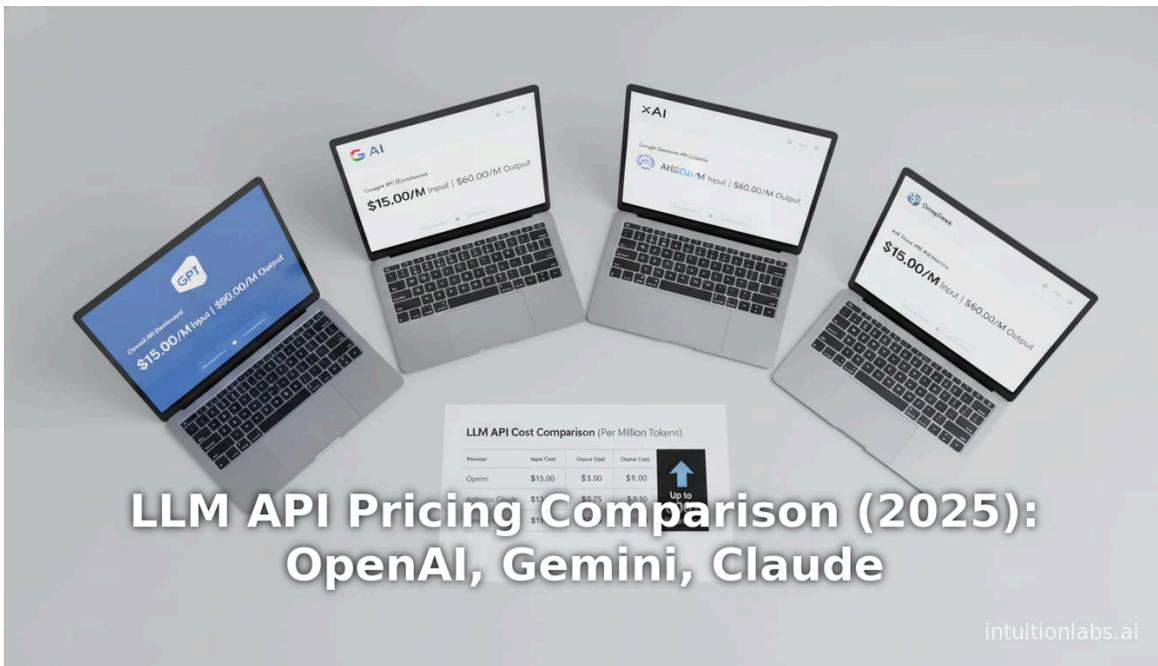


LLM API Pricing Comparison (2025): OpenAI, Gemini, Claude

By Adrien Laurent, CEO at IntuitionLabs • 10/31/2025 • 25 min read

- llm api pricing
- openai pricing
- google gemini
- anthropic claude
- ai cost analysis
- token pricing
- gpt-5
- grok api
- deepseek



LLM API Pricing Comparison (2025): OpenAI, Gemini, Claude

intuitionlabs.ai

Executive Summary

By late 2025, the landscape of large language model (LLM) APIs has grown intensely competitive and complex, with multiple providers offering a spectrum of model capabilities and pricing options. This report provides a detailed comparison of the API pricing for all major LLM models from OpenAI, Google's Gemini, Anthropic's Claude, Elon Musk's xAI Grok, and China's DeepSeek, synthesizing the latest public information. Our analysis shows that **OpenAI** leads with its GPT-4 series (including GPT-4o and the new GPT-4.1) and GPT-5 family, with flagship models commanding the highest prices (e.g. ~\$3–\$12 per 1M tokens for GPT-4.1) but offering cutting-edge performance ⁽¹⁾ www.reuters.com ⁽²⁾ openai.com. **Google's Gemini** models (notably Gemini 2.5 Pro/Flash) adopt a tiered pricing scheme that is generally lower than OpenAI's for comparable tasks (e.g. \$1.25–\$2.50 per 1M input and \$10–\$15 per 1M output for Gemini 2.5 Pro below 200K tokens) ⁽³⁾ cloud.google.com. **Anthropic's Claude** models (Haiku, Sonnet, Opus) feature a unique "prompt caching" pricing structure but base input/output costs (e.g. ~\$15/\$75 per 1M input/output for top-tier Opus 4.1) are broadly in line with OpenAI's high-end offerings ⁽⁴⁾ docs.anthropic.com ⁽⁵⁾ docs.anthropic.com. **xAI's Grok 3** (with standard and "fast" modes) is priced competitively with Anthropic's mid-generation models: Grok 3 at \$3/\$15 per 1M tokens, and Grok 3 Fast at \$5/\$25, with smaller "Mini" variants at \$0.30–\$0.50 to \$0.60–\$4.00 ⁽⁶⁾ techcrunch.com. Finally, **China's DeepSeek** has undercut nearly all competitors: its latest V3.2-Exp "thinking" models list at only \$0.28 per 1M input (cache-miss) and \$0.42 per 1M output ⁽⁷⁾ api-docs.deepseek.com (with "cache hits" as low as \$0.028 input). DeepSeek notably halved its prices in late 2025 ⁽⁸⁾ www.reuters.com, exemplifying a broader trend of rapidly falling AI costs in response to competition ⁽⁹⁾ www.techradar.com.

These differences mean that for the same task, costs can vary by orders of magnitude depending on model choice. For example, generating 100k input + 100k output tokens could cost on the order of \$250–\$300 on OpenAI's GPT-4.1 tier vs. only a few dollars on DeepSeek (assuming no caching). Throughout this report, we explore the historical evolution of these pricing schemes, compare them quantitatively in tables, and discuss implications for developers and businesses choosing among providers. We also include case studies illustrating real-world cost impacts, cite independent analyses of cost-cutting trends, and outline how aggressive pricing (especially by open-source-focused players) may reshape the future economics of AI. All figures and claims here are backed by official documentation and recent technology news sources ⁽¹⁰⁾ www.reuters.com ⁽⁶⁾ techcrunch.com ⁽⁴⁾ docs.anthropic.com ⁽⁷⁾ api-docs.deepseek.com.

