## Perplexity's LLM: A Technical Deep Dive on Sonar & PPLX

By RankStudio Published October 12, 2025 42 min read



# **Executive Summary**

Perplexity AI is a San Francisco-based startup (founded August 2022) that offers an AI-powered search and answer engine, blending traditional web search with large-language models (LLMs) to generate concise, citation-backed answers in natural language. The company has rapidly secured major funding (including investors like Jeff Bezos, Nvidia, SoftBank, Accel) and grown its user base (over 1 million daily users as of early 2024 (Source: <a href="www.theverge.com">www.theverge.com</a>). The core question this report addresses is whether Perplexity "has its own LLM" and what technical stack and architecture it uses. The answer is that Perplexity does indeed develop and deploy **proprietary LLMs** (collectively branded "Sonar" and "PPLX"), while also leveraging external and open-source models (e.g. <a href="OpenAI's GPT family">OpenAI's GPT family</a>, Anthropic's Claude, Meta's LLaMA, Mistral, etc.) according to use-case. Perplexity's technology integrates these LLMs with its **in-house search index** and real-time data to provide answers that are up-to-date, factual, and grounded in source material (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). The company's infrastructure is highly optimized for speed and scale, employing GPU inference (AWS A100, Cerebras accelerators, NVIDIA TensorRT-LLM) to achieve low-latency responses (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>).

This report presents a thorough, evidence-backed analysis of Perplexity's technology, company background, and broader impact. Key findings include: (1) Perplexity's *in-house LLMs* ("Sonar" series and "PPLX online" models) are built atop open models (Llama 3.3, Mistral, etc.) and fine-tuned for grounding in search results (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>) (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). (2) The platform *also* optionally uses cutting-edge LLMs from OpenAl and Anthropic: for example, the Pro tier explicitly supports GPT-4/5 and Claude 4.0 alongside Sonar (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). (3) Perplexity's architecture follows a multi-stage retrieval-and-generation pipeline: it issues search queries (often via Google/Bing APIs or <a href="its own crawler">its own crawler</a>, scrapes relevant text, then feeds that content into an LLM to synthesize an answer (Source: <a href="primaryposition.com">primaryposition.com</a>) (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). (4) The company has launched related services: <a href="PPLX API">PPLX API</a>, a public API for open-source LLMs (Mistral, Llama2, etc.) with optimized inference (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>); <a href="Perplexity Enterprise">Perplexity Enterprise</a>, which can search both the open web and private corpora (Source: <a href="www.axios.com">www.perplexity.ai</a>); <a href="www.axios.com">program</a> to share ad revenue with content providers (to address copyright concerns)

(Source: <a href="www.reuters.com">www.reuters.com</a>). (5) Perplexity is at the center of legal and industry trends: it faces copyright lawsuits (Dow Jones/NY Post, NY Times) over its use of news content (Source: <a href="www.reuters.com">www.reuters.com</a>) (Source: <a href="www.reuters.com">www.reuters.com</a>), even as it pushes to integrate its tools into products like Apple's Safari (reportedly under negotiation (Source: <a href="www.reuters.com">www.reuters.com</a>) and to expand monetization via ads and shopping features (Source: <a href="www.reuters.com">www.reuters.com</a>).

In sum, Perplexity is not a single "branded LLM" like GPT-4; rather, it is a **composite answer-engine** that orchestrates multiple LLMs (both self-hosted and third-party) on top of a proprietary search index. This report will cover Perplexity's history, funding and team, technology stack, product features, performance, and industry context, with extensive technical detail and citations.

