultats: Précision de 95.6% sur le test set (500 samples)
  - ▶ Résultats drastiquement meilleurs que le MLP, mais entraînement plus lent.

## EXEMPLE

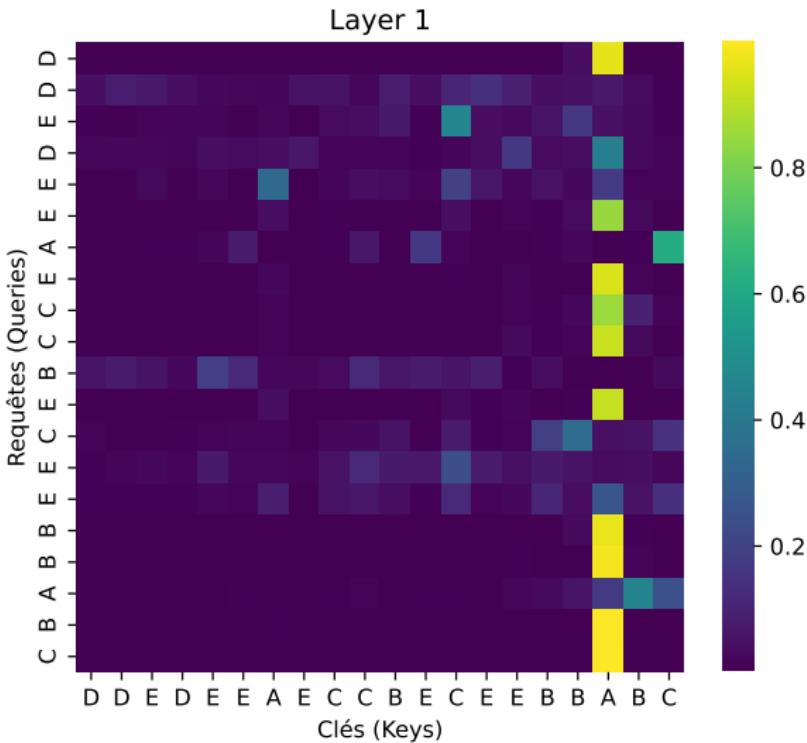
- ▶ Modèle 2: Mini Transformer
  - ▶ Embedding dim = 32, nb têtes d'attention = 1, nb de transformer blocks = 3
  - ▶ Nb de paramètres = 39'077 (à peu près comme le MLP)
  - ▶ Entraînement: 300 époques sur le train set (1500 samples)
  - ▶ Résultats: Précision de 95.6% sur le test set (500 samples)
  - ▶ Résultats drastiquement meilleurs que le MLP, mais entraînement plus lent.

## EXEMPLE

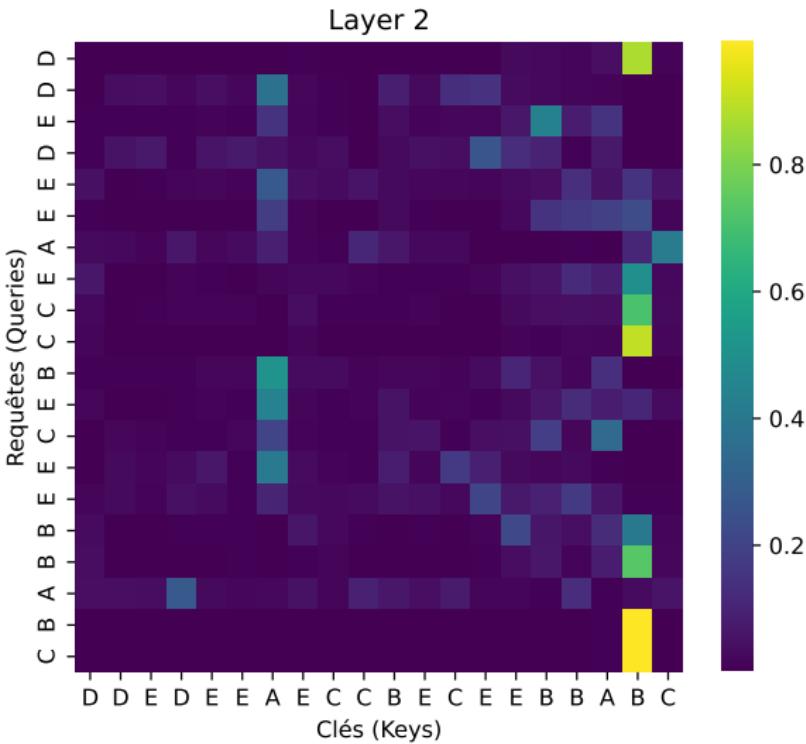
```
class MiniTransformer(nn.Module):
    def __init__(self, vocab_size=VOCAB_SIZE, seq_len=SEQ_LEN, embed_dim=EMBED_DIM, heads=1): # 1 heads only
        super().__init__()
        self.embed = nn.Embedding(vocab_size, embed_dim)
        self.pos_embed = nn.Embedding(seq_len, embed_dim)
        encoder_layer = nn.TransformerEncoderLayer(embed_dim,
                                                    heads,
                                                    dim_feedforward=HIDDEN_DIM*2, # augmented dim
                                                    batch_first=True)
        self.encoder = nn.TransformerEncoder(encoder_layer, num_layers=3) # 3 blocks!
        self.fc = nn.Linear(embed_dim, vocab_size)

    def forward(self, x):
        pos_ids = torch.arange(x.shape[1], device=x.device).unsqueeze(0)
        emb = self.embed(x) + self.pos_embed(pos_ids)
        out = self.encoder(emb)
        pooled = out[:, -1, :]
        return self.fc(pooled)
```

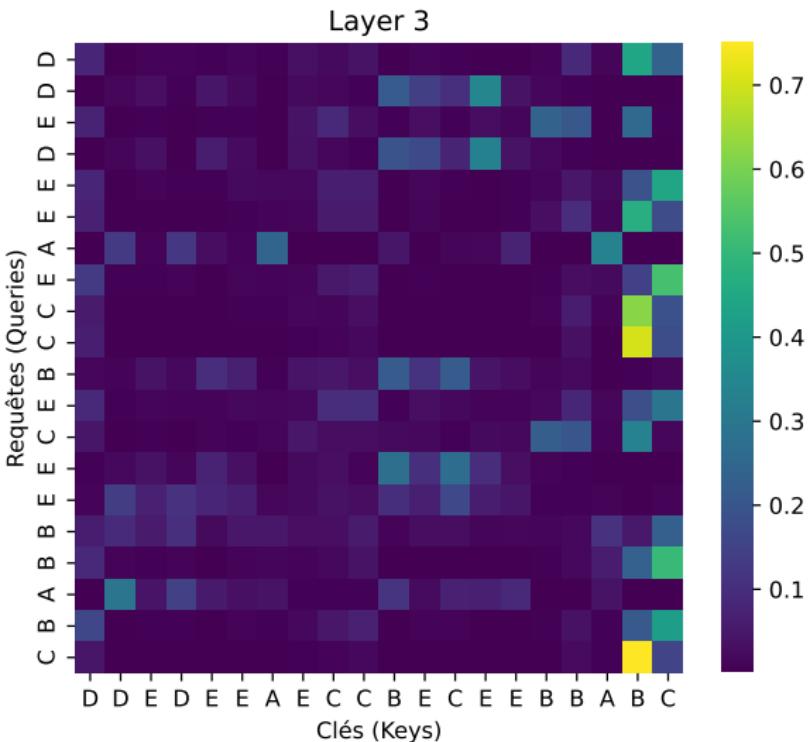
## EXEMPLE: VISUALISATION DES POIDS D'ATTENTION