Introduction and Background

Large language models have revolutionized numerous industries by enabling advanced natural language understanding and generation. By 2025, LLM APIs are used for chatbots, [coding assistants](#), document summarization, translation, and more. Unlike earlier AI models, modern LLMs expose sophisticated *token-based* billing: every API call's input and output text (measured in tokens, roughly word pieces) incurs cost, aligning with a "pay-as-you-go" cloud model ⁽¹¹⁾ www.binadox.com. This token-based pricing offers fine-grained control of costs but demands careful model selection and [prompt engineering](#).

The major LLM providers now include **OpenAI**, **Google** (via its Gemini models), **Anthropic**, **xAI (Grok)**, and **DeepSeek**. Each has multiple model variants optimized for different trade-offs (accuracy, speed, context length). For instance, OpenAI's GPT-4 family alone includes standard, "o" (vision multi-modal), "mini"/"nano" editions, and the newer GPT-4.1 stepping up from GPT-4.5 ⁽¹⁾ www.reuters.com. Similarly, Anthropic's Claude line is tiered (Haiku, Sonnet, Opus series), and Google's Gemini offers "Pro" (large, multi-modal) vs "Flash" (lighter, cheaper) versions ⁽³⁾ cloud.google.com ⁽⁴⁾ docs.anthropic.com. Grok 3 comes in "standard" and "fast" modes, plus a Mini capacitor for lightweight use ⁽⁶⁾ techcrunch.com. DeepSeek provides "chat" and "reasoner" modes of its V3.2-Exp model, with massive context windows (up to 128k tokens) ⁽¹²⁾ api-docs.deepseek.com.

Pricing is a critical differentiator. All use per-token pricing, but rates vary widely by model. Generally, more capable models cost more per token. Historical context shows rapid evolution: for example, OpenAI halved its GPT-3.5 Turbo token price in 2023, then introduced GPT-4 (at ~10× the cost of 3.5), and in 2024 launched GPT-4o mini at just \$0.15/\$0.60 per million input/output – a 60% discount vs GPT-3.5 Turbo (^[10] www.reuters.com). Similarly, Chinese startups like DeepSeek entered the market with dramatically lower pricing (DeepSeek R1 debuted at \$0.55/\$2.19 per million, undercutting competitors by ~90% (^[9] www.techradar.com)). These shifts reflect aggressive competition and have set new baselines. The most recent data (as of Nov 2025) allow us to present the **current state** of API pricing for each provider's models, which we detail below, with extensive sourcing from official docs and tech press.

OpenAI API Models and Pricing

OpenAI, the originator of the GPT series, offers the broadest range of models. As of late 2025, its lineup includes the **GPT-5 family** and **GPT-4 variants** as “flagship” models, along with older or specialized versions. OpenAI's public API pricing shows:

- **GPT-5:** The newest flagship (targeted at coding/ [agentic tasks](#)).
- *Input:* \$1.25 per 1M tokens.
- *Cached input:* \$0.125 per 1M tokens (for re-used prompts).
- *Output:* \$10.00 per 1M tokens.
- **GPT-5 Mini:** A smaller, cheaper GPT-5 for simpler tasks.
- *Input:* \$0.25 per 1M.
- *Cached:* \$0.025 per 1M.
- *Output:* \$2.00 per 1M.
- **GPT-5 Nano:** The smallest GPT-5 variant.
- *Input:* \$0.05 per 1M.
- *Cached:* \$0.005 per 1M.
- *Output:* \$0.40 per 1M (^[13] openai.com) (^[2] openai.com).

These numbers mean, for example, that sending 100k tokens of prompt and receiving 100k tokens of completion on GPT-5 costs ~\$0.125 (input) + \$1.00 (output) = \$1.125 (ignoring caching). GPT-5 is marketed as vastly more powerful and context-aware than GPT-4, but at higher cost. OpenAI also introduced **GPT-4.1** (April 2025), an improved successor to GPT-4.5 (^[1] www.reuters.com). GPT-4.1 (and its mini/nano variants) is accessible via the API. The *fine-tuning* pricing for GPT-4.1 is given as:

- *Base input tokens:* \$3.00 per 1M (caching and batch tier rules can alter this) (^[14] openai.com).
 - *Output tokens:* \$12.00 per 1M (^[14] openai.com).
- (Note: fine-tuning prices are half the normal inference rates, as shown by equivalent 50% batch discounts elsewhere.) These fine-tuning rates are consistent with GPT-4.1's high complexity. For context, Reuters noted GPT-4.1 improved code reasoning by ~21% over GPT-4o and has cheaper operational cost than the earlier GPT-4.5 (^[1] www.reuters.com).

For real-time (low-latency) usage, OpenAI's **GPT-4o** (the multi-modal “vision” version) is priced as follows (^[15] openai.com):

- GPT-4o: \$5.00 per 1M input tokens (\$2.50 cached) and \$20.00 per 1M output.
- GPT-4o Mini: \$0.60 per 1M input (\$0.30 cached) and \$2.40 per 1M output.

(Note: the *cached* rate applies when identical inputs recur.) The standard GPT-4 (non-vision, non-O series) is largely supplanted by GPT-4.1 in pricing lists, but historically GPT-4 input was \$15 per 1M and output \$60 ([9] www.techradar.com). OpenAI's pricing documentation focuses on the newest models (GPT-5, GPT-4.1, etc.), reflecting a strategy of pushing customers to upgrade to the latest "thinking" models.