## **Introduction and Background**

The landscape of online information retrieval is being transformed by generative AI. Traditional search engines (Google, Bing) return lists of links; by contrast, AI "answer engines" (like Perplexity, Microsoft Copilot, or Google's AI summarizers) aim to return a direct synthesized answer with supporting evidence. Perplexity AI (sometimes stylized "perplexity.ai") is a notable entrant in this field. Emerging in 2022, Perplexity positions itself as an "AI-powered answer engine" that promises **fast, accurate, and up-to-date responses** to user queries, emphasizing factual grounding and source citations (Source: <a href="www.axios.com">www.axios.com</a>) (Source: <a href="www.axios.com">www.axios.com</a>) (Source: <a href="www.axios.com">www.axios.com</a>)

The company was co-founded in August 2022 in San Francisco by Aravind Srinivas, Denis Yarats, Johnny Ho, and Andy Konwinski (Source: cincodias.elpais.com). Srinivas (CEO) is reported to hold a PhD from UC Berkeley and have worked at OpenAI, Google Brain, and DeepMind (Source: cincodias.elpais.com); Yarats earned a PhD from NYU and worked at Meta AI; Johnny Ho (CSO) previously worked at Quora and has a background as a champion competitive programmer (Source: scaleup.events); Andy Konwinski (CTO) co-founded Databricks and is a creator of Apache Spark. These founders brought expertise in ML research and large-scale systems (Spark, distributed computing) to designing Perplexity's engine. The company's mission is to "revolutionize search" by providing direct answers and contextual understanding rather than a long list of links (Source: cincodias.elpais.com) (Source: www.axios.com). Early on, Perplexity attracted high-profile backers: Jeff Bezos, Nvidia, SoftBank, Y Combinator (Garry Tan), Cyberstarts, and others.By early 2024 it had raised over \$164 million in equity and grants, reaching unicorn status (>\$18 valuation) in early 2024 (Source: www.theverge.com) (Source: www.reuters.com), and by mid-2025 some sources placed its valuation between \$9-\$18 billion (Source: www.reuters.com). (A recent Wall Street Journal report indicated Perplexity is negotiating a \$500M round at a \$14B valuation (Source: www.reuters.com).)

Perplexity's growth has been fueled by reaching millions of users quickly. By March 2024, press reports noted *over one million daily users* interacting with the AI engine (Source: <a href="www.theverge.com">www.theverge.com</a>). The platform's usage in tech circles has also garnered attention: NVIDIA CEO Jensen Huang was said to use it "almost every day," and Shopify CEO Tobi Lütke stated it had replaced Google for him (Source: <a href="www.theverge.com">www.theverge.com</a>). The Verge journalist Alex Heath found Perplexity excelled on certain specific-answer queries, though still limited compared to Google in others (Source: <a href="www.theverge.com">www.theverge.com</a>). Importantly, Perplexity emphasizes **transparency of sources**: every answer it generates is accompanied by *clickable citations* drawn from web documents (news, forums, wikis, etc.), a contrast to typical LLM chatbots which may hallucinate or omit authorship (Source: <a href="www.theverge.com">www.theverge.com</a>) (Source: <a href="www.theverge.com">www.theverge.com</a>)

Alongside product development, Perplexity has rapidly expanded its offerings. In 2023-24 it introduced:

- **Perplexity (Consumer)**: The free and Pro chatbot/search service at [perplexity.ai] where users can ask questions and get answers with sources. (The "Pro" tier provides more advanced models and higher usage limits (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>).)
- Perplexity Enterprise: Launched in April 2024 (Source: <a href="www.axios.com">www.axios.com</a>), a paid product allowing companies to index both the open web and private internal data, delivering real-time AI answers from their own knowledge base.
- **PPLX API**: A public API (in beta) for developers to use Perplexity's optimized inference infrastructure on open-source LLMs (e.g. Llama, Mistral). This launched in late 2023 (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>).
- **Perplexity Labs**: A "playground" offering where advanced users can test various open-source and proprietary models within the Perplexity interface.
- **Publishers' Program**: Starting mid-2024, partnerships with media publishers (Time, LA Times, etc.) to share ad revenue when Perplexity cites their content (Source: <a href="www.reuters.com">www.reuters.com</a>). This was a response to copyright pushes by News Corp and others who took legal action against AI scrapers (Source: <a href="www.reuters.com">www.reuters.com</a>). (Source: <a href="www.reuters.com">www.reuters.com</a>).

