

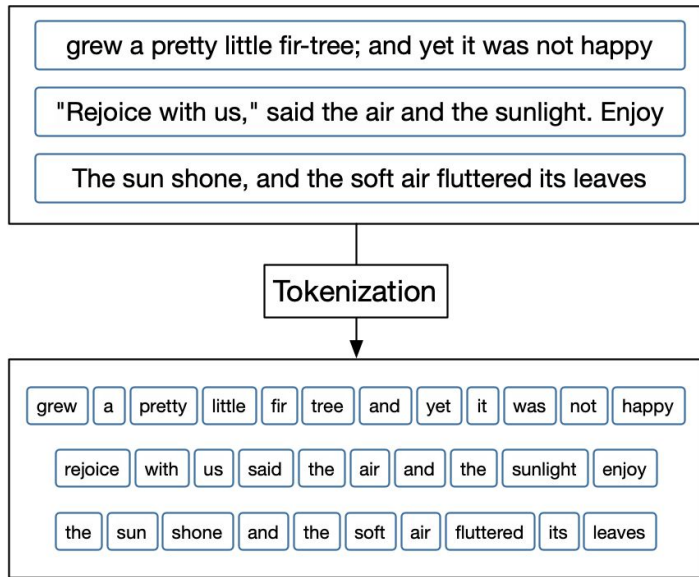
What are Tokens?

Examine token limits what are and how to gracefully handle limits.

What Is Tokenization?

— — —

- Sentences are transformed into lots of individual tokens.



Token Limits

— — —

- Restricted by the models that we choose to use.
- Each model has a unique token limit.

3,569 tokens in prompt

Up to 1,563 tokens in response

This model can only process a maximum of 4,001 tokens in a single request, please reduce your prompt or response length.

[Learn more about pricing](#)

Standard Token Limits

— — —

LATEST MODEL	DESCRIPTION	MAX REQUEST	TRAINING DATA
text-davinci-003	Most capable GPT-3 model. Can do any task the other models can do, often with higher quality, longer output and better instruction-following. Also supports inserting completions within text.	4,000 tokens	Up to Jun 2021
text-curie-001	Very capable, but faster and lower cost than Davinci.	2,048 tokens	Up to Oct 2019
text-babbage-001	Capable of straightforward tasks, very fast, and lower cost.	2,048 tokens	Up to Oct 2019
text-ada-001	Capable of very simple tasks, usually the fastest model in the GPT-3 series, and lowest cost.	2,048 tokens	Up to Oct 2019

How to get the token limit for ChatGPT and GPT-3

— — —

<https://platform.openai.com/tokenizer>

Tokenizer

The GPT family of models process text using **tokens**, which are common sequences of characters found in text. The models 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 would be tokenized by the API, and the total count of tokens in that piece of text.

GPT-3 Codex

This is some text
|

Clear

Show example

Tokens

5

Characters

18

This is some text

How to get the token limit - Automatically

— — —

```
import tiktoken
enc = tiktoken.get_encoding("gpt2")
assert enc.decode(enc.encode("hello world")) == "hello world"

# To get the tokeniser corresponding to a specific model in the OpenAI API:
enc = tiktoken.encoding_for_model("text-davinci-003")
```

Approaches To Avoid Hitting Token Limits

— — —

- Shortening the prompt
- Chunking
- Windowed chunks
- Summarisation