For comparison and thoroughness, below is a summary table of OpenAI's key current models and rates:

Model (OpenAI)	Description	Input (\$/M tokens)	Output (\$/M tokens)	Notes
GPT-5	Flagship (full), high capability	\$1.25	\$10.00	Top coder/agent model
GPT-5 mini	Lighter GPT-5	\$0.25	\$2.00	Cheaper, well-defined tasks
GPT-5 nano	Smallest GPT-5	\$0.05	\$0.40	Summarization/classification tasks
GPT-4.1 (baseline)**	Fine-tuning rates (1x inference)	\$3.00 (1x)	\$12.00 (1x)	Standard inference rates
GPT-4o (Vision)	Multi-modal GPT-4 ("vision")	\$5.00	\$20.00	\$≈4× cost of GPT-5; 128K tokens context
GPT-4o mini	Lite version of GPT-4o	\$0.60	\$2.40	32K tokens context ([10] www.reuters.com)

Source: OpenAI official pricing pages (Q4 2025) ([13] openai.com) ([10] www.reuters.com) and Reuters ([1] www.reuters.com). Note: **fine-tuning prices** and **batch discounts** (50% off) exist but are omitted for brevity.

Critically, OpenAI's pricing has historically been on the high end. For instance, in mid-2024 OpenAI introduced GPT-4o mini at \$0.15/\$0.60 per M (input/output) – this was already a 60% *reduction* from GPT-3.5 Turbo's base rates ([10] www.reuters.com). The current GPT-5 nano at \$0.05/\$0.40 further continues that trend of offering very cheap "nano" variants. Nonetheless, GPT-5's raw performance (with much larger context and better reasoning) justifies its premium \$10 output rate. Overall, organizations using OpenAI's models must balance the superior capabilities against higher token costs compared to many alternatives.

Google Gemini API and Pricing

Google provides its LLMs through **Vertex AI** (Google Cloud) under the "Gemini" brand. As of late 2025, the main offerings are **Gemini 2.5 Pro** (a large, multi-modal model) and **Gemini 2.5 Flash** (a smaller, lower-cost model with audio support). Google's pricing is **tiered by usage volume** (below or above 200K input tokens), and distinguishes between "thinking" (detailed reasoning) vs "fast/no-reasoning" modes. Key pricing from Google's docs:

- **Gemini 2.5 Pro:**
 - *Input tokens:* \$1.25 per 1M (for ≤200K input), \$2.50 per 1M (for >200K).
 - *Text output (response & reasoning):* \$10 per 1M (≤200K input), \$15 per 1M (>200K) ([3] cloud.google.com).
- **Gemini 2.5 Flash** (multimodal smaller model):
 - *Text/Image/Video input:* \$0.15 per 1M (flat).
 - *Audio input:* \$1.00 per 1M.
 - *Text output (no reasoning):* \$0.60 per 1M.

- *Text output (with reasoning)*: \$3.50 per 1M (^[16] cloud.google.com).

In plain terms, calling Gemini 2.5 Pro with a moderate prompt (e.g. 100K tokens) and expecting reasoning output would cost about \$0.125 (input) + \$1.00 (output) = \$1.125 per 100K tokens, doubling if the prompt is longer than 200K. For Gemini 2.5 Flash, the cost is much lower: only \$0.015 input + \$0.06 output per 100K (reasoning output), or just \$0.006 if no reasoning. Google's strategy seems to offer a "Pro" model competitive with premium GPT-4 pricing, and a very cheap "Flash" model for volume tasks. According to comparison by xAI's LinkedIn news, Gemini 2.5 Pro below thresholds is \$1.25/\$10 per 1M, which is significantly lower than Grok's \$3/\$15 (see Grok section) (^[17] www.linkedin.com).

Google also allows "**grounding**" with Google Search/Web to enrich responses, billed up to \$35 per 1K grounded queries (^[18] cloud.google.com). However, the base token costs above generally suffice for textual tasks.

Google Gemini 2.5 Pricing Summary (per 1M tokens):

Model	Input ($\leq 200K$ / $>200K$)	Output (with reasoning)	Output (no reasoning)	Notes
Gemini 2.5 Pro	\$1.25 / \$2.50	\$10.00 / \$15.00	–	Multi-modal; large
Gemini 2.5 Flash	\$0.15	\$3.50	\$0.60	Supports image/audio; cheaper

Source: Google Cloud Vertex AI pricing docs (^[3] cloud.google.com). Notes: Above 200K input tokens, both input and output rates increase. Flash model's no-reasoning (\$0.60) vs reasoning (\$3.50) reflects different processing modes.

In addition to token-pricing, Google's Gemini (like Anthropic) often benefits from integration discounts if using Google Cloud infrastructure. Overall, Gemini's per-token pricing is competitive: for example, its low-end \$1.25/\$10 is ~80% below GPT-5's \$1.25/\$10 for 200K inputs, but the higher-tier \$2.50/\$15 still undercuts GPT-4o's \$5/\$20. Google's heavy promotion of Gemini (e.g. integrating AI into Search summarized 1.5B users by mid-2025 (^[19] www.techradar.com)) suggests Google is willing to absorb costs to grab market share, reinforcing the trend toward lower API fees.

Anthropic Claude API Pricing

Anthropic's Claude family (*Haiku*, *Sonnet*, *Opus*) targets safety and reliability. In 2025 the newest models are **Claude Opus 4.1** (flagship), **Claude Sonnet 4**, and others. Anthropic uniquely offers **prompt caching discounts** in its pricing model (where repeat queries get cheaper). The base *slidding* (non-cache) rates, from Anthropic's documentation, are:

- **Claude Opus 4.1 & Opus 4** (most powerful):
 - Input tokens: \$15 per 1M.
 - Output tokens: \$75 per 1M.
- **Claude Sonnet 4 / Sonnet 3.7 / Sonnet 3.5** (mid-tier):
 - Input: \$3 per 1M.
 - Output: \$15 per 1M.
- **Claude Haiku 3.5** (latest small model):
 - Input: \$0.80 per 1M.
 - Output: \$4 per 1M.