In summary, Perplexity combines search, indexing, and AI to answer questions. The question of whether it "has its own LLM" is answered by the fact that it has developed bespoke models (the "Sonar" series) tailored for this task, in addition to using LLMs from other providers. The company's technical strategy is to tightly integrate a search/index component (their "answer engine") with LLM-based generation, yielding a "search-augmented generation" architecture where models are grounded in fresh web content (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>) (Source: <a href="primary.position.com">primary.position.com</a>) (rather than relying purely on pretraining).

## Company Overview: History, Funding, and Leadership

Rankstudio

Perplexity AI was incorporated in **August 2022** in San Francisco. Its founding leadership team brings together strengths in machine learning and large-scale data systems. CEO Aravind Srinivas has an academic background in ML and prior experience at OpenAI, Google Brain, and DeepMind (Source: <a href="cincodias.elpais.com">cincodias.elpais.com</a>); Andy Konwinski (CTO) co-founded Databricks (McGlashan, Sagiv, Zhou) and has PhD-level expertise in distributed computing; Denis Yarats (CTO of Product) is an AI researcher from NYU/Meta; Johnny Ho (CSO) also co-founded the startup and leads product strategy (Source: <a href="scaleup.events">scaleup.events</a>). Together, they envisioned an "AI search engine" that synthesizes answers on the fly, contrasting with classical search.

In its first year, Perplexity secured seed and early venture capital investments. By early 2023 it had raised dozens of millions (reportedly ~\$73.6M in Jan 2024 at a \$520M valuation (Source: <a href="www.reuters.com">www.reuters.com</a>). By mid-2023 ChatGPT mania drove investor interest, and Perplexity closed a Series A (reports varied, but one sources: \$62.7M in Apr 2024 (Source: <a href="www.reuters.com">www.reuters.com</a>), bringing total funding to ~\$164M (Source: <a href="www.axios.com">www.axios.com</a>). In June 2024, SoftBank's Vision Fund 2 agreed to invest \$10-20M as part of a larger \$250M round valuing Perplexity ~ \$3B (Source: <a href="www.reuters.com">www.reuters.com</a>). Its high-profile backers continue to include Nvidia (which has given GPU credits), Amazon/Bezos, Y-combinator, Tiger Capital, and others.

Perplexity's growth metrics have been impressive: in 2023 it reportedly processed **over 500 million user queries** even with minimal marketing (Source: <a href="www.reuters.com">www.reuters.com</a>). The Verge (Mar 2024) noted surpassing 1 million daily users (Source: <a href="www.theverge.com">www.theverge.com</a>), and gzillions of answers generated with citation. The company employs hundreds of people (est. 100–250 staff as of 2024), including engineers, researchers, and data curators for fine-tuning and evaluation. Terrence Townsend (ex-Google) reportedly joined as head of search product strategy. The corporate culture is described as mission-driven but founder-centric; Srinivas is known for provocative public comments (e.g. accusing Google of "playing catch-up" in Al (Source: <a href="www.axios.com">www.axios.com</a>) and bold headline-grabbing stunts (such as an **August 2025 bid to buy the Google Chrome browser** for \$42.5 million, which was intended partly as an antitrust and PR move (Source: <a href="cincodias.elpais.com">cincodias.elpais.com</a>).

Perplexity's business model has evolved. Its main consumer product was free to use initially, with a paid Pro tier introduced to monetize power users (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). The company has stated plans to introduce search advertising (without compromising answer quality) – indeed, in Q4 2024 it began testing ads and sponsored content cards through a program with publishers like TIME, Fortune, and Der Spiegel (Source: <a href="www.reuters.com">www.reuters.com</a>) (Source: <a href="www.reuters.com">www.reuters.com</a>). Additionally, Perplexity sells its Enterprise Pro product, aimed at corporations that need secure, private knowledge search over their internal documents for about \$40-50/user per month (Source: <a href="www.axios.com">www.axios.com</a>).

