



# 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)**

```
respuesta = response['choices'][0]['message']['content'].strip()
print ( "Los Clientes de la Sucursal X se pueden definir como: ",
respuesta )
```

Clear

Show example

**Tokens**

**Characters**

**35**

**143**

```
respuesta = response['choices'][0]['message']['content'].strip()
print ( "Los Clientes de la Sucursal X se pueden definir como: ",
respuesta )
```



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.