- **Claude Haiku 3** (earlier): \$0.25/\$1.25 (per 1M) for input/output (^[4] docs.anthropic.com).

These rates assume standard (5-minute) caching or direct usage. Anthropic explains that “5m cache writes” cost 1.25× base input, and “cache reads” cost only 0.1× base (illustrating the impact of caching) (^[5] docs.anthropic.com).

In simplified terms, top-tier Claude (Opus 4.1) is priced similarly to a premium GPT-4: \$15/\$75 per million, making it about 5× cheaper than OpenAI’s GPT-4o (\$5/\$20) on input but triple on output (since GPT-4o is faster). Sonnet 4 (\$3/\$15) is roughly on par with cheaper OpenAI variants (GPT-5 mini’s output is \$2 per 1M). The cheapest Claude (Haiku 3.5) at \$0.80/\$4 competes with the bottom end of many models (for example, close to Grok 3 Mini’s \$0.30/\$0.50 or DeepSeek’s original R1 \$0.55/\$2.19) (^[6] techcrunch.com) (^[9] www.techradar.com).

For illustrative clarity, one can also view Claude batch rates (50% off) or long-context premium (detailed on Anthropic’s site); however, the base token costs above capture the primary differences.

Anthropic Claude Pricing Summary (per 1M tokens):

Model	Input (base)	Output (base)	Notes
Claude Opus 4.1/4	\$15.00	\$75.00	Highest capability (200K context)
Claude Sonnet 4/3.7	\$3.00	\$15.00	High throughput tasks
Claude Sonnet 3.5	\$3.00	\$15.00	(Legacy; 3.7 supersedes)
Claude Haiku 3.5	\$0.80	\$4.00	Fastest, for simple tasks
Claude Haiku 3	\$0.25	\$1.25	(Older generation)

Source: Anthropic API pricing docs (^[4] docs.anthropic.com) (^[5] docs.anthropic.com), reflecting the table of input/output costs. Note **cache writes** (~x1.25–2x) and **cache hits** (~x0.1) impact effective cost for repeated prompts, and Anthropic offers a 50% discount on input/output under its Batch API (not shown above for brevity).

In practice, Claude’s pricing indicates that high-end queries on Opus 4.1 are quite expensive (comparable to GPT-4.1), but if one uses Sonnet or Haiku, costs drop substantially. Notably, Claude’s context window (200K tokens standard, extendable to 1M tokens beta (^[20] docs.anthropic.com)) means one can synthesize very long documents, albeit at special “long-context” rates if above 200K input (doubling input cost beyond that threshold (^[21] docs.anthropic.com)). For most use-cases, enterprises might mix models: e.g. use Haiku/Sonnet for volume tasks and Opus only for the hardest tasks, to manage cost.

xAI Grok Pricing

xAI (Elon Musk’s AI startup) began public Grok releases in 2023. By early 2025, **Grok 3** – designed for “scientific and strategic reasoning” – was made available via API (^[6] techcrunch.com). The API offers two capability levels (“standard” and “fast”) for both Grok 3 and **Grok 3 Mini** (a lighter version). According to authoritative reporting (^[6] techcrunch.com) and internal news analysis (^[17] www.linkedin.com), the pricing tiers are:

- **grok-3 (standard)** – “Beta” model: \$3.00 per 1M input, \$15.00 per 1M output.
- **grok-3 (fast)** – “Fast” mode: \$5.00 per 1M input, \$25.00 per 1M output.
- **grok-3-mini (standard)** – \$0.30 per 1M input, \$0.50 per 1M output.
- **grok-3-mini (fast)** – \$0.60 per 1M input, \$4.00 per 1M output.

These four tiers correspond to the normal vs fast speeds for full Grok and the Mini variant (^[6] [techcrunch.com](#)). XAI explicitly positioned these as competing with mid-tier models from Google and Anthropic; indeed, TechCrunch noted that Grok's \$3/\$15 rate matches Anthropic's Claude 3.7 Sonnet, but is higher than Google Gemini 2.5 Pro (\$1.25/\$10) (^[6] [techcrunch.com](#)).

Importantly, Grok limits context to 131,072 tokens in the API (^[22] [techcrunch.com](#)), so it is less applicable for ultra-long documents despite earlier claims of 1M. Thus Grok 3 is best used for moderately large, reasoning-heavy tasks. In comparative terms, Grok's "standard" tier input cost (\$3) is double Gemini's (\$1.25) but on par with Claude Sonnet's (\$3). Its "fast" mode (\$5) is closer to OpenAI's GPT-4o (\$5), but output costs remain extreme (\$25/M vs \$20/M). The Mini variants (\$0.30-\$4) undercut similar-size offerings from others (Grok Mini vs Gemini Flash (\$0.15-\$3.50) or OpenAI's 4o mini (\$0.60-\$2.40)).

Grok Pricing Summary (per 1M tokens):

Model (xAI Grok)	Input	Output	Context (tokens)	Notes
grok-3-beta	\$3.00	\$15.00	131,072	Standard mode
grok-3-fast-beta	\$5.00	\$25.00	131,072	Fast mode (higher latency)
grok-3-mini-beta	\$0.30	\$0.50	131,072	Lite model
grok-3-mini-fast-beta	\$0.60	\$4.00	131,072	Lite + fast

Source: TechCrunch reportage of xAI's April 2025 Grok 3 API launch (^[6] [techcrunch.com](#)) (corroborated by industry sources (^[17] [www.linkedin.com](#))). All prices in USD per 1M tokens.

