

- Context Window: 모델이 한 번에 처리할 수 있는 최대 입출력 토큰 수를 말합니다.
- max\_tokens: 모델이 답변으로 생성할 수 있는 최대 토큰 수를 말합니다.

## Models

### Flagship models

GPT-4o	GPT-4o mini	o1-preview & o1-mini <span>Beta</span>
Our high-intelligence flagship model for complex, multi-step tasks	Our affordable and intelligent small model for fast, lightweight tasks	A new series of reasoning models for solving hard problems
Text and image input, text output	Text and image input, text output	Text input, text output
128k context length	128k context length	128k context length
Smarter model, higher price per token	Faster model, lower price per token	Higher latency, uses tokens to think

## Context Window

- **Context Window**는 AI 모델이 한 번에 처리할 수 있는 최대 토큰 수를 의미합니다.
- 이 값은 주어진 텍스트 길이를 제한하며, 이 창 안에 들어가는 모든 입력과 출력의 합이 해당 창 크기를 넘지 않아야 합니다.
- 예를 들어, GPT-4의 경우 컨텍스트 윈도우가 8,192 토큰 또는 32,768 토큰 등 다양한 크기로 제공됩니다. 이 창이 클수록 긴 문장이나 문단을 더 잘 이해하고 처리할 수 있습니다.
- 더 큰 Context Window는 긴 대화나 복잡한 명령어를 처리할 때 유리하지만, 일반적으로 모델의 메모리 사용량도 증가하게 됩니다.

## max\_tokens

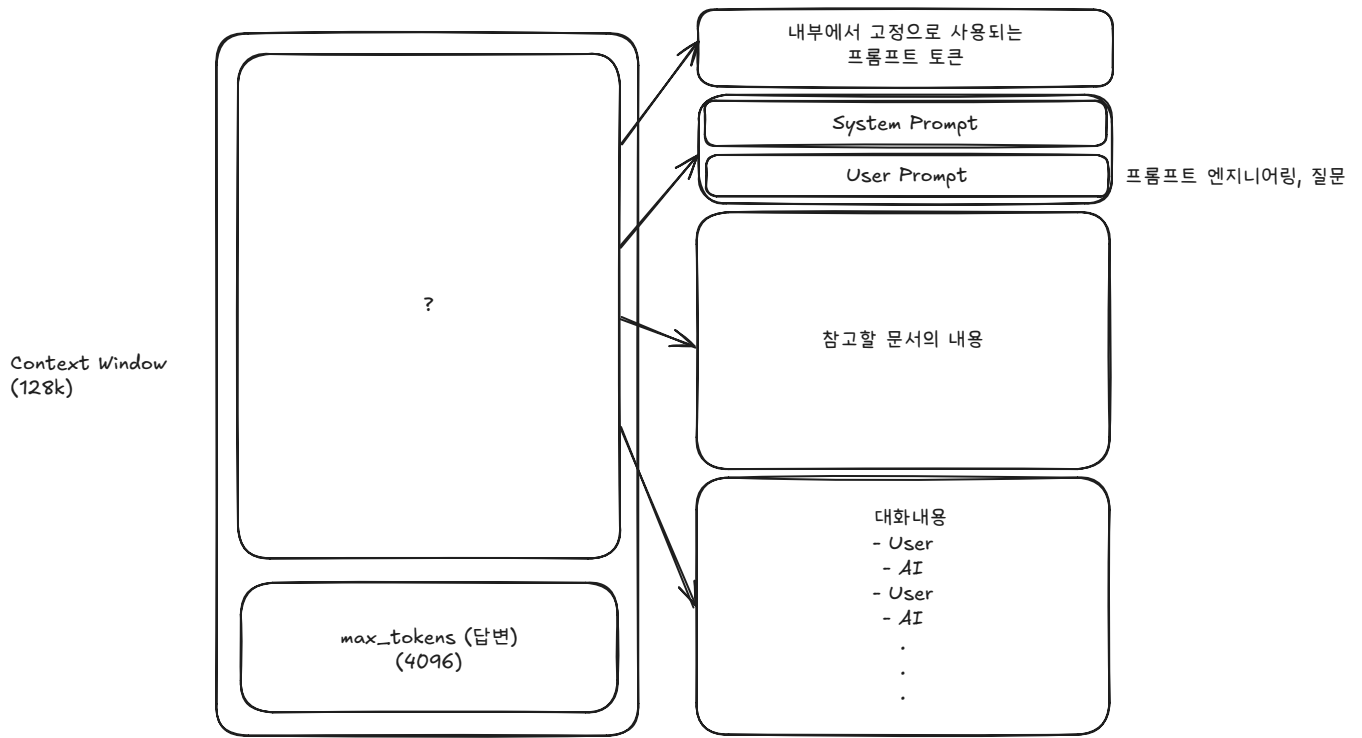
- **max\_tokens**는 모델이 출력으로 생성할 수 있는 최대 토큰 수를 지정하는 설정입니다.
- 예를 들어, max\_tokens를 100으로 설정하면 모델은 최대 100개의 토큰으로만 응답을 생성합니다.
- max\_tokens와 Context Window는 상호작용합니다. 예를 들어, Context Window가 4,000 토큰이고 이미 3,000개의 토큰이 입력으로 사용되었다면, 최대 1,000개의 토큰만 출력으로 생성할 수 있습니다.

# GPT-4o

GPT-4o (“o” for “omni”) is our most advanced GPT model. It is multimodal (accepting text or image inputs and outputting text), and it has the same high intelligence as GPT-4 Turbo but is much more efficient—it generates text 2x faster and is 50% cheaper. Additionally, GPT-4o has the best vision and performance across non-English languages of any of our models. GPT-4o is available in the OpenAI API to paying customers. Learn how to use GPT-4o in our [text generation guide](#).

MODEL	CONTEXT WINDOW	MAX OUTPUT TOKENS	TRAINING DATA
<div><code>gpt-4o</code></div> <div>Our high-intelligence flagship model for complex, multi-step tasks. GPT-4o is cheaper and faster than GPT-4 Turbo. Currently points to <code>gpt-4o-2024-08-06</code>.</div>	128,000 tokens	16,384 tokens	Up to Oct 2023
<div><code>gpt-4o-2024-08-06</code></div> <div>Latest snapshot that supports <a href="#">Structured Outputs</a>. <code>gpt-4o</code> currently points to this version.</div>	128,000 tokens	16,384 tokens	Up to Oct 2023
<div><code>gpt-4o-2024-05-13</code></div> <div>Original <code>gpt-4o</code> snapshot from May 13, 2024.</div>	128,000 tokens	4,096 tokens	Up to Oct 2023
<div><code>chatgpt-4o-latest</code></div> <div>The <code>chatgpt-4o-latest</code> model version continuously points to the version of GPT-4o used in ChatGPT, and is updated frequently, when there are significant changes.</div>	128,000 tokens	16,384 tokens	Up to Oct 2023

## 구조



문서들을 참고해서 넣어줘야 하는 문서의 내용이 많으면 많을 수록 토큰의 비용이 크게 나올 것입니다. RAG 시스템을 만들때도 필요한 부분만 발췌해서 넣어주는 이유입니다.(비용과 성능적으로 유리)