## EXEMPLE: VISUALISATION DES POIDS D'ATTENTION

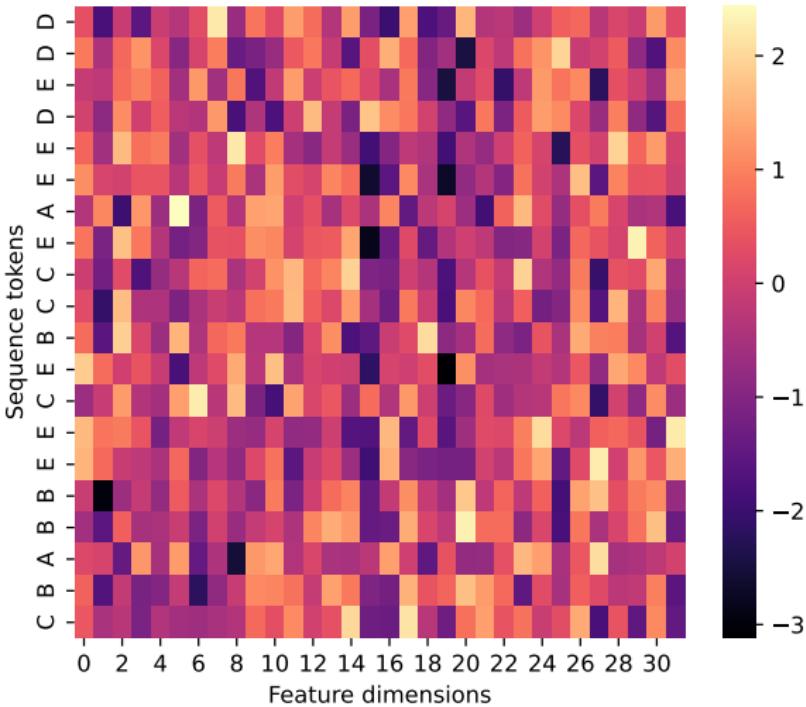


## EXEMPLE: VISUALISATION DES POIDS D'ATTENTION

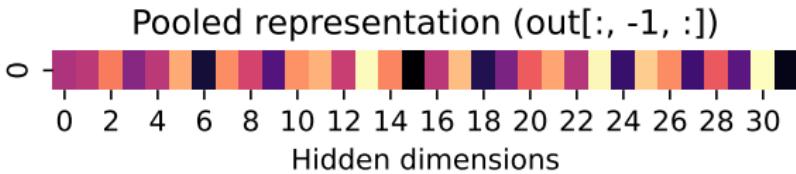


## EXEMPLE: VISUALISATION DES POIDS D'ATTENTION

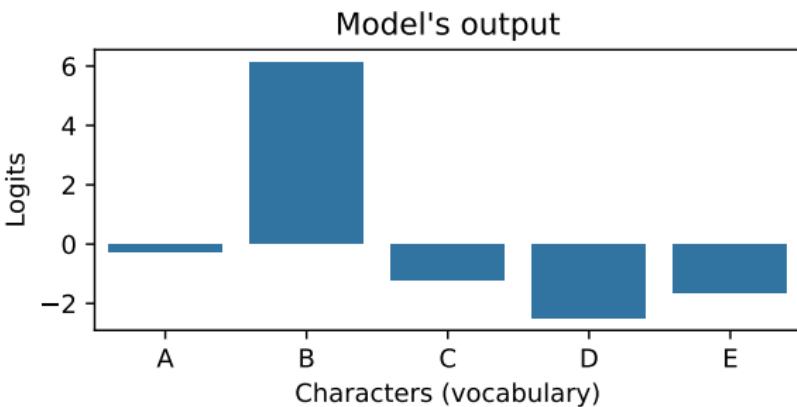
### Attention 'Values' (V) from last layer



## EXEMPLE: VISUALISATION DES POIDS D'ATTENTION



## EXEMPLE: VISUALISATION DES POIDS D'ATTENTION



# ARCHITECTURES

- ▶ Le premier Transformer était une architecture encodeur-décodeur. En pratique, on utilise d'autres architectures également...
- ▶ Architecture **encoder-only**, composée de *blocs d'encodeurs* empilés les uns sur les autres.
- ▶ Architecture **encoder-decoder**, composée de *blocs d'encodeurs* empilés suivis de *blocs de décodeurs* empilés.
- ▶ Architecture **decoder-only**, architecture composée de *blocs de décodeurs* empilés.

# ARCHITECTURES

- ▶ Le premier Transformer était une architecture encodeur-décodeur. En pratique, on utilise d'autres architectures également...
- ▶ Architecture **encoder-only**, composée de *blocs d'encodeurs* empilés les uns sur les autres.
- ▶ Architecture **encoder-decoder**, composée de *blocs d'encodeurs* empilés suivis de *blocs de décodeurs* empilés.
- ▶ Architecture **decoder-only**, architecture composée de *blocs de décodeurs* empilés.

# ARCHITECTURES

- ▶ Le premier Transformer était une architecture encodeur-décodeur. En pratique, on utilise d'autres architectures également...
- ▶ Architecture **encoder-only**, composée de *blocs d'encodeurs* empilés les uns sur les autres.
- ▶ Architecture **encoder-decoder**, composée de *blocs d'encodeurs* empilés suivis de *blocs de décodeurs* empilés.
- ▶ Architecture **decoder-only**, architecture composée de *blocs de décodeurs* empilés.

## ARCHITECTURES

- ▶ Le premier Transformer était une architecture encodeur-décodeur. En pratique, on utilise d'autres architectures également...
  - ▶ Architecture **encoder-only**, composée de *blocs d'encodeurs* empilés les uns sur les autres.
  - ▶ Architecture **encoder-decoder**, composée de *blocs d'encodeurs* empilés suivis de *blocs de décodeurs* empilés.
  - ▶ Architecture **decoder-only**, architecture composée de *blocs de décodeurs* empilés.

## ENCODER-ONLY

- ▶ **Modèles d'encodage ou Embedders**, composés de *blocks d'encodeurs* (BERT, RoBERTa, ChemBERTa, etc.).
  - ▶ Apprennent des représentations contextuelles de tokens ou de séquences complètes.
  - ▶ Pour des tâches de compréhension (et non de génération).
    - Sequence classification: sentiment analysis, spam detection.
    - Token classification: named entity recognition (NER), POS tagging
    - Semantic similarity / retrieval: recherche sémantique, embeddings.
    - Property prediction: prédiction de propriétés moléculaires.

## ENCODER-ONLY

- ▶ **Modèles d'encodage ou Embedders**, composés de *blocks d'encodeurs* (BERT, RoBERTa, ChemBERTa, etc.).
  - ▶ Apprennent des représentations contextuelles de tokens ou de séquences complètes.
  - ▶ Pour des tâches de compréhension (et non de génération).
    - Sequence classification: sentiment analysis, spam detection.
    - Token classification: named entity recognition (NER), POS tagging
    - Semantic similarity / retrieval: recherche sémantique, embeddings.
    - Property prediction: prédiction de propriétés moléculaires.

# ENCODER-ONLY

- ▶ **Modèles d'encodage** ou **Embedders**, composés de *blocks d'encodeurs* (BERT, RoBERTa, ChemBERTa, etc.).
- ▶ Apprennent des représentations contextuelles de tokens ou de séquences complètes.
- ▶ Pour des tâches de compréhension (et non de génération).
  - Sequence classification: sentiment analysis, spam detection.
  - Token classification: named entity recognition (NER), POS tagging
  - Semantic similarity / retrieval: recherche sémantique, embeddings.
  - Property prediction: prédiction de propriétés moléculaires.

## ENCODER-ONLY

- ▶ Modèles d'encodage ou **Embedders**, composés de *blocks d'encodeurs* (BERT, RoBERTa, ChemBERTa, etc.).
  - ▶ Apprennent des représentations contextuelles de tokens ou de séquences complètes.
  - ▶ Pour des tâches de compréhension (et non de génération).
    - Sequence classification: sentiment analysis, spam detection.
    - Token classification: named entity recognition (NER), POS tagging
    - Semantic similarity / retrieval: recherche sémantique, embeddings.
    - Property prediction: prédiction de propriétés moléculaires.

## ENCODER-ONLY

- ▶ Modèles d'encodage ou **Embedders**, composés de *blocks d'encodeurs* (BERT, RoBERTa, ChemBERTa, etc.).
  - ▶ Apprennent des représentations contextuelles de tokens ou de séquences complètes.
  - ▶ Pour des tâches de compréhension (et non de génération).
    - Sequence classification: sentiment analysis, spam detection.
    - Token classification: named entity recognition (NER), POS tagging
    - Semantic similarity / retrieval: recherche sémantique, embeddings.
    - Property prediction: prédiction de propriétés moléculaires.