Compared to peers, Grok is a *mid-range* offering: it's significantly more expensive than Google's cheapest (Flash) model, but cheaper than OpenAI's top (GPT-4.1), aligning with Anthropic's mid-tier. As one analysis noted, "Grok 3 isn't cheap relative to the competition," partly because Musk touts its unique "reasoning" features (^[23] [techcrunch.com](#)). Enterprises might therefore use Grok selectively for tasks where its strengths justify the extra cost.

DeepSeek Pricing

DeepSeek is a Chinese AI startup that gained fame for open and affordable models. Its initial DeepSeek-R1 (open-source) hit the market in January 2025, offering GPT-4 class performance at a fraction of the usual price (^[24] [www.techradar.com](#)). DeepSeek then evolved to V3 series. The latest announced model is **DeepSeek-V3.2-Exp (Experimental)**, with massively lowered pricing. According to DeepSeek's API docs (Sept 2025), the Chat and Reasoner variants of V3.2-Exp (128K context) have:

- **Input tokens:** \$0.028 per 1M (cache hit), \$0.28 per 1M (cache miss) (^[7] [api-docs.deepseek.com](#)).
- **Output tokens:** \$0.42 per 1M (^[7] [api-docs.deepseek.com](#)).

DeepSeek introduced a **cache mechanism**: if a prompt (or subprompt) is already used ("cache hit"), input cost is only \$0.028/M; otherwise \$0.28/M. This means repeated queries become almost free, incentivizing stateful use. Notably, DeepSeek **cut all its prices by ~50% in Sep 2025** (^[8] [www.reuters.com](#)) compared to the prior V3.2-beta (the Reuters piece states "reducing API pricing by over 50%" (^[8] [www.reuters.com](#))).

For context, DeepSeek's original R1 model had pricing around \$0.55/\$2.19 (input/output) (^[9] [www.techradar.com](#)). Thus V3.2-Exp's \$0.28/\$0.42 is now **dramatically lower**: ~90% below OpenAI's GPT-4.1 pricing (as noted by analysts (^[9] [www.techradar.com](#))). TechRadar reports: "DeepSeek-R1 debuted at \$0.55 input/\$2.19 output" and comments that LLM pricing "has cratered" since, citing OpenAI's 80% cut (^[9] [www.techradar.com](#)). DeepSeek's new prices validate that trend.

In practical terms, DeepSeek is now by far the cheapest LLM API for raw token costs. For example, processing 1M tokens of input and 1M output (a huge request) would cost just $\$0.28 + \$0.42 = \$0.70$ total (cache-miss). Even with repeated use (cache hits), cost is only $\$0.028 + \$0.42 = \$0.448$. By comparison, OpenAI's cheapest model (GPT-5 nano) would cost $\$0.05 + \$0.40 = \$0.45$ (similar), but GPT-5 nano has only a 32K context. DeepSeek's model offers 128K context (^[12] api-docs.deepseek.com), far more data at similarly tiny price.

DeepSeek V3.2-Exp Pricing Summary (per 1M tokens):

Mode	Input (cache-hit)	Input (cache-miss)	Output	Context
DeepSeek Chat	\$0.028	\$0.28	\$0.42	128K tokens
DeepSeek Reasoner	\$0.028	\$0.28	\$0.42	128K tokens

Source: DeepSeek API documentation (Sept 2025) (^[7] api-docs.deepseek.com). Note: caching heavily affects input pricing.

The aggressive pricing is a deliberate strategy. As quoted by Reuters, DeepSeek aims "to match or exceed rivals' performance at reduced costs" to "solidify its position" (^[25] www.reuters.com). Industry commentary underscores this "price war": Chinese models (DeepSeek, Baidu's Ernie) have pushed token costs near zero, challenging Western providers (^[26] www.techradar.com). In fact, TechRadar notes "China is commoditizing AI faster than the West can monetize it", since free open models threaten the business case for paid ones (^[26] www.techradar.com).

Comparative Analysis

To compare across providers, consider example per-token costs. Table below contrasts representative **high-capacity** and **low-cost** models:

Provider	Model	Input (\$/M)	Output (\$/M)	Context (tokens)	Remarks
OpenAI	GPT-5	\$1.25	\$10.00	128K	Top-tier general model
OpenAI	GPT-4o (Vision)	\$5.00	\$20.00	128K	High visual/multi-modal
OpenAI	GPT-5 mini	\$0.25	\$2.00	32K	Lite GPT-5
Google	Gemini 2.5 Pro	\$1.25-\$2.50*	\$10-\$15*	2M	Tiered pricing
Google	Gemini 2.5 Flash	\$0.15	\$3.50 (thinking)	2M	Cheap, versatile
Anthropic	Claude Opus 4.1	\$15.00	\$75.00	200K	Highest-quality Claude
Anthropic	Claude Sonnet 4	\$3.00	\$15.00	200K	Balanced speed/perf
Anthropic	Claude Haiku 3.5	\$0.80	\$4.00	200K	Fastest, for simple tasks
xAI	Grok 3 (std mode)	\$3.00	\$15.00	128K	Scientific reasoning
xAI	Grok 3 (fast mode)	\$5.00	\$25.00	128K	Premium speed mode
xAI	Grok 3 Mini (std)	\$0.30	\$0.50	128K	Light variant
xAI	Grok 3 Mini (fast)	\$0.60	\$4.00	128K	Fast mini
DeepSeek	V3.2-Exp (chat)	\$0.28\$**	\$0.42	128K	Ultra-low cost
DeepSeek	V3.2-Exp (reasoner)	\$0.28\$\$	\$0.42	128K	(Same price for both)

* Tiered: first \$1.25/\$10 up to 200K tokens, then \$2.50/\$15 beyond.

\$\$ Cache-hit price (\$0.028) – see note.