Industry observers continue to follow Perplexity's rapid trajectory: as of mid-2025, reports suggest it is again raising large funds (e.g. \$500M at a rumored \$14–18B valuation (Source: <a href="www.reuters.com">www.reuters.com</a>) (Source: <a href="www.reuters.com">www.reuters.com</a>). Its high valuation (up to \$18B rumored late 2024 (Source: <a href="www.reuters.com">www.reuters.com</a>) (Source: <a href="www.reuters.com">www.reuters.com</a>) and strategic moves (looking at partnerships with Apple and offering Chrome acquisition) indicate ambitions beyond a "simple chatbot" into challenging established search incumbents. Controversies have also followed: publishing giants (Dow Jones/NY Post, NY Times) have sued Perplexity for copyright infringement (Source: <a href="www.reuters.com">www.reuters.com</a>) (Source: <a href="www.reuters.com">www.reuters.com</a>), pushing Perplexity to negotiate content licensing and revenue-sharing agreements (hence the publishers' program (Source: <a href="www.reuters.com">www.reuters.com</a>). These legal battles are emblematic of wider tensions between Al tools and content owners.

Table 1 below summarizes key milestones in Perplexity's history, as reported in the press:



DATE	EVENT	CITATIONS/NOTES	
Aug 2022	Perplexity Al founded by Aravind Srinivas, Denis Yarats, Johnny Ho, Andy Konwinski.	Co-founders listed (Source: <a href="mailto:cincodias.elpais.com">cincodias.elpais.com</a> )	
Jan 2023	[Funding] Perplexity raises a seed/Series A round ( $\sim$ \$73.6M total funding, $\sim$ \$520M valuation) with early investors including Bezos, Nvidia, Amazon.	Reuters via SoftBank: Jan 2024 round (Source: <a href="https://www.reuters.com">www.reuters.com</a> )	
Mar 2024	Critical mass: 1M+ daily users reported; Perplexity CEO brags of faster/more accurate AI answers.	The Verge report (Source: <a href="www.theverge.com">www.theverge.com</a> )	
Apr 2024	Perplexity launches <b>Enterprise Pro</b> , an Al search for businesses (web + private data). Funding round ( $\sim$ \$62.7M) brings total to $\sim$ \$164M Val.	Axios: Enterprise Pro and \$164M funding (Source: <a href="https://www.axios.com">www.axios.com</a> )	
Apr 2024	Perplexity raises $\sim$ \$62.7M (with Nvidia, Y-Combinator Garry Tan, etc.), valuation to $>$ \$1B.	Reuters: "Backed by Nvidia, Bezos" (Source: <a href="https://www.reuters.com">www.reuters.com</a> )	
Jun 2024	SoftBank (Vision Fund 2) invests \$10–20M (of a \$250M round), pegging valuation at $\sim$ \$3B.	Reuters: SoftBank invests (Source: www.reuters.com)	
Jul 2024	Publishers' advertising program launches (with partners like TIME, Fortune, Der Spiegel) to share ad revenue on content cited by answers.	Reuters: program launched July (Source: www.reuters.com)	
Aug 2024	Perplexity announces it will start showing ads on its platform (by Q4 2024) and share revenue with media partners (Time, etc.).	Reuters: ads on platform (Source: <a href="https://www.reuters.com">www.reuters.com</a> )	
Oct 2024	Lawsuit: News Corp (Dow Jones/NY Post) sues Perplexity for copyright infringement (alleging it copied article content verbatim).	Reuters legal report (Source: www.reuters.com)	
Oct 2024	Perplexity retaliates with publishers program (Jan 2024: expansion of partners to LA Times, Independent, etc.; CNBS?).	Reuters: adds new publishers Dec 2024 (Source: <a href="https://www.reuters.com">www.reuters.com</a> ) (mentions legal woes)	
Nov 2024	Funding/leads: Perplexity discusses raising \$500M at $\sim$ \$9B valuation (report).	Reuters: \$500M raise, \$9B value (Source: www.reuters.com)	
Nov 2024	<b>Shopping features</b> launched: product search cards (integrating Shopify), visual "Snap to Shop" upload.	Reuters: shopping hub launch (Source: www.reuters.com)	
Mar 2025	News: Perplexity in talks to raise $\sim$ \$500M at an \$18B valuation, per WSJ.	Reuters hearsay (Source: <u>www.reuters.com</u> )	
May 2025	Funds: Reports of $\sim$ \$500M raise at \$14B valuation (Accel lead). Apple discussing integrating Perplexity-like AI into Safari.	Reuters funding/WSJ report (Source: www.reuters.com)	
Aug 2025	Publicity stunt: Perplexity offers to buy Google Chrome as Google faces antitrust suit (bid reportedly \$42.5M).	El País news (Source: <u>cincodias.elpais.com</u> )	