## ENCODER-ONLY

- ▶ Modèles d'encodage ou **Embedders**, composés de *blocks d'encodeurs* (BERT, RoBERTa, ChemBERTa, etc.).
  - ▶ Apprennent des représentations contextuelles de tokens ou de séquences complètes.
  - ▶ Pour des tâches de compréhension (et non de génération).
    - Sequence classification: sentiment analysis, spam detection.
    - Token classification: named entity recognition (NER), POS tagging
    - Semantic similarity / retrieval: recherche sémantique, embeddings.
    - Property prediction: prédiction de propriétés moléculaires.

## ENCODER-ONLY

- ▶ Modèles d'encodage ou **Embedders**, composés de *blocks d'encodeurs* (BERT, RoBERTa, ChemBERTa, etc.).
  - ▶ Apprennent des représentations contextuelles de tokens ou de séquences complètes.
  - ▶ Pour des tâches de compréhension (et non de génération).
    - Sequence classification: sentiment analysis, spam detection.
    - Token classification: named entity recognition (NER), POS tagging
    - Semantic similarity / retrieval: recherche sémantique, embeddings.
    - Property prediction: prédiction de propriétés moléculaires.

## ENCODER-ONLY

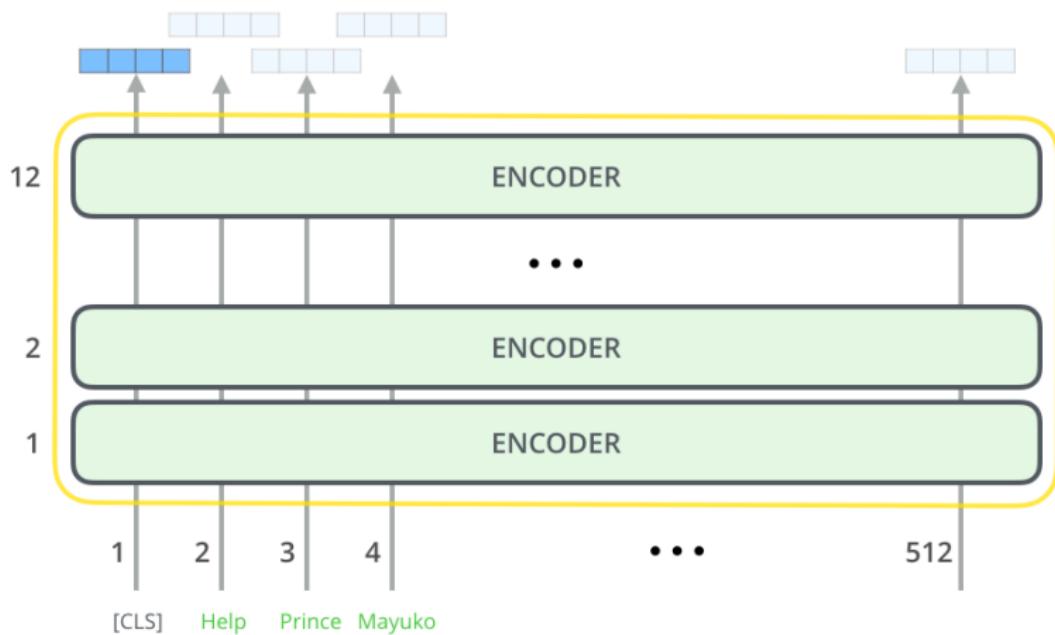


Figure taken from [\[Alammar, 2018\]](#).

## ENCODER-ONLY: ENTRAÎNEMENT

- Le modèle est entraîné en mode **supervisé** classique.
    1. L'encodeur produit une représentation contextuelle de la séquence d'input (input embedding), transmise à un **MLP** de sortie: tête de classification ou de régression.
    2. Le modèle calcule la loss entre les prédictions et les labels.
    3. Les paramètres sont mis à jour par **backpropagation**.
  - **Avantages:** simplicité d'entraînement et convergence rapide.

## ENCODER-ONLY: ENTRAÎNEMENT

- Le modèle est entraîné en mode **supervisé** classique.
    1. L'encodeur produit une représentation contextuelle de la séquence d'input (**input embedding**), transmise à un **MLP de sortie**: tête de classification ou de régression.
    2. Le modèle calcule la loss entre les prédictions et les labels.
    3. Les paramètres sont mis à jour par **backpropagation**.
  - **Avantages:** simplicité d'entraînement et convergence rapide.

## ENCODER-ONLY: ENTRAÎNEMENT

- Le modèle est entraîné en mode **supervisé** classique.
    1. L'encodeur produit une représentation contextuelle de la séquence d'input (input embedding), transmise à un **MLP de sortie**: tête de classification ou de régression.
    2. Le modèle calcule la loss entre les prédictions et les labels.
    3. Les paramètres sont mis à jour par **backpropagation**.
  - **Avantages:** simplicité d'entraînement et convergence rapide.

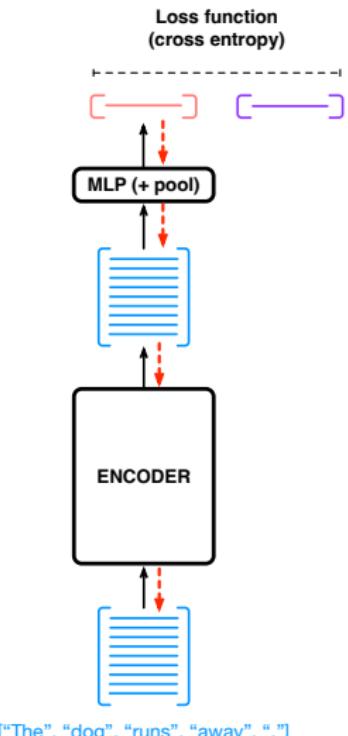
## ENCODER-ONLY: ENTRAÎNEMENT

- Le modèle est entraîné en mode **supervisé** classique.
    1. L'encodeur produit une représentation contextuelle de la séquence d'input (input embedding), transmise à un **MLP de sortie**: tête de classification ou de régression.
    2. Le modèle calcule la loss entre les prédictions et les labels.
    3. Les paramètres sont mis à jour par **backpropagation**.
  - **Avantages:** simplicité d'entraînement et convergence rapide.

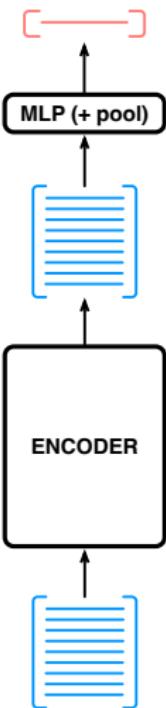
## ENCODER-ONLY: ENTRAÎNEMENT

- Le modèle est entraîné en mode **supervisé** classique.
    1. L'encodeur produit une représentation contextuelle de la séquence d'input (input embedding), transmise à un **MLP de sortie**: tête de classification ou de régression.
    2. Le modèle calcule la loss entre les prédictions et les labels.
    3. Les paramètres sont mis à jour par **backpropagation**.
  - **Avantages:** simplicité d'entraînement et convergence rapide.

## ENCODER-ONLY: ENTRAÎNEMENT



## ENCODER-ONLY: INFÉRENCE



[“The”, “dog”, “runs”, “away”, “.”]

# ENCODER-DECODER

- ▶ **Modèles sequence-to-sequence**, composés de *blocs d'encodeurs et de décodeurs* (T5, BART, mBART, etc.).
- ▶ Apprennent une transformation de séquences en séquences.
- ▶ Pour des tâches de transformation de texte à texte.
  - Machine translation: traduction automatique entre langues.
  - Summarization: génération de résumés de documents.
  - Question answering: réponse générative à une question.
  - Paraphrase generation: reformulation de texte.

# ENCODER-DECODER

- ▶ **Modèles sequence-to-sequence**, composés de *blocs d'encodeurs et de décodeurs* (T5, BART, mBART, etc.).
- ▶ Apprennent une transformation de séquences en séquences.
- ▶ Pour des tâches de transformation de texte à texte.
  - Machine translation: traduction automatique entre langues.
  - Summarization: génération de résumés de documents.
  - Question answering: réponse générative à une question.
  - Paraphrase generation: reformulation de texte.

# ENCODER-DECODER

- ▶ **Modèles sequence-to-sequence**, composés de *blocs d'encodeurs et de décodeurs* (T5, BART, mBART, etc.).
- ▶ Apprennent une transformation de séquences en séquences.
- ▶ Pour des tâches de transformation de texte à texte.
  - Machine translation: traduction automatique entre langues.
  - Summarization: génération de résumés de documents.
  - Question answering: réponse générative à une question.
  - Paraphrase generation: reformulation de texte.

