

# Transformers & Hugging Face

Jour 4 — Après-midi

Julien Rolland

Formation M2 Développement Fullstack

Jour 4

- 1 Du contexte à l'attention
- 2 L'architecture Transformer
- 3 Hugging Face en pratique
- 4 Évaluation

## La polysémie en pratique

- « La **souris** mange le fromage. »
- « Ma **souris** d'ordinateur est cassée. »

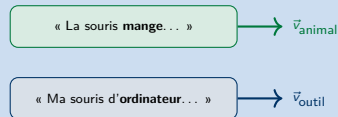
Dans Word2Vec, « souris » n'a **qu'un seul vecteur** — une moyenne des deux sens.

## Les autres limites (rappel J4-AM)

- Embeddings **figés** après entraînement
- Ignorent l'**ordre** des mots dans la phrase
- Biais hérités du corpus d'entraînement

## Le besoin : embeddings contextuels

Le vecteur d'un mot doit être calculé **à la volée** en fonction de ses voisins.



## Solution

Le **mécanisme d'attention** — cœur du Transformer.

## Principe & exemple

Chaque mot « vote » pour ceux qui l'aident à se définir.

« L'animal ... car **il** était fatigué. »

Distribution d'attention de « il » :

- **animal** : 90 %
- fatigué : 6 %
- autres mots : 4 %

## Avantage : parallélisme

Contrairement aux RNN, l'attention traite **tous les mots simultanément** — GPU.

## Distribution d'attention de « il »

**90 % animal** — co-référence directe

**6 % fatigué** — trait partagé

**4 % autres mots** (négligeable)

## Interprétation

Le modèle « sait » que « il » réfère à « animal » **grâce au contexte** global — sans règle grammaticale explicite.

## Mécanisme Query / Key / Value

- $Q$  (Query) : « Ce que je **cherche**. »
- $K$  (Key) : « Ce que j'**offre**. »
- $V$  (Value) : « L'information que je **contiens**. »

## Équation

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

- $QK^T$  : score de pertinence (mot, mot)
- $\sqrt{d_k}$  : normalisation (gradients stables)
- Softmax  $\rightarrow$  poids  $\in [0, 1]$ , somme = 1

## Multi-Head Attention

On lance  $h$  attentions en **parallèle** sur des sous-espaces distincts :

- Tête 1 : relations **grammaticales**
- Tête 2 : relations **sémantiques**
- Tête 3 : relations de **coréférence**
- ...

Les sorties sont **concaténées** puis projetées.

## Note « Geek »

C'est une **recherche floue** dans une base de données vectorielle :  $Q$  est la requête,  $K$  l'index,  $V$  le résultat.

## Encoder (ex: BERT)

- + Lit la phrase **dans les deux sens**
- + Embeddings contextuels riches
- Idéal : **classification**, NER, question-answering

BERT, RoBERTa

Encoder  
Self-Attention (bidirectionnel)

Classification, NER

## Decoder (ex: GPT)

- + Génère mot après mot (*autoregressive*)
- Attention **masquée** : ne voit que le passé
- Idéal : **génération** de texte, chatbots

GPT-2/3/4, LLaMA

Decoder  
Masked Self-Attention

Génération

## Encoder-Decoder (ex: T5)

- + Encode la source, décode la cible
- Idéal : **traduction**, résumé, reformulation

T5, BART, mBART

Encoder + Decoder  
Cross-Attention

Traduction, Résumé

## Pre-training — la fondation

- Entraînement massif sur **tout Internet** (Wikipédia, livres, code...)
- Objectif : **prédire le token masqué** (BERT) ou le token suivant (GPT)
- Coût : millions d'euros de calcul GPU
- Résultat : le modèle apprend la « **langue** » — grammaire, faits, raisonnement

## Fine-tuning — l'adaptation

On ajuste le modèle sur une **tâche métier** :

- Classification de tickets support
- Résumé automatique de rapports
- Détection de sentiment client

**Quelques dizaines de lignes** de code, quelques **minutes** de calcul → niveau état de l'art.

## Pourquoi ça marche ?

Les couches inférieures encodent des connaissances **générales** réutilisables ; seule la tête de classification change.

## The Hub

Le « GitHub » des modèles d'IA :

- > 100 000 modèles prêts à l'emploi
- Hébergement gratuit des poids
- Versionning, fiches modèles, benchmarks intégrés

## transformers

Interface **unifiée** pour tout modèle :

- BERT, RoBERTa, CamemBERT
- GPT-2, LLaMA, Mistral
- T5, BART, Whisper, CLIP...

Même API quel que soit l'architecture.

## datasets

- Accès instantané à des milliers de bases de données (IMDb, Wikipedia, GLUE, SQuAD...)
- Chargement en une ligne :  
`load_dataset("imdb")`

## Écosystème complet

- tokenizers : tokeniseurs rapides (Rust)
- evaluate : métriques standardisées
- accelerate : multi-GPU/TPU
- PEFT : fine-tuning efficace (LoRA)



## Concept

Abstraction totale : téléchargement du modèle, tokenisation, inférence, décodage — tout en **une seule ligne**.

```
1 from transformers import pipeline
2
3 clf = pipeline("sentiment-analysis")
4 print(clf("Ce cours est dense !"))
5 # [{'label': 'POSITIVE', 'score': 0.98}]
6
7 trad = pipeline("translation",
8                 model="Helsinki-NLP/opus-mt-fr-en")
9 print(trad("Bonjour le monde"))
10 # [{'translation_text': 'Hello World'}]
```

## Tâches disponibles

- "sentiment-analysis"
- "text-generation"
- "summarization"
- "translation"
- "question-answering"
- "ner" (entités nommées)

## Intérêt terrain

C'est ce qu'un dev Fullstack utilisera **90 % du temps** en entreprise. Pas besoin de connaître l'architecture interne.

## Les 3 composants obligatoires

- 1 **Tokenizer** : identique au modèle pré-entraîné
- 2 **Model Head** : couche de sortie remplacée (ex: 30 000 → 2 classes)
- 3 **Trainer** : boucle d'entraînement optimisée

## Pourquoi le Tokenizer doit coïncider ?

Le modèle a appris des embeddings pour des IDs précis. Si l'on utilise un tokenizer différent, les IDs ne correspondent plus aux poids — le modèle produit du bruit.

## Alternatives légères

- **LoRA / QLoRA** : < 1 % des params
- **Prompt-tuning** : tokens « soft »
- **Zero-shot** : sans fine-tuning

```
1 from transformers import (  
2     AutoTokenizer,  
3     AutoModelForSequenceClassification,  
4 )  
5 model = AutoModelForSequenceClassification\  
6     .from_pretrained("camembert-base",  
7                     num_labels=2)  
8 tokenizer = AutoTokenizer\  
9     .from_pretrained("camembert-base")
```

## Pourquoi l'accuracy ment

99 % de mails normaux, 1 % de spams.

Prédire « normal » partout → **99 %** — mais le modèle est **inutile**.

## Précision & Rappel

$$\text{Précision} = \frac{TP}{TP+FP} \quad \text{Rappel} = \frac{TP}{TP+FN}$$

## F1-Score — moyenne harmonique

$$F_1 = 2 \cdot \frac{\text{Précision} \times \text{Rappel}}{\text{Précision} + \text{Rappel}}$$

## Matrice de confusion

Visualiser **où** le modèle se trompe :

		Prédit +	Prédit -
Réal	+	TP vrais +	FP faux +
	-	FN faux -	TN vrais -

## Exemple

Confondre « colère » et « tristesse » est moins grave que « positif » et « négatif » — la matrice le révèle.