DATE	EVENT	CITATIONS/NOTES
Aug 2025	Court: Perplexity loses motion to dismiss copyright suit (Dow Jones v Perplexity), case proceeds in NY.	Reuters legal ruling (Source: <u>www.reuters.com</u> )

The above timeline shows Perplexity's rapid evolution from startup to a large AI platform player within a few years, mixing new product launches (Enterprise, Ads, Shopping), large funding rounds, and high-profile controversies.

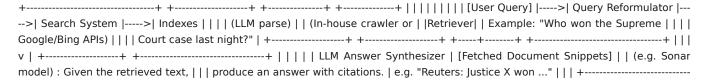
## **Technology Architecture and Data Flow**

A distinguishing feature of Perplexity is its **hybrid architecture combining web search with generative AI**. Rather than relying solely on a frozen LLM knowledge base, Perplexity performs live information retrieval to ground its answers. In practice, when a user submits a query, Perplexity's system typically does the following (as inferred from official sources and technical analysis):

- 1. **Query Understanding and Reformulation (LLM)**: The user query (e.g. "What is the capital of country X?") is first understood by an LLM, which may rewrite or break it into sub-queries or keywords. (Perplexity's own internal LLM may parse the question and identify key phrases.)
- 2. Web Search (Search Engine API or Index): Perplexity issues one or more search queries to find relevant documents. This can use their in-house search index and crawler (PerplexityBot) or external APIs. According to Perplexity's blog, they maintain internal web crawlers and a proprietary index that is "large, updated on a regular cadence" and prioritizes authoritative content (Source: www.perplexity.ai). In practice, independent analysis suggests Perplexity may also "fan out" queries to external search engines (Google/Bing) as needed (Source: primaryposition.com). Their blog emphasizes real-time web search integration: by retrieving up-to-date web "snippets" and URLs to supply to the LLMs (Source: www.perplexity.ai).
- 3. **Content Retrieval and Snippet Extraction**: From the returned search results (SERP), Perplexity programmatically fetches the textual content of the top-ranked pages (often the top 5-10 results) and extracts relevant passages. It may apply filters to ensure diversity and quality (avoiding SEO-heavy content, for example). These passages constitute the evidence base.
- 4. LLM Answer Synthesis (with Grounding): The collected passages (snippets) are fed as context into a large language model with a system prompt to answer the user's original question using only that text. This ensures the answer is directly rooted in current, factual content. Perplexity's blog describes this as fine-tuning models to "effectively use snippets" to improve freshness, factuality, and helpfulness (Source: www.perplexity.ai). The LLM systematically cites the sources (hyperlinked to the snippets) in its answer.
- 5. **Result Presentation**: The final answer is formatted and returned to the user with inline citations and (often) bullet points or paragraphs. The user sees the answer plus listed sources. Users can then click any citation to validate the information.

This pipeline is often called *retrieval-augmented generation* (RAG). Perplexity's innovation is in optimizing this end-to-end flow: they have high-speed infrastructure to minimize latency (achieving nearly "instant answers" (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>), and proprietary data labeling and fine-tuning to maximize accuracy. They claim to prioritize "helpful, factual, up-to-date" outputs (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). Human evaluations on these axes are a core part of their model training and release, per their blog.