# ENCODER-DECODER

- ▶ **Modèles sequence-to-sequence**, composés de *blocs d'encodeurs et de décodeurs* (T5, BART, mBART, etc.).
- ▶ Apprennent une transformation de séquences en séquences.
- ▶ Pour des tâches de transformation de texte à texte.
  - Machine translation: traduction automatique entre langues.
  - Summarization: génération de résumés de documents.
  - Question answering: réponse générative à une question.
  - Paraphrase generation: reformulation de texte.

# ENCODER-DECODER

- ▶ **Modèles sequence-to-sequence**, composés de *blocs d'encodeurs et de décodeurs* (T5, BART, mBART, etc.).
- ▶ Apprennent une transformation de séquences en séquences.
- ▶ Pour des tâches de transformation de texte à texte.
  - Machine translation: traduction automatique entre langues.
  - Summarization: génération de résumés de documents.
  - Question answering: réponse générative à une question.
  - Paraphrase generation: reformulation de texte.

# ENCODER-DECODER

- ▶ **Modèles sequence-to-sequence**, composés de *blocs d'encodeurs et de décodeurs* (T5, BART, mBART, etc.).
- ▶ Apprennent une transformation de séquences en séquences.
- ▶ Pour des tâches de transformation de texte à texte.
  - Machine translation: traduction automatique entre langues.
  - Summarization: génération de résumés de documents.
  - Question answering: réponse générative à une question.
  - Paraphrase generation: reformulation de texte.

# ENCODER-DECODER

- ▶ **Modèles sequence-to-sequence**, composés de *blocs d'encodeurs et de décodeurs* (T5, BART, mBART, etc.).
- ▶ Apprennent une transformation de séquences en séquences.
- ▶ Pour des tâches de transformation de texte à texte.
  - Machine translation: traduction automatique entre langues.
  - Summarization: génération de résumés de documents.
  - Question answering: réponse générative à une question.
  - Paraphrase generation: reformulation de texte.

# ENCODER-DECODER

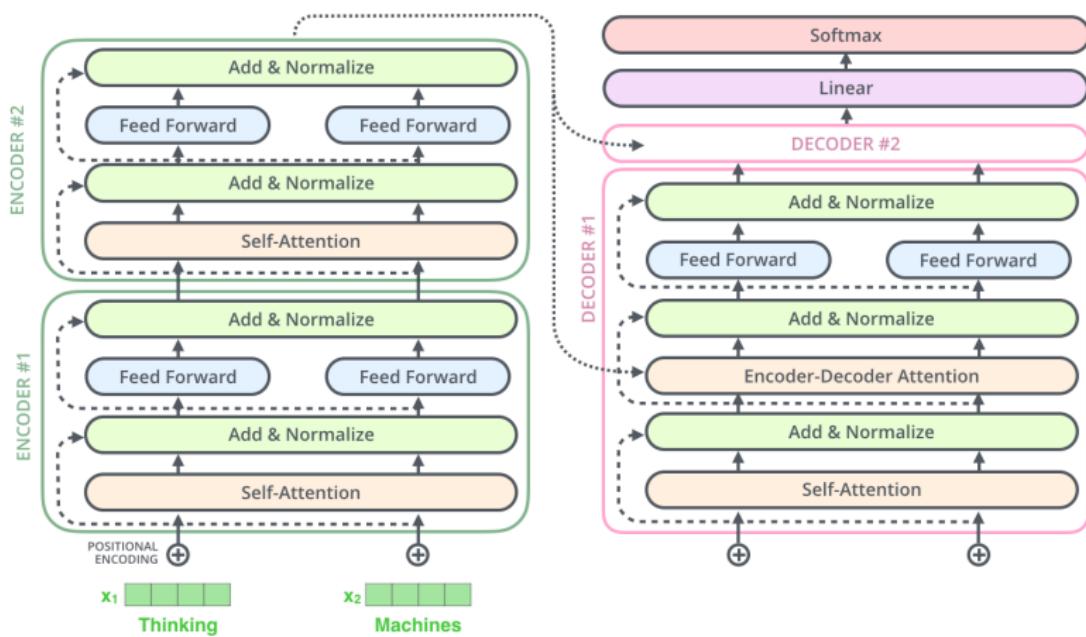


Figure taken from [Alammar, 2018].

# ENCODER-DECODER: ENTRAÎNEMENT

- ▶ **Le modèle est entraîné en mode **teacher forcing**.**
- 1. L'encodeur produit une représentation contextuelle de la séquence d'input (input embedding).
- 2. À chaque étape  $i$ , le décodeur reçoit l'input embedding en entier, ainsi que le token [START] suivi des  $i - 1$  premiers tokens de la séquence cible, et génère un  $i$ -ème token.
- 3. Le modèle calcule la loss entre le  $i$ -ème token généré et le  $i$ -ème token cible.
- 4. Les paramètres sont mis à jour par **backpropagation**.
- ▶ **Avantages:** entraînement parallélisable sur toute la séquence et absence d'accumulation d'erreurs pendant l'apprentissage.

# ENCODER-DECODER: ENTRAÎNEMENT

- ▶ Le modèle est entraîné en mode **teacher forcing**.
- 1. L'encodeur produit une représentation contextuelle de la séquence d'input (input embedding).
- 2. À chaque étape  $i$ , le décodeur reçoit l'input embedding en entier, ainsi que le token [START] suivi des  $i - 1$  premiers tokens de la séquence cible, et génère un  $i$ -ème token.
- 3. Le modèle calcule la loss entre le  $i$ -ème token généré et le  $i$ -ème token cible.
- 4. Les paramètres sont mis à jour par **backpropagation**.
- ▶ **Avantages:** entraînement parallélisable sur toute la séquence et absence d'accumulation d'erreurs pendant l'apprentissage.

# ENCODER-DECODER: ENTRAÎNEMENT

- ▶ Le modèle est entraîné en mode **teacher forcing**.
- 1. L'encodeur produit une représentation contextuelle de la séquence d'input (input embedding).
- 2. À chaque étape  $i$ , le décodeur reçoit l'input embedding en entier, ainsi que le token [START] suivi des  $i - 1$  premiers tokens de la séquence cible, et génère un  $i$ -ème token.
- 3. Le modèle calcule la loss entre le  $i$ -ème token généré et le  $i$ -ème token cible.
- 4. Les paramètres sont mis à jour par **backpropagation**.
- ▶ **Avantages:** entraînement parallélisable sur toute la séquence et absence d'accumulation d'erreurs pendant l'apprentissage.

# ENCODER-DECODER: ENTRAÎNEMENT

- ▶ Le modèle est entraîné en mode **teacher forcing**.
- 1. L'encodeur produit une représentation contextuelle de la séquence d'input (input embedding).
- 2. À chaque étape  $i$ , le décodeur reçoit l'input embedding en entier, ainsi que le token [START] suivi des  $i - 1$  premiers tokens de la séquence cible, et génère un  $i$ -ème token.
- 3. Le modèle calcule la loss entre le  $i$ -ème token généré et le  $i$ -ème token cible.
- 4. Les paramètres sont mis à jour par **backpropagation**.
- ▶ **Avantages:** entraînement parallélisable sur toute la séquence et absence d'accumulation d'erreurs pendant l'apprentissage.

