



<b>GPT-4.1</b> Flagship coding & instruction model	25 M	6.25 M
<b>GPT-4.1 mini</b> Fast, cost-effective	125 M	31.25 M
<b>GPT-4.1 nano</b> Fastest & cheapest	500 M	125 M
<b>GPT-4o</b> Multimodal flagship	20 M	5 M
<b>GPT-4o mini</b> Lightweight multimodal	333.33 M	83.33 M
<b>GPT-5 nano</b> Fastest for summarization & classification	1 B	125 M
<b>GPT-5 pro</b> Smartest and most precise model	3.33 M	416.67 K
<b>GPT-5.1</b> Best for coding and agentic tasks	40 M	5 M
<b>o3</b> Advanced reasoning	25 M	6.25 M
<b>o4-mini</b> Efficient reasoning	45.45 M	11.36 M



**Anthropic**

Claude Opus, Sonnet, and Haiku model families

MODEL	INPUT TOKENS / 100 CR	OUTPUT TOKENS / 100 CR
<b>Claude Haiku 3</b> Fastest & cheapest	200 M	40 M
<b>Claude Haiku 3.5</b> Efficient	62.5 M	12.5 M
<b>Claude Haiku 4.5</b> Near-frontier performance	50 M	10 M
<b>Claude Opus 4.5</b> Maximum intelligence	10 M	2 M
<b>Claude Sonnet 4</b> Balanced capability	16.67 M	3.33 M
<b>Claude Sonnet 4.5</b> Best for coding	16.67 M	3.33 M

 **Google Gemini**  
Gemini Pro, Flash, and Flash-Lite models

MODEL	INPUT TOKENS / 100 CR	OUTPUT TOKENS / 100 CR
<b>Gemini 2.0 Flash</b> Balanced multimodal	500 M	125 M
<b>Gemini 2.0 Flash-Lite</b> Ultra-low cost	666.67 M	166.67 M
<b>Gemini 2.5 Flash</b> Hybrid reasoning, 1M context	166.67 M	20 M

### Gemini 2.5 Flash-Lite

Most cost-effective

500 M

125 M

### Gemini 2.5 Pro

Coding & complex reasoning

40 M

5 M

### Gemini 3 Pro

Most powerful agentic model

25 M

4.17 M

P

## Perplexity

Sonar search-grounded models

MODEL	INPUT TOKENS / 100 CR	OUTPUT TOKENS / 100 CR
<b>Sonar</b> Lightweight search	50 M	50 M
<b>Sonar Deep Research</b> Exhaustive research reports	25 M	6.25 M
<b>Sonar Pro</b> Deep content understanding	16.67 M	3.33 M
<b>Sonar Reasoning</b> Step-by-step logic	50 M	10 M
<b>Sonar Reasoning Pro</b> Enhanced multi-step reasoning	25 M	6.25 M

Understanding This Rate Card

- Higher-capability models consume more credits per token
- Efficient models stretch your credits further
- All values based on 100 Credits
- GPT-5 mini is the baseline standard model
- **Green values** indicate more tokens than GPT-5 mini
- Rates are subject to change by providers
- Actual usage varies based on conversation patterns