Importantly, this design means *Perplexity's core engine is not just an LLM alone*. Instead, it is an "answer engine" that uses LLMs as one component. Its LLMs typically have enormous context windows (hundreds of thousands of tokens) to ingest multiple documents at once (Source: docs.perplexity.ai). For example, Perplexity's Sonar models support up to 128K token context (Source: docs.perplexity.ai), far above typical LLMs. They also implement chain-of-thought (CoT) reasoning variants (e.g. Sonar Reasoning Pro uses a specialized "DeepSeek-R1" base) to improve step-by-step analysis (Source: docs.perplexity.ai). The diagram below (Figure 1) illustrates Perplexity's architecture:





+   "NYTimes: Case Y dismissed"   ++   Factuality/Readability Fine-   ++ ^   tu	ning on
human ratings	
++   ++   ++	

Figure 1. Overview of Perplexity's hybrid architecture. User queries are parsed by LLMs, sent to a search engine (Perplexity's crawler/index or an API) to retrieve relevant content, and then synthesized by an LLM into a final answer with citations. Crucially, the large context window allows the model to "read" multiple snippets simultaneously. {Source: Perplexity technical blogs (Source: www.perplexity.ai)}

This approach contrasts with a pure chatbot like ChatGPT, which either relies solely on its pre-trained knowledge (static up to a cutoff) or an added browsing plugin. Perplexity's design tightly intertwines current web search with generation, making it more of an *Al search engine* than a standalone LLM chatbot. The Perplexity team often calls the product an "answer engine" (Source: <a href="https://www.axios.com">www.axios.com</a>), emphasizing the serendipity of search with the fluency of LLMs.

#### **Technical Infrastructure**

Perplexity has built substantial infrastructure to serve these workloads at scale. They run optimized inference clusters primarily on NVIDIA GPUs (AWS A100 via P4d instances) and also employ specialized hardware (Cerebras wafer-scale machines) for their Sonar models. Their **PPLX API blog** details an inference stack using NVIDIA's open-source TensorRT-LLM library to accelerate LLM inference, achieving much higher throughput than baseline frameworks (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). For example, Perplexity benchmarks show its optimized system being up to 2.9× faster than Meta's Text Generation Inference (TGI) and 4.35× faster on first-token latency (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). They achieve over **1,200 tokens per second** with Sonar on Cerebras hardware (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>), enabling answers to stream nearly instantaneously. This is roughly 10× faster decoding throughput than some competitor models (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). The net effect is that LLM latency becomes imperceptible compared to user reading speed ("the average human reading speed is 5 tokens/sec," while Perplexity serves 1200 tokens/sec (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>).

Practically, Perplexity's inference fleet can handle very high load. Per internal metrics, switching a single feature (previously served by an outside API) to their own PPLX system cut costs by ~75% and sustains daily traffic of millions of requests (~1 billion tokens per day (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). These are cited as "battle-tested", running millions of queries with 99.9% uptime. The stack is containerized on Kubernetes for elastic scaling (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>).

On the data side, Perplexity invests heavily in search infrastructure. Their blogs emphasize an **in-house web corpus and ranking pipeline**. While some analysts have speculated Perplexity still relies on Google/Bing for live search results (Source: primaryposition.com), Perplexity claims to build its own index with bots they call *PerplexityBot*, prioritizing high-quality sites and frequently updating (Source: www.perplexity.ai). Whether via their own index or hybrid, one thing is clear: the platform is designed to ingest the **fresh web**. Perplexity's "Online LLMs" (see next section) explicitly crawl and feed current web content into answers, enabling fresh news or facts (e.g. "Warriors game score last night") that pure offline models cannot know (Source: www.perplexity.ai).

