



# Tokenizer

## Learn about language model tokenization

OpenAI's large language models (sometimes referred to as GPT's) process text using **tokens**, which are common sequences of characters found in a set of text. The models learn to understand the statistical relationships between these tokens, and excel at producing the next token in a sequence of tokens.

You can use the tool below to understand how a piece of text might be tokenized by a language model, and the total count of tokens in that piece of text.

It's important to note that the exact tokenization process varies between models. Newer models like GPT-3.5 and GPT-4 use a different tokenizer than previous models, and will produce different tokens for the same input text.

**GPT-3.5 & GPT-4**   **GPT-3 (Legacy)**

- Cliente2: Multi-canal (Backoffice 40%, Caja 40%, Comercial 20%)
- Cliente3: Caja (Backoffice 20%, Caja 60%, Comercial 20%)
- Cliente4: Comercial (Backoffice 20%, Caja 10%, Comercial 70%)
- Cliente5: Comercial (Backoffice 27%, Caja 27%, Comercial 45%)

Espero que esta información resumida te sea de ayuda. ¡Si necesitas más clarificaciones o tienes alguna otra pregunta, no dudes en decírmelo!

Clear

Show example

Tokens

Characters

172

554

assistant:Según las reglas definidas, la clasificación de los clientes sería la siguiente:

- Cliente1: Multi-canal (Backoffice 42%, Caja 33%, Comercial 25%)
- Cliente2: Multi-canal (Backoffice 40%, Caja 40%, Comercial 20%)



---

- Clientes: Comercial (Backoffice 21%, Caja 21%, Comercial 43%)

Espero que esta información resumida te sea de ayuda. ¡Si necesitas más clarificaciones o tienes alguna otra pregunta, no dudes en decírmelo!

Text   Token IDs

A helpful rule of thumb is that one token generally corresponds to ~4 characters of text for common English text. This translates to roughly  $\frac{3}{4}$  of a word (so 100 tokens  $\approx$  75 words).

If you need a programmatic interface for tokenizing text, check out our [tiktoken](#) package for Python. For JavaScript, the community-supported [@dbdq/tiktoken](#) package works with most GPT models.