# Model Overview

| Feature | GPT-4o Mini | Gemini 2.0 Flash-Lite |
|---|---|---|
| **Input Context Window**<br>The number of tokens supported by the input context window. | **128K**<br>tokens | **1M**<br>tokens |
| **Maximum Output Tokens**<br>The number of tokens that can be generated by the model in a single request. | **16.4K**<br>tokens | **8,192**<br>tokens |
| **Open Source**<br>Whether the model's code is available for public use. | **No** | **No** |
| **Release Date**<br>When the model was first released. | **July 18, 2024**<br>10 months ago | **December 11, 2024**<br>5 months ago |
| **Knowledge Cut-off Date**<br>When the model's knowledge was last updated. | **October 2023** | **June 2024** |

# Pricing Comparison

Compare costs for input and output tokens between GPT-4o Mini and Gemini 2.0 Flash-Lite.

| Price Type | GPT-4o Mini | Gemini 2.0 Flash-Lite |
|---|---|---|
| **Input**<br>Cost for processing tokens in your prompts | **$0.15**<br>per million tokens | **$0.07**<br>per million tokens |
| **Output**<br>Cost for tokens generated by the model | **$0.60**<br>per million tokens | **$0.30**<br>per million tokens |