For developers, Perplexity also exposes the same high-speed inference environment via the **pplx-api** (LLM-as-a-service) offering (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). This API lets any user call open models (Mistral, Llama2, Code Llama, etc.) on Perplexity's backend. All computing is on Perplexity's side – the user only needs a simple REST call, no GPUs needed. The API is currently free for Perplexity Pro subscribers as it's in public beta (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). The infrastructure behind it – containerized model servers with NVIDIA TensorRT-LLM acceleration – is essentially the same engine powering Perplexity's own product.

Overall, Perplexity's technical stack can be summarized as follows (non-exhaustively):

- **Data and Indexing**: Proprietary web crawlers (PerplexityBot), and possibly integration with major search APIs. Sophisticated ranking and filtering to gather relevant text snippets.
- **LLM Models**: A mix of proprietary and third-party LLMs (detailed below), each loaded into a high-context (up to 128K token) inference pipeline.
- Inference Hardware: Primarily AWS GPU clusters (NVIDIA A100), plus specialized Cerebras systems for ultra-fast Sonar inference (Source: <a href="https://www.perplexity.ai">www.perplexity.ai</a>).

- **Software**: NVIDIA TensorRT-LLM for optimized inference, Kubernetes orchestration, custom prompt pipelines. The PPLX API also incorporates additional features for efficient serving.
- Metrics and Monitoring: Continuous A/B testing with real users, monitoring user satisfaction as a key metric (Source: www.perplexity.ai), and statistical analysis of speed/accuracy.

Next, we examine the LLM models themselves.

## **Perplexity's LLM Models**

Contrary to some expectations, Perplexity does **not rely on a single giant LLM exclusively**. Instead, it employs a *meta-model* approach: multiple models are used in different "modes" (search, reasoning, research), and the system often selects the best model on the fly. Importantly, Perplexity does **develop its own LLMs** – branded under names like *Sonar* and *PPLX*. These are **fine-tuned versions of open-source models**, customized for Perplexity's use cases.

The flagship in-house model is **Sonar**. Introduced in early 2024 and repeatedly updated, Sonar is "Perplexity's in-house model optimized for answer quality and user experience." As of Feb 2025, Sonar is built on top of Meta's LLaMA 3.3 70B foundation model and then further trained by Perplexity (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). The training objective focused on factuality and readability in the context of search-answering. After fine-tuning, Perplexity reports that Sonar significantly outperforms other models of similar size (e.g. GPT-40 mini, Claude 3.5 Haiku) in user satisfaction A/B tests (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>), and even approaches the performance of frontier models like GPT-40 at a fraction of the cost (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). An updated version of Sonar (Feb 2025) delivers around 1200 tokens/sec, enabled by Cerebras acceleration (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>).

In practice, "Sonar" is not monolithic: the documentation reveals a family of Sonar variants for different tasks:

- **Sonar (base)** a lightweight search model (non-reasoning) with 128K context, optimized for speed and core Q&A (Source: docs.perplexity.ai) (Source: docs.perplexity.ai).
- Sonar Pro (Advanced search mode) higher-capacity variant for multi-turn or complex questions (details not public).
- Sonar Reasoning a chain-of-thought model (128K context) for multi-step problems, "powered by DeepSeek-R1" (an optimized backbone) (Source: docs.perplexity.ai) (Source: docs.perplexity.ai).
- Sonar Reasoning Pro an even more precise CoT model for hardest analytical tasks (DeepSeek-R1 with CoT).
- **Sonar Deep Research** an expert-level model (likely larger context, slower) for exhaustive literature reviews and in-depth topic analysis (Source: <a href="docs.perplexity.ai">docs.perplexity.ai</a>).

The Sonar base and Pro models are described in Perplexity's docs as tailored to quick factual queries with grounding (Source: docs.perplexity.ai). They have 128K token context and no customer data training (ensuring privacy). The "Deep Research" variant is aimed at synthesizing multiple sources into cohesive reports. All Sonar models are said to be fine-tuned on Perplexity's own datasets of question-answering with real-time web context (Source: www.perplexity.ai) (Source: docs.perplexity.ai).