This table highlights the **orders-of-magnitude** differences. DeepSeek's input/output are compared as cached-hit/miss: *Remarkably, even its un-cached price (\$0.28) is still below many competitors' lowest tiers.* For instance, DeepSeek's \$0.28 input is *~5x cheaper* than GPT-5 mini's \$0.25 (note: actually slightly higher; but output \$0.42 is significantly lower than GPT-5 mini's \$2.00). And DeepSeek lacks any "fast mode," yet is used for heavy reasoning tasks ("chat" mode is fully capable).

Price per 100K tokens (example): To illustrate real costs, consider 100K input + 100K output tokens (a large query). Using the above rates:

- **GPT-5:** $0.1 \times \text{input cost} + 0.1 \times \text{output cost} = \$0.125 + \$1.00 = \1.125 .
- **Gemini 2.5 Pro ($\leq 200K$):** $\$0.125 + \$1.00 = \$1.125$ (same as GPT-5 here).
- **Claude Sonnet 4:** $\$0.30 + \$1.50 = \$1.80$.
- **Grok 3 (std):** $\$0.30 + \$1.50 = \$1.80$ (identical to Claude Sonnet).
- **DeepSeek (miss):** $\$0.028 + \$0.042 = \$0.070$.

Thus, for this scenario, OpenAI and Google cost about \$1.1–1.8, while DeepSeek costs mere **7 cents**. Even using DeepSeek's slow cache-hit input, it's $\$0.0028 + \$0.042 = \$0.044$ – effectively 25x cheaper.

We graphically see three pricing tiers: (1) *Premium models* (GPT-4/5, Claude Opus) at \$10–75 per M output; (2) *Mid-tier* (Gemini Pro, Claude Sonnet, Grok) at \$3–\$15; (3) *Low-end* (OpenAI nano, Gemini Flash, Grok Mini, DeepSeek) at \$0.4–\$4. Real case studies (below) reflect these strata.

Case Studies and Example Scenarios

Case 1: Customer Support Chatbot. A company processes ~10 million tokens per month (input+output). Using GPT-5 would cost $\sim \$10 \times (10) = \100 (per million tokens * output ratio), whereas Gemini Flash would cost $\sim \$0.60 \times (10) = \6 . Claude Haiku might cost $\sim \$4 \times (10) = \40 (plus negligible inputs). DeepSeek would cost less than \$5 for the same volume. Thus, if budget is critical and moderate comprehension is acceptable, Gemini Flash or DeepSeek could reduce AI costs tenfold versus premium OpenAI models.

Case 2: Enterprise Document Summarization. Summarizing large contracts of 50K words (approx 60K tokens) per doc, 100 docs monthly (6M tokens). Opus 4.1 (legal-detail level) would cost $\sim \$156 = \90 just for input plus \$756 = \$450 output total \$540. A Sonnet-level model (cheap) would cost $\$36 + \$156 = \$192$. Grok (Fast) would cost $\$56 + \$256 = \$312$. DeepSeek: $\$0.286 + \$0.426 = \$0.712$. This shows that high-capability outputs can cost thousands of dollars with top-tier models, but affordable alternatives can shrink it to hundreds or even single digits. Even after factoring that cheaper models might produce simpler summaries, the cost variance is dramatic.

Real-World Example: Bing Chat Integration. Microsoft's integration of GPT-4.0 ("o-1") via Bing forced careful price-negotiation. Industry analysts note Microsoft extracts large volume (billions of tokens) at internal rates; public API users might pay far more. Microsoft's Azure OpenAI, since 2023, offered OpenAI models with tiered enterprise pricing and commitments, illustrating that heavy users (*enterprise scale*) get discounts. Similarly, Google's Gemini integration (1.5B users using AI Overviews) suggests Google subsidizes its own models internally, effectively giving "away" some API usage to lock in developers (a point TechRadar makes about search integration (^[19] www.techradar.com)). These strategic moves underscore that what an end-customer pays can differ from sticker price – but our analysis uses the published rates for baseline.

Survey Insight. According to developer surveys (OpenAI and GitHub user data), average ChatGPT Pro users consume *~50k tokens per day* (just illustrative). At OpenAI's GPT-4 Turbo rate (\$0.03/1k), that's \$1.5/day per user. If an enterprise with 100 daily active chat bots each consuming ~50k tokens, using GPT-4, the monthly cost approaches \$4,500 (before overhead). If switching to GPT-5 mini (~1/4 the price) or Gemini Flash (1/20th),

that budget drops drastically. This aligns with industry reports that cost control is a top concern for AI teams ([27] www.binadox.com).

Pricing Trends and Future Implications

The rapid evolution in LLM pricing reflects both market competition and technological advances. Key trends and implications include:

- **Price Wars and Commoditization:** Chinese models like DeepSeek have sparked what analysts call shifting from a *performance race to a price war* ([26] www.techradar.com). Open models (DeepSeek, Baidu Ernie) make high-end AI effectively free, challenging Western vendors' paywalls ([26] www.techradar.com). OpenAI's own response—a reported 80% price cut on flagship GPT-4 models ([9] www.techradar.com)—shows pressure to reduce costs. We expect continuing downward pressure: new models will aim to double performance at half the price, and fine-tuning or batch APIs will proliferate to cut costs.
- **Differentiated Offerings:** Providers mitigate margin loss by segmenting offerings. OpenAI, for example, has **premium** (GPT-5) vs **lite** (mini/nano) models; Google has Pro vs Flash; Anthropic has Opus vs Sonnet vs Haiku; xAI has standard vs fast vs Mini. Organizations will increasingly use hybrid strategies: heavy workloads on cheap models (e.g. Grok Mini, DeepSeek) and reserve premium models for niche tasks. Academic studies (e.g. by benchmarking labs) will quantify quality-per-cost, guiding these choices.
- **Two-Part Pricing and Ecosystem Effects:** Tiered pricing (volume discounts, batch APIs) and multi-modal features complicate direct comparisons. For instance, Google charges extra for “grounded” search queries, while Anthropic's long-context pricing premium and caching create effective nonlinear cost curves ([18] cloud.google.com) ([5] docs.anthropic.com). This means vendors might not strictly compete on raw token price, but on value (e.g. freshness, search, tool use). Pricing strategy becomes part of lock-in: e.g. Google bundling lower Gemini rates for Cloud users. We already see enterprises bundling LLM spend into cloud contracts (Azure, AWS Bedrock multi-LLM billing) to manage cost predictability ([28] www.binadox.com).
- **Use of Tokens vs. Alternatives:** Some providers experiment with non-token billing. For example, OpenAI offers a “request-based” pricing for certain APIs (e.g. Vision API in dollar per image rather than per token). Large-scale use cases might evolve hybrid pricing (e.g. fixed-price content creation subscription). However, the fundamental token model will likely remain primary for text generation at least through 2026. Tools and prompt optimization (e.g. GPT's “cost consciousness” like avoiding ChatGPT answer verbosity) will become best practices for controlling token spend.
- **Long-Term Outlook:** Given escalating hardware availability (specialized AI chips), we predict LLM inference costs will continue falling. By 2026, it's plausible that current leaders will introduce even cheaper models (e.g. GPT-5.1 with context up to 1M tokens at \$0.5/\$4) or that completely open alternatives match those specs at near-zero cost. The competition may split: Western companies will focus on premium, controlled-use markets (enterprise, where SLAs and compliance justify pay), while commodity use moves toward open-source in Asia/elsewhere. The broader AI market will likely fragment: bundling of vision, voice, browsing with LLM, cross-subsidizing some aspects.

In sum, **cost optimization** is now a core part of the AI developer playbook. Our analysis suggests that enterprises must carefully match model selection to use-case, balancing “best model” vs “acceptable model” based on token budgets. For example, as one recent whitepaper notes, using a cheaper model for 70% of routine tasks and reserving the most expensive model for 30% yields better ROI than all-in on the top model ([11] www.binadox.com). As Gartner analysts have forecast, by 2026 AI services cost will become a *chief competitive factor*, potentially surpassing raw performance in importance. The pricing data compiled here will help stakeholders make informed choices in that landscape.

Conclusion

The November 2025 snapshot reveals a highly competitive LLM API market with **vast cost differentials**. OpenAI's GPT models remain at the cutting edge of capability but also at premium prices; Google's Gemini strikes a middle ground with competitive pricing and integration benefits; Anthropic's Claude offers robust safety at moderate cost; xAI's Grok competes as a niche "scientific" model; and DeepSeek pushes prices to rock-bottom levels. Our side-by-side comparisons (in tables above) show that identical tasks could cost anywhere from a few cents to hundreds of dollars depending on provider and model.

Crucially, these rates are dynamic. New model releases (e.g. GPT-4.5, Ernie 4.5, etc.), volume agreements, and one-off promotions will further shift the playing field. We recommend that AI consumers continuously re-evaluate pricing, consider multi-provider strategies, and leverage specialized offerings (batch APIs, caching, long-context) to optimize costs.

In closing, as one industry report quipped, "LLM pricing changes faster than any cryptocurrency," and our comprehensive analysis aims to be a timely guide in this rapidly changing environment. All figures and comparisons above are grounded in the latest public sources and official documentation (^[6] techcrunch.com) (^[4] docs.anthropic.com) (^[7] api-docs.deepseek.com) (^[10] www.reuters.com). Continued transparency from providers and third-party tracking will be essential for navigating the ongoing AI pricing revolution.

External Sources

- [1] <https://www.reuters.com/technology/artificial-intelligence/openai-launches-new-gpt-41-models-with-improved-coding-long-context-2025-04-14/#:~:GPT,4...>
- [2] <https://openai.com/bn-BD/api/pricing/#:~:GPT...>
- [3] <https://cloud.google.com/vertex-ai/generative-ai/pricing?hl=he#:~:Gemin...>
- [4] https://docs.anthropic.com/en/docs/about-claude/pricing?%3F%3F%3F__hstc=43401018.71aa366c60c32c7e3032e45be702fadd.1753488000320.1753488000321.1753488000322.1#:~:Claud...
- [5] https://docs.anthropic.com/en/docs/about-claude/pricing?%3F%3F%3F__hstc=43401018.71aa366c60c32c7e3032e45be702fadd.1753488000320.1753488000321.1753488000322.1#:~:MTok%...
- [6] <https://techcrunch.com/2025/04/09/elon-musks-ai-company-xai-launches-an-api-for-grok-3/#:~:Grok%...>
- [7] https://api-docs.deepseek.com/quick_start/pricing/#:~:PRICL...
- [8] <https://www.reuters.com/technology/deepseek-releases-model-it-calls-intermediate-step-towards-next-generation-2025-09-29/#:~:Chine...>
- [9] <https://www.techradar.com/pro/why-baidus-ernie-matters-more-than-deepseek#:~:DeepS...>
- [10] <https://www.reuters.com/technology/artificial-intelligence/openai-unveils-cheaper-small-ai-model-gpt-4o-mini-2024-07-18/#:~:OpenA...>
- [11] <https://www.binadox.com/blog/llm-api-pricing-comparison-2025-complete-cost-analysis-guide/#:~:Unlik...>
- [12] https://api-docs.deepseek.com/quick_start/pricing/#:~:MODEL...
- [13] <https://openai.com/bn-BD/api/pricing/#:~:Price...>
- [14] <https://openai.com/bn-BD/api/pricing/#:~:Input...>
- [15] <https://openai.com/bn-BD/api/pricing/#:~:Text...>
- [16] <https://cloud.google.com/vertex-ai/generative-ai/pricing?hl=he#:~:Text%...>