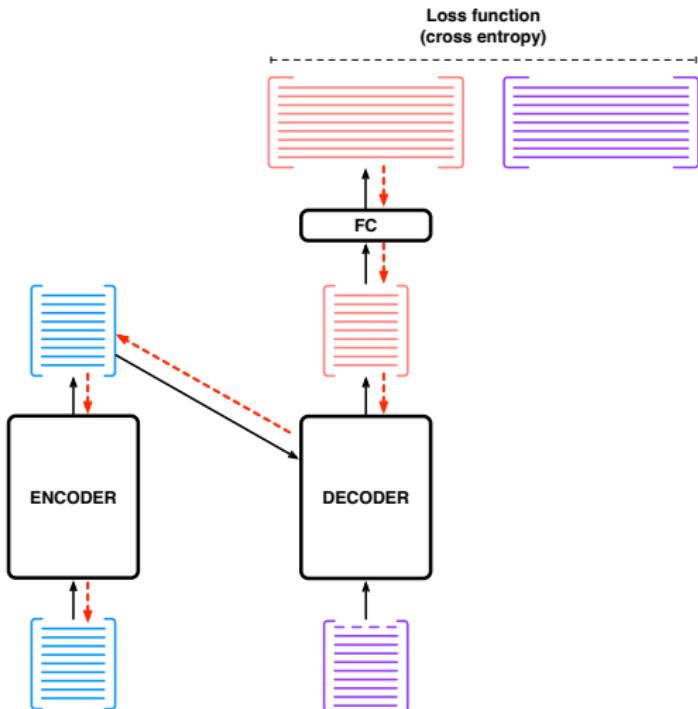
# ENCODER-DECODER: ENTRAÎNEMENT

- ▶ Le modèle est entraîné en mode **teacher forcing**.
- 1. L'encodeur produit une représentation contextuelle de la séquence d'input (input embedding).
- 2. À chaque étape  $i$ , le décodeur reçoit l'input embedding en entier, ainsi que le token [START] suivi des  $i - 1$  premiers tokens de la séquence cible, et génère un  $i$ -ème token.
- 3. Le modèle calcule la loss entre le  $i$ -ème token généré et le  $i$ -ème token cible.
- 4. Les paramètres sont mis à jour par **backpropagation**.
- ▶ Avantages: entraînement parallélisable sur toute la séquence et absence d'accumulation d'erreurs pendant l'apprentissage.

# ENCODER-DECODER: ENTRAÎNEMENT

- ▶ Le modèle est entraîné en mode **teacher forcing**.
  1. L'encodeur produit une représentation contextuelle de la séquence d'input (input embedding).
  2. À chaque étape  $i$ , le décodeur reçoit l'input embedding en entier, ainsi que le token [START] suivi des  $i - 1$  premiers tokens de la séquence cible, et génère un  $i$ -ème token.
  3. Le modèle calcule la loss entre le  $i$ -ème token généré et le  $i$ -ème token cible.
  4. Les paramètres sont mis à jour par **backpropagation**.
- ▶ **Avantages:** entraînement parallélisable sur toute la séquence et absence d'accumulation d'erreurs pendant l'apprentissage.

# ENCODER-DECODER: ENTRAÎNEMENT



**["The", "dog", "runs", "away", "."]** **[[START], "Le", "chien", "s", "enfuit", "."]**

# ENCODER-DECODER: INFÉRENCE

► **Mode auto-régressif:** le modèle génère une séquence de texte *pas à pas*, un token à la fois.

1. L'encodeur traite la séquence d'input une seule fois et produit une représentation contextuelle (input embedding).
  2. À chaque étape, le décodeur reçoit tout l'input embedding ainsi les tokens d'output déjà générés, et prédit le token suivant.
  3. Ce nouveau token est ajouté aux tokens d'output déjà générés, qui sont réinjectés dans le décodeur.
  4. Le processus continue jusqu'à la production du token [END].
- **Stratégies de décodage:** greedy decoding, beam search, sampling.

# ENCODER-DECODER: INFÉRENCE

- ▶ **Mode auto-régressif:** le modèle génère une séquence de texte *pas à pas*, un token à la fois.
- 1. L'encodeur traite la séquence d'input une seule fois et produit une représentation contextuelle (input embedding).
- 2. À chaque étape, le décodeur reçoit tout l'input embedding ainsi les tokens d'output déjà générés, et prédit le token suivant.
- 3. Ce nouveau token est ajouté aux tokens d'output déjà générés, qui sont réinjectés dans le décodeur.
- 4. Le processus continue jusqu'à la production du token [END].
- ▶ **Stratégies de décodage:** greedy decoding, beam search, sampling.

# ENCODER-DECODER: INFÉRENCE

- ▶ **Mode auto-régressif:** le modèle génère une séquence de texte *pas à pas*, un token à la fois.
- 1. L'encodeur traite la séquence d'input une seule fois et produit une représentation contextuelle (input embedding).
- 2. À chaque étape, le décodeur reçoit tout l'input embedding ainsi les tokens d'output déjà générés, et prédit le token suivant.
- 3. Ce nouveau token est ajouté aux tokens d'output déjà générés, qui sont réinjectés dans le décodeur.
- 4. Le processus continue jusqu'à la production du token [END].
- ▶ **Stratégies de décodage:** greedy decoding, beam search, sampling.

# ENCODER-DECODER: INFÉRENCE

- ▶ **Mode auto-régressif:** le modèle génère une séquence de texte *pas à pas*, un token à la fois.
- 1. L'encodeur traite la séquence d'input une seule fois et produit une représentation contextuelle (input embedding).
- 2. À chaque étape, le décodeur reçoit tout l'input embedding ainsi les tokens d'output déjà générés, et prédit le token suivant.
- 3. Ce nouveau token est ajouté aux tokens d'output déjà générés, qui sont réinjectés dans le décodeur.
- 4. Le processus continue jusqu'à la production du token [END].
- ▶ **Stratégies de décodage:** greedy decoding, beam search, sampling.

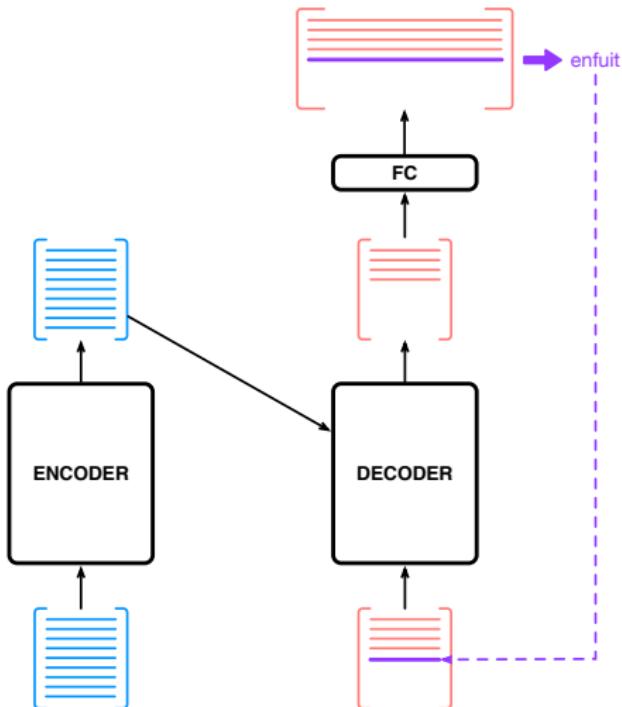
# ENCODER-DECODER: INFÉRENCE

- ▶ **Mode auto-régressif:** le modèle génère une séquence de texte *pas à pas*, un token à la fois.
  1. L'encodeur traite la séquence d'input une seule fois et produit une représentation contextuelle (input embedding).
  2. À chaque étape, le décodeur reçoit tout l'input embedding ainsi les tokens d'output déjà générés, et prédit le token suivant.
  3. Ce nouveau token est ajouté aux tokens d'output déjà générés, qui sont réinjectés dans le décodeur.
  4. Le processus continue jusqu'à la production du token [END].
- ▶ **Stratégies de décodage:** greedy decoding, beam search, sampling.

# ENCODER-DECODER: INFÉRENCE

- ▶ **Mode auto-régressif:** le modèle génère une séquence de texte *pas à pas*, un token à la fois.
- 1. L'encodeur traite la séquence d'input une seule fois et produit une représentation contextuelle (input embedding).
- 2. À chaque étape, le décodeur reçoit tout l'input embedding ainsi les tokens d'output déjà générés, et prédit le token suivant.
- 3. Ce nouveau token est ajouté aux tokens d'output déjà générés, qui sont réinjectés dans le décodeur.
- 4. Le processus continue jusqu'à la production du token [END].
- ▶ **Stratégies de décodage:** greedy decoding, beam search, sampling.

## ENCODER-DECODER: INFÉRENCE



[“The”, “dog”, “runs”, “away”, “.”] [[START], “Le”, “chien”, “s”, “enfuit”]

# DECODER-ONLY

- ▶ **Modèles génératifs auto-régressifs**, composés de *blocs de décodeurs* uniquement (GPT, LLaMA, Mistral, etc.).
- ▶ Apprennent à générer des séquences token par token, en se basant sur les tokens précédents: mode auto-régressif.
- ▶ Pour des tâches de génération libre.
  - Language modeling: prédiction du prochain mot dans une séquence.
  - Text generation: génération libre de texte, complétion de phrase.
  - Dialogue / chat: agents conversationnels, assistants.
  - Code generation: génération automatique de code.

# DECODER-ONLY

- ▶ **Modèles génératifs auto-régressifs**, composés de *blocs de décodeurs* uniquement (GPT, LLaMA, Mistral, etc.).
- ▶ Apprennent à générer des séquences token par token, en se basant sur les tokens précédents: mode auto-régressif.
- ▶ Pour des tâches de génération libre.
  - Language modeling: prédiction du prochain mot dans une séquence.
  - Text generation: génération libre de texte, complétion de phrase.
  - Dialogue / chat: agents conversationnels, assistants.
  - Code generation: génération automatique de code.

# DECODER-ONLY

- ▶ **Modèles génératifs auto-régressifs**, composés de *blocs de décodeurs* uniquement (GPT, LLaMA, Mistral, etc.).
- ▶ Apprennent à générer des séquences token par token, en se basant sur les tokens précédents: mode auto-régressif.
- ▶ Pour des tâches de génération libre.
  - Language modeling: prédiction du prochain mot dans une séquence.
  - Text generation: génération libre de texte, complétion de phrase.
  - Dialogue / chat: agents conversationnels, assistants.
  - Code generation: génération automatique de code.

# DECODER-ONLY

- ▶ **Modèles génératifs auto-régressifs**, composés de *blocs de décodeurs* uniquement (GPT, LLaMA, Mistral, etc.).
- ▶ Apprennent à générer des séquences token par token, en se basant sur les tokens précédents: mode auto-régressif.
- ▶ Pour des tâches de génération libre.
  - Language modeling: prédiction du prochain mot dans une séquence.
  - Text generation: génération libre de texte, complétion de phrase.
  - Dialogue / chat: agents conversationnels, assistants.
  - Code generation: génération automatique de code.

# DECODER-ONLY

- ▶ **Modèles génératifs auto-régressifs**, composés de *blocs de décodeurs* uniquement (GPT, LLaMA, Mistral, etc.).
- ▶ Apprennent à générer des séquences token par token, en se basant sur les tokens précédents: mode auto-régressif.
- ▶ Pour des tâches de génération libre.
  - Language modeling: prédiction du prochain mot dans une séquence.
  - Text generation: génération libre de texte, complétion de phrase.
    - Dialogue / chat: agents conversationnels, assistants.
    - Code generation: génération automatique de code.

# DECODER-ONLY

- ▶ **Modèles génératifs auto-régressifs**, composés de *blocs de décodeurs* uniquement (GPT, LLaMA, Mistral, etc.).
- ▶ Apprennent à générer des séquences token par token, en se basant sur les tokens précédents: mode auto-régressif.
- ▶ Pour des tâches de génération libre.
  - Language modeling: prédiction du prochain mot dans une séquence.
  - Text generation: génération libre de texte, complétion de phrase.
  - Dialogue / chat: agents conversationnels, assistants.
  - Code generation: génération automatique de code.

# DECODER-ONLY

- ▶ **Modèles génératifs auto-régressifs**, composés de *blocs de décodeurs* uniquement (GPT, LLaMA, Mistral, etc.).
- ▶ Apprennent à générer des séquences token par token, en se basant sur les tokens précédents: mode auto-régressif.
- ▶ Pour des tâches de génération libre.
  - Language modeling: prédiction du prochain mot dans une séquence.
  - Text generation: génération libre de texte, complétion de phrase.
  - Dialogue / chat: agents conversationnels, assistants.
  - Code generation: génération automatique de code.

# DECODER-ONLY

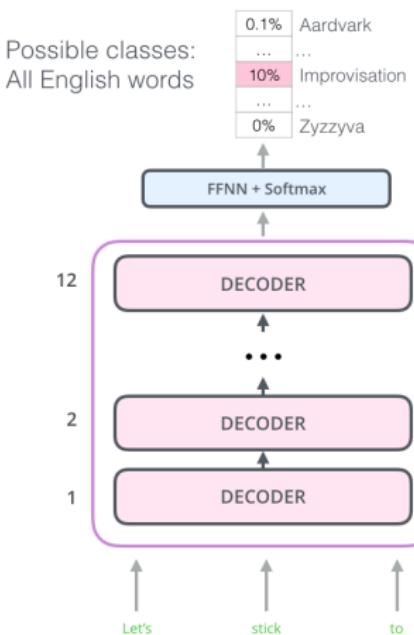


Figure taken from [Alammar, 2018].

## DECODER-ONLY: ENTRAÎNEMENT

- ▶ Le modèle est entraîné en mode **auto-régressif (causal language modeling)**.

1. Le décodeur reçoit un contexte (e.g., un prompt) ainsi que le token [START] suivi de la séquence cible en entier.
  2. Le modèle prédit alors tous les token de la séquence cible, en appliquant une **masked self-attention**, pour que chaque  $i$ -ème token ne voit que les  $i - 1$  tokens précédents.
  3. Le modèle calcule la loss entre les tokens générés et les tokens cibles.
  4. Les paramètres sont mis à jour par **backpropagation**.
- ▶ **Avantages:** entraînement parallélisé sur toute la séquence et apprentissage direct du mode auto-régressif.

## DECODER-ONLY: ENTRAÎNEMENT

- ▶ Le modèle est entraîné en mode **auto-régressif (causal language modeling)**.
- 1. Le décodeur reçoit un contexte (e.g., un prompt) ainsi que le token [START] suivi de la séquence cible en entier.
- 2. Le modèle prédit alors tous les token de la séquence cible, en appliquant une **masked self-attention**, pour que chaque  $i$ -ème token ne voit que les  $i - 1$  tokens précédents.
- 3. Le modèle calcule la loss entre les tokens générés et les tokens cibles.
- 4. Les paramètres sont mis à jour par **backpropagation**.
- ▶ **Avantages:** entraînement parallélisé sur toute la séquence et apprentissage direct du mode auto-régressif.

## DECODER-ONLY: ENTRAÎNEMENT

- Le modèle est entraîné en mode **auto-régressif (causal language modeling)**.
    1. Le décodeur reçoit un contexte (e.g., un prompt) ainsi que le token [START] suivi de la séquence cible en entier.
    2. Le modèle prédit alors tous les token de la séquence cible, en appliquant une **masked self-attention**, pour que chaque  $i$ -ème token ne voit que les  $i - 1$  tokens précédents.
    3. Le modèle calcule la loss entre les tokens générés et les tokens cibles.
    4. Les paramètres sont mis à jour par **backpropagation**.
  - **Avantages:** entraînement parallélisé sur toute la séquence et apprentissage direct du mode auto-régressif.

## DECODER-ONLY: ENTRAÎNEMENT

- Le modèle est entraîné en mode **auto-régressif (causal language modeling)**.
    1. Le décodeur reçoit un contexte (e.g., un prompt) ainsi que le token [START] suivi de la séquence cible en entier.
    2. Le modèle prédit alors tous les token de la séquence cible, en appliquant une **masked self-attention**, pour que chaque  $i$ -ème token ne voit que les  $i - 1$  tokens précédents.
    3. Le modèle calcule la loss entre les tokens générés et les tokens cibles.
    4. Les paramètres sont mis à jour par **backpropagation**.
  - **Avantages:** entraînement parallélisé sur toute la séquence et apprentissage direct du mode auto-régressif.

