

Fiche Technique : L'Architecture Transformer

1. L'Entree (Input)

Avant même d'entrer dans les couches, le texte subit deux transformations :

- Embedding (Plongement) : Chaque mot est transformé en un vecteur numérique dense. C'est le 'sens' brut du mot.
- Positional Encoding : On ajoute un motif mathématique aux vecteurs pour donner l'ordre des mots (1er, 2eme...).

2. L'Encodeur (Le Lecteur)

Rôle : Analyser le contexte et comprendre les relations entre les mots.

A. Self-Attention (Auto-Attention)

Chaque mot se demande : 'Quelle importance les autres mots ont-ils pour moi ?'

Mécanisme Q, K, V (Query, Key, Value). Exemple : Dans 'La pomme est rouge', 'rouge' porte attention sur 'pomme'.

B. Add & Norm

Connexion Residuelle : On ajoute l'entrée à la sortie pour que l'info circule bien.

Normalisation : On stabilise les chiffres.

C. Feed-Forward Network

Un petit réseau de neurones classique pour digérer l'information.

3. Le Décodeur (L'Ecrivain)

Rôle : Générer la séquence de sortie mot par mot.

A. Masked Self-Attention (Attention Masquée)

Interdiction de voir le futur. On applique un masque pour cacher les mots qui n'ont pas encore été écrits.

B. Cross-Attention (Attention Croisée)

Le pont entre l'Encodeur et le Décodeur. Le Décodeur regarde la phrase source complète (Encodeur) pour savoir quoi écrire.

C. Sortie (Linear + Softmax)

Transforme le vecteur final en probabilités pour choisir le mot le plus probable dans le dictionnaire.

Résumé des différences

- Encodeur seul (ex: BERT) : Pour comprendre, classer, extraire (SQuAD).
- Décodeur seul (ex: GPT) : Pour générer du texte, chater.
- Encodeur-Décodeur (ex: T5) : Pour traduire, résumer.

Ton projet actuel : Tu as fait du 'From Scratch' (Encodeur) pour de l'Extraction de réponse.