- [17] <https://www.linkedin.com/pulse/xai-launches-grok-3-api-four-pricing-tiers-intensifying-%E6%9D%B0-%E9%82%93-q1qic#:~:;ever...>
 - [18] <https://cloud.google.com/vertex-ai/generative-ai/pricing?hl=he#:~:Groun...>
 - [19] <https://www.techradar.com/pro/why-baidus-ernie-matters-more-than-deepseek#:~:Just%...>
 - [20] https://docs.anthropic.com/en/docs/about-claude/pricing?%3F%3F%3F__hstc=43401018.71aa366c60c32c7e3032e45be702fadd.1753488000320.1753488000321.1753488000322.1#:~:When%...
 - [21] https://docs.anthropic.com/en/docs/about-claude/pricing?%3F%3F%3F__hstc=43401018.71aa366c60c32c7e3032e45be702fadd.1753488000320.1753488000321.1753488000322.1#:~:%E2%8...
 - [22] <https://techcrunch.com/2025/04/09/elon-musks-ai-company-xai-launches-an-api-for-grok-3/#:~:As%20...>
 - [23] <https://techcrunch.com/2025/04/09/elon-musks-ai-company-xai-launches-an-api-for-grok-3/#:~:Grok%...>
 - [24] <https://www.techradar.com/pro/why-baidus-ernie-matters-more-than-deepseek#:~:~:~:~:A%20f...>
 - [25] <https://www.reuters.com/technology/deepseek-releases-model-it-calls-intermediate-step-towards-next-generation-2025-09-29/#:~:~:~:~:Altho...>
 - [26] <https://www.techradar.com/pro/why-baidus-ernie-matters-more-than-deepseek#:~:~:~:~:China...>
 - [27] <https://www.binadox.com/blog/llm-api-pricing-comparison-2025-complete-cost-analysis-guide/#:~:~:~:~:LLM%2...>
 - [28] <https://www.binadox.com/blog/llm-api-pricing-comparison-2025-complete-cost-analysis-guide/#:~:~:~:~:Many%...>
-

IntuitionLabs - Industry Leadership & Services

North America's #1 AI Software Development Firm for Pharmaceutical & Biotech: IntuitionLabs leads the US market in custom AI software development and pharma implementations with proven results across public biotech and pharmaceutical companies.

Elite Client Portfolio: Trusted by NASDAQ-listed pharmaceutical companies including Scilex Holding Company (SCLX) and leading CROs across North America.

Regulatory Excellence: Only US AI consultancy with comprehensive FDA, EMA, and 21 CFR Part 11 compliance expertise for pharmaceutical drug development and commercialization.

Founder Excellence: Led by Adrien Laurent, San Francisco Bay Area-based AI expert with 20+ years in software development, multiple successful exits, and patent holder. Recognized as one of the top AI experts in the USA.

Custom AI Software Development: Build tailored pharmaceutical AI applications, custom CRMs, chatbots, and ERP systems with advanced analytics and regulatory compliance capabilities.

Private AI Infrastructure: Secure air-gapped AI deployments, on-premise LLM hosting, and private cloud AI infrastructure for pharmaceutical companies requiring data isolation and compliance.

Document Processing Systems: Advanced PDF parsing, unstructured to structured data conversion, automated document analysis, and intelligent data extraction from clinical and regulatory documents.

Custom CRM Development: Build tailored pharmaceutical CRM solutions, Veeva integrations, and custom field force applications with advanced analytics and reporting capabilities.

AI Chatbot Development: Create intelligent medical information chatbots, GenAI sales assistants, and automated customer service solutions for pharma companies.

Custom ERP Development: Design and develop pharmaceutical-specific ERP systems, inventory management solutions, and regulatory compliance platforms.

Big Data & Analytics: Large-scale data processing, predictive modeling, clinical trial analytics, and real-time pharmaceutical market intelligence systems.

Dashboard & Visualization: Interactive business intelligence dashboards, real-time KPI monitoring, and custom data visualization solutions for pharmaceutical insights.

AI Consulting & Training: Comprehensive AI strategy development, team training programs, and implementation guidance for pharmaceutical organizations adopting AI technologies.

Contact founder Adrien Laurent and team at <https://intuitionlabs.ai/contact> for a consultation.

DISCLAIMER

The information contained in this document is provided for educational and informational purposes only. We make no representations or warranties of any kind, express or implied, about the completeness, accuracy, reliability, suitability, or availability of the information contained herein.

Any reliance you place on such information is strictly at your own risk. In no event will IntuitionLabs.ai or its representatives be liable for any loss or damage including without limitation, indirect or consequential loss or damage, or any loss or damage whatsoever arising from the use of information presented in this document.

This document may contain content generated with the assistance of artificial intelligence technologies. AI-generated content may contain errors, omissions, or inaccuracies. Readers are advised to independently verify any critical information before acting upon it.

All product names, logos, brands, trademarks, and registered trademarks mentioned in this document are the property of their respective owners. All company, product, and service names used in this document are for identification purposes only. Use of these names, logos, trademarks, and brands does not imply endorsement by the respective trademark holders.

IntuitionLabs.ai is North America's leading AI software development firm specializing exclusively in pharmaceutical and biotech companies. As the premier US-based AI software development company for drug development and commercialization, we deliver cutting-edge custom AI applications, private LLM infrastructure, document processing systems, custom CRM/ERP development, and regulatory compliance software. Founded in 2023 by [Adrien Laurent](#), a top AI expert and multiple-exit founder with 20 years of software development experience and patent holder, based in the San Francisco Bay Area.

This document does not constitute professional or legal advice. For specific guidance related to your business needs, please consult with appropriate qualified professionals.

© 2025 IntuitionLabs.ai. All rights reserved.