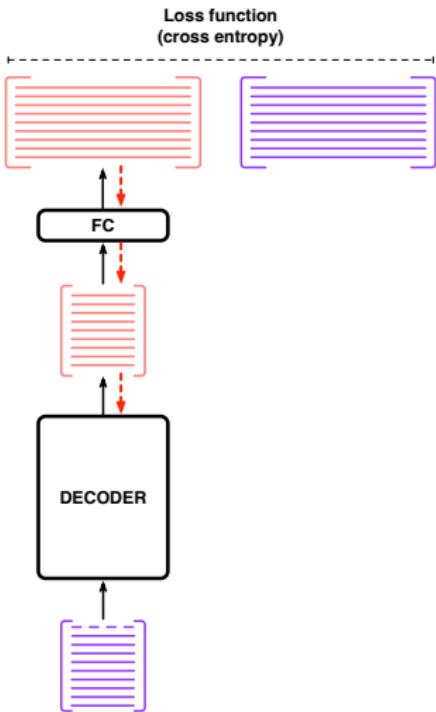
## DECODER-ONLY: ENTRAÎNEMENT

- Le modèle est entraîné en mode **auto-régressif (causal language modeling)**.
    1. Le décodeur reçoit un contexte (e.g., un prompt) ainsi que le token [START] suivi de la séquence cible en entier.
    2. Le modèle prédit alors tous les token de la séquence cible, en appliquant une **masked self-attention**, pour que chaque  $i$ -ème token ne voit que les  $i - 1$  tokens précédents.
    3. Le modèle calcule la loss entre les tokens générés et les tokens cibles.
    4. Les paramètres sont mis à jour par **backpropagation**.
  - **Avantages:** entraînement parallélisé sur toute la séquence et apprentissage direct du mode auto-régressif.

## DECODER-ONLY: ENTRAÎNEMENT

- Le modèle est entraîné en mode **auto-régressif (causal language modeling)**.
    1. Le décodeur reçoit un contexte (e.g., un prompt) ainsi que le token [START] suivi de la séquence cible en entier.
    2. Le modèle prédit alors tous les token de la séquence cible, en appliquant une **masked self-attention**, pour que chaque  $i$ -ème token ne voit que les  $i - 1$  tokens précédents.
    3. Le modèle calcule la loss entre les tokens générés et les tokens cibles.
    4. Les paramètres sont mis à jour par **backpropagation**.
  - **Avantages:** entraînement parallélisé sur toute la séquence et apprentissage direct du mode auto-régressif.

## DECODER-ONLY: ENTRAÎNEMENT



[[START], "Le", "chien", "s'", "enfuit", "."]

## DECODER-ONLY: INFÉRENCE

- ▶ **Mode auto-régressif:** le modèle génère une séquence de texte *pas à pas*, un token à la fois.
    1. À chaque étape, le décodeur reçoit les tokens d'output déjà générés, et prédit le token suivant.
    2. Ce nouveau token est ajouté aux tokens d'output déjà générés, qui sont réinjectés dans le décodeur.
    3. Le processus continue jusqu'à la production du token [END].
  - ▶ **Stratégies de décodage:** greedy decoding, beam search, sampling.

## DECODER-ONLY: INFÉRENCE

- **Mode auto-régressif:** le modèle génère une séquence de texte *pas à pas*, un token à la fois.
    1. À chaque étape, le décodeur reçoit les tokens d'output déjà générés, et prédit le token suivant.
    2. Ce nouveau token est ajouté aux tokens d'output déjà générés, qui sont réinjectés dans le décodeur.
    3. Le processus continue jusqu'à la production du token [END].
  - **Stratégies de décodage:** greedy decoding, beam search, sampling.

## DECODER-ONLY: INFÉRENCE

- **Mode auto-régressif:** le modèle génère une séquence de texte *pas à pas*, un token à la fois.
    1. À chaque étape, le décodeur reçoit les tokens d'output déjà générés, et prédit le token suivant.
    2. Ce nouveau token est ajouté aux tokens d'output déjà générés, qui sont réinjectés dans le décodeur.
    3. Le processus continue jusqu'à la production du token [END].
  - **Stratégies de décodage:** greedy decoding, beam search, sampling.

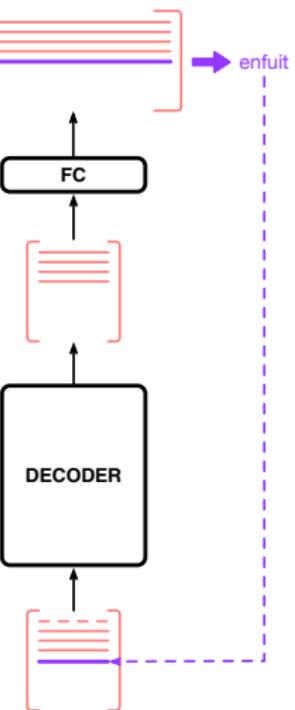
## DECODER-ONLY: INFÉRENCE

- ▶ **Mode auto-régressif:** le modèle génère une séquence de texte *pas à pas*, un token à la fois.
    1. À chaque étape, le décodeur reçoit les tokens d'output déjà générés, et prédit le token suivant.
    2. Ce nouveau token est ajouté aux tokens d'output déjà générés, qui sont réinjectés dans le décodeur.
    3. Le processus continue jusqu'à la production du token [END].
  - ▶ **Stratégies de décodage:** greedy decoding, beam search, sampling.

## DECODER-ONLY: INFÉRENCE

- ▶ **Mode auto-régressif:** le modèle génère une séquence de texte *pas à pas*, un token à la fois.
    1. À chaque étape, le décodeur reçoit les tokens d'output déjà générés, et prédit le token suivant.
    2. Ce nouveau token est ajouté aux tokens d'output déjà générés, qui sont réinjectés dans le décodeur.
    3. Le processus continue jusqu'à la production du token [END].
  - ▶ **Stratégies de décodage:** greedy decoding, beam search, sampling.

## DECODER-ONLY: INFÉRENCE



[[START], “Le”, “chien”, “s”, “enfuit”]

# TRAINING TASKS: PRE-TRAINING

## (1) Causal Language Modeling (Causal LM)

- ▶ Tâche principale pour entraîner les LLMs *génératifs* de type *decoder-only*, tels que GPT.
- ▶ **Objectif:** “Next-token prediction”: étant donnée une suite de tokens  $x_{<t} = (x_0, x_1, \dots, x_{t-1})$ , il faut prédire le token suivant  $x_t$ . On n’utilise que le *contexte uni-directionnel* du passé: *mode auto-régressif*.
- ▶ Formellement, le modèle apprend la distribution de langage:

$$P(x_0, x_1, \dots, x_n) = \prod_{t=0}^n P(x_t | x_{<t})$$

- ▶ **Exemple:** ‘The cat sat on the [MASK]’ → ‘mat’
- ▶ **Loss:** cross-entropy loss  $\mathcal{L} = -\sum_t \log P(x_t | x_{<t})$

# TRAINING TASKS: PRE-TRAINING

## (1) Causal Language Modeling (Causal LM)

- ▶ Tâche principale pour entraîner les LLMs *génératifs* de type *decoder-only*, tels que GPT.
- ▶ Objectif: “Next-token prediction”: étant donnée une suite de tokens  $x_{<t} = (x_0, x_1, \dots, x_{t-1})$ , il faut prédire le token suivant  $x_t$ . On n’utilise que le *contexte uni-directionnel* du passé: *mode auto-régressif*.
- ▶ Formellement, le modèle apprend la distribution de langage:

$$P(x_0, x_1, \dots, x_n) = \prod_{t=0}^n P(x_t | x_{<t})$$

- ▶ Exemple: ‘The cat sat on the [MASK]’ → ‘mat’
- ▶ Loss: cross-entropy loss  $\mathcal{L} = -\sum_t \log P(x_t | x_{<t})$

# TRAINING TASKS: PRE-TRAINING

## (1) Causal Language Modeling (Causal LM)

- ▶ Tâche principale pour entraîner les LLMs *génératifs* de type *decoder-only*, tels que GPT.
- ▶ **Objectif:** “Next-token prediction”: étant donnée une suite de tokens  $x_{<t} = (x_0, x_1, \dots, x_{t-1})$ , il faut prédire le token suivant  $x_t$ . On n'utilise que le *contexte uni-directionnel* du passé: *mode auto-régressif*.
- ▶ Formellement, le modèle apprend la distribution de langage:

$$P(x_0, x_1, \dots, x_n) = \prod_{t=0}^n P(x_t | x_{<t})$$

- ▶ **Exemple:** ‘The cat sat on the [MASK]’ → ‘mat’
- ▶ **Loss:** cross-entropy loss  $\mathcal{L} = -\sum_t \log P(x_t | x_{<t})$

## TRAINING TASKS: PRE-TRAINING

## (1) Causal Language Modeling (Causal LM)

- ▶ Tâche principale pour entraîner les LLMs *génératifs* de type *decoder-only*, tels que GPT.
  - ▶ **Objectif:** “Next-token prediction”: étant donnée une suite de tokens  $x_{<t} = (x_0, x_1, \dots, x_{t-1})$ , il faut prédire le token suivant  $x_t$ . On n’utilise que le *contexte uni-directionnel* du passé: *mode auto-régressif*.
  - ▶ Formellement, le modèle apprend la distribution de langage:

$$P(x_0, x_1, \dots, x_n) = \prod_{t=0}^n P(x_t \mid x_{<t})$$

# TRAINING TASKS: PRE-TRAINING

## (1) Causal Language Modeling (Causal LM)

- ▶ Tâche principale pour entraîner les LLMs *génératifs* de type *decoder-only*, tels que GPT.
- ▶ **Objectif:** “Next-token prediction”: étant donnée une suite de tokens  $x_{<t} = (x_0, x_1, \dots, x_{t-1})$ , il faut prédire le token suivant  $x_t$ . On n'utilise que le *contexte uni-directionnel* du passé: *mode auto-régressif*.
- ▶ Formellement, le modèle apprend la distribution de langage:

$$P(x_0, x_1, \dots, x_n) = \prod_{t=0}^n P(x_t | x_{<t})$$

- ▶ **Exemple:** ‘The cat sat on the [MASK]’ → ‘mat’
- ▶ **Loss:** cross-entropy loss  $\mathcal{L} = -\sum_t \log P(x_t | x_{<t})$

## TRAINING TASKS: PRE-TRAINING

## (1) Causal Language Modeling (Causal LM)

- ▶ Tâche principale pour entraîner les LLMs *génératifs* de type *decoder-only*, tels que GPT.
  - ▶ **Objectif:** “Next-token prediction”: étant donnée une suite de tokens  $x_{<t} = (x_0, x_1, \dots, x_{t-1})$ , il faut prédire le token suivant  $x_t$ . On n’utilise que le *contexte uni-directionnel* du passé: *mode auto-régressif*.
  - ▶ Formellement, le modèle apprend la distribution de langage:

$$P(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{t=0}^n P(\mathbf{x}_t \mid \mathbf{x}_{<t})$$

- ▶ **Exemple:** ‘The cat sat on the [MASK]’  $\rightarrow$  ‘mat’
  - ▶ **Loss:** cross-entropy loss  $\mathcal{L} = -\sum_t \log P(\mathbf{x}_t \mid \mathbf{x}_{<t})$ .

## TRAINING TASKS: PRE-TRAINING

## (2) Mask-Language Modeling (MLM)

## TRAINING TASKS: PRE-TRAINING

## (2) Mask-Language Modeling (MLM)

- ▶ Pour entraîner les LLMs *non génératifs* de type *encoder-only*, tels que BERT, RoBERTa, etc.
  - ▶ Objectif: étant donnée une suite de tokens, on masque des tokens au hasard (15%) et on demande de prédire ces tokens en utilisant le *contexte bi-directionnel* (passé et futur).
  - ▶ Exemple: ‘The [MASK] sat on the [MASK]’ → ‘cat’, ‘mat’
  - ▶ Loss: cross-entropy loss

## TRAINING TASKS: PRE-TRAINING

## (2) Mask-Language Modeling (MLM)

- ▶ Pour entraîner les LLMs *non génératifs* de type *encoder-only*, tels que BERT, RoBERTa, etc.
  - ▶ **Objectif:** étant donnée une suite de tokens, on masque des tokens au hasard (15%) et on demande de prédire ces tokens en utilisant le *contexte bi-directionnel* (passé et futur).
  - ▶ **Exemple:** ‘The [MASK] sat on the [MASK]’ → ‘cat’, ‘mat’
  - ▶ **Loss:** cross-entropy loss

## TRAINING TASKS: PRE-TRAINING

## (2) Mask-Language Modeling (MLM)

- ▶ Pour entraîner les LLMs *non génératifs* de type *encoder-only* , tels que BERT, RoBERTa, etc.
  - ▶ **Objectif:** étant donnée une suite de tokens, on masque des tokens au hasard (15%) et on demande de prédire ces tokens en utilisant le *contexte bi-directionnel* (passé et futur).
  - ▶ **Exemple:** ‘The [MASK] sat on the [MASK]’ → ‘cat’, ‘mat’
  - ▶ **Loss:** cross-entropy loss

## TRAINING TASKS: PRE-TRAINING

## (2) Mask-Language Modeling (MLM)

- ▶ Pour entraîner les LLMs *non génératifs* de type *encoder-only* , tels que BERT, RoBERTa, etc.
  - ▶ **Objectif:** étant donnée une suite de tokens, on masque des tokens au hasard (15%) et on demande de prédire ces tokens en utilisant le *contexte bi-directionnel* (passé et futur).
  - ▶ **Exemple:** ‘The [MASK] sat on the [MASK]’ → ‘cat’, ‘mat’
  - ▶ **Loss:** cross-entropy loss

# TRAINING TASKS: FINE-TUNING

## (3) Fine-tuning supervisé (supervised fine-tuning, SFT)

- ▶ Pour les modèles de type ChatGPT.
  - ▶ Objectif: enseigner au modèle à suivre des instructions humaines et à générer des réponses utiles et polies.
- Données: paires de *prompts* et *réponses humaines*.

## TRAINING TASKS: FINE-TUNING

### (3) Fine-tuning supervisé (supervised fine-tuning, SFT)

- ▶ Pour les modèles de type ChatGPT.
  - ▶ Objectif: enseigner au modèle à suivre des instructions humaines et à générer des réponses utiles et polies.

Données: paires de *prompts* et *réponses humaines*.

## TRAINING TASKS: FINE-TUNING

### (3) Fine-tuning supervisé (supervised fine-tuning, SFT)

- ▶ Pour les modèles de type ChatGPT.
  - ▶ **Objectif:** enseigner au modèle à suivre des instructions humaines et à générer des réponses utiles et polies.

**Données:** paires de *prompts* et *réponses humaines*.

## TRAINING TASKS: FINE-TUNING

## (4) Reinforcement Learning from Human Feedback (RLHF)

## TRAINING TASKS: FINE-TUNING

## (4) Reinforcement Learning from Human Feedback (RLHF)

- ▶ Pour les modèles de type ChatGPT.

## TRAINING TASKS: FINE-TUNING

## (4) Reinforcement Learning from Human Feedback (RLHF)

- ▶ Pour les modèles de type ChatGPT.
  - ▶ **Objectif:** Pour un prompt donné:
    - ▶ on collect multiple de responses possibles du modèle;
    - ▶ on les classifie de manière humaine;
    - ▶ on fine-tune le modèle par renforcement (RL with PPO/DPO) à prédire les réponses préférées.

## TRAINING TASKS: FINE-TUNING

## (4) Reinforcement Learning from Human Feedback (RLHF)

- ▶ Pour les modèles de type ChatGPT.
  - ▶ **Objectif:** Pour un prompt donné:
    - ▶ on collect multiple de responses possibles du modèle;
    - ▶ on les classifie de manière humaine;
    - ▶ on fine-tune le modèle par renforcement (RL with PPO/DPO) à prédire les responses préférées.

## TRAINING TASKS: FINE-TUNING

## (4) Reinforcement Learning from Human Feedback (RLHF)

- ▶ Pour les modèles de type ChatGPT.
  - ▶ **Objectif:** Pour un prompt donné:
    - ▶ on collect multiple de responses possibles du modèle;
    - ▶ on les classifie de manière humaine;
    - ▶ on fine-tune le modèle par renforcement (RL with PPO/DPO) à prédire les réponses préférées.

## TRAINING TASKS: FINE-TUNING

## (4) Reinforcement Learning from Human Feedback (RLHF)

- ▶ Pour les modèles de type ChatGPT.
  - ▶ **Objectif:** Pour un prompt donné:
    - ▶ on collect multiple de responses possibles du modèle;
    - ▶ on les classifie de manière humaine;
    - ▶ on fine-tune le modèle par renforcement (RL with PPO/DPO) à prédire les réponses préférées.

# BIBLIOGRAPHIE

-  **Alammar, J. (2018).**  
The illustrated transformer.
-  **Fleuret, F. (2022).**  
Deep Learning Course.
-  **Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017).**  
**Attention is all you need.**  
In Guyon, I., von Luxburg, U., Bengio, S., Wallach, H. M., Fergus, R., Vishwanathan, S. V. N., and Garnett, R., editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.