


Tokenization vs Embedding : Différence

 29 April
2026Mis à jour le 29 April
202648 min de
lecture

Tokenization vs embedding expliqué : BPE, WordPiece, SentencePiece, tiktoken, Cohere, BGE. Impact RAG, context window, pricing.

La tokenization et l'embedding constituent les deux piliers fondamentaux sur lesquels repose l'IA artificielle moderne appliquée au langage naturel. Chaque mot que vous tapez dans un moteur de recherche ou soumettez à un moteur de recherche sémantique, chaque document que vous indexez passe inévitablement par ces deux étapes de transformation. Pourtant, la confusion entre les deux concepts persiste chez les praticiens de l'IA. La tokenization découpe le texte brut en unités discrètes que l'embedding transforme ces unités en vecteurs numériques dans un espace multidimensionnel où la similarité sémantique est mesurée. Comprendre cette distinction, ainsi que la relation profonde entre les deux, est indispensable pour optimiser les systèmes de RAG (Retrieval-Augmented Generation), gérer efficacement les fenêtres de contexte et construire des pipelines NLP performants. Nous explorerons les algorithmes de tokenization (BPE, WordPiece, SentencePiece, Unigram), les modèles de comptage de tokens avec tiktoken, l'impact sur le RAG et la tarification, et les défis de la tokenization multilingue.

À RETENIR

A retenir : La tokenization est le découpage du texte en unités (tokens). L'embedding transforme les tokens ou de séquences de tokens en vecteurs numériques. La tokenization précède l'embedding dans le pipeline. Le choix du tokenizer affecte directement la qualité des embeddings.

Qu'est-ce que la tokenization et pourquoi est-elle nécessaire ?

Les modèles de langage, qu'ils soient des LLM génératifs ou des modèles d'embedding, nécessitent une représentation numérique du texte. Un ordinateur ne comprend que des nombres. La tokenization est le processus qui convertit une séquence de caractères en une séquence d'entiers, chacun correspondant à une entrée dans un vocabulaire prédéfini. C'est une étape essentielle du pipeline de traitement du langage naturel.

L'histoire de la tokenization en NLP a connu plusieurs approches successives. La plus simple est la tokenization par mots (word-level), qui découpe le texte aux espaces et à la ponctuation. Simple mais problématique : le vocabulaire est limité (des milliers de mots dans une langue comme le français), les mots rares ou inconnus sont traités comme des OOV (Out-Of-Vocabulary), et les variations morphologiques (manger, mangeon, mangent) sont traitées comme des entités complètement différentes.

La tokenization par caractères (character-level) résout le problème de vocabulaire (tous les caractères sont dans le vocabulaire) mais perd toute notion de sens au niveau du token et produit des séquences très longues. C'est pourquoi elle n'est pas utilisée pour les modèles Transformer dont la complexité est quadratique en fonction de la longueur de la séquence.

La tokenization en sous-mots (subword tokenization) est le compromis qui a dominé. Elle découpe les mots fréquents en tokens complets et les mots rares en sous-unités. Par exemple, "anticonstitutionnellement" pourrait être découpé en ["anti", "constitu", "tion", "nellement"]. Elle permet un vocabulaire de taille gérée (typiquement 32 000 à 128 000 tokens), une couverture étendue et une préservation partielle de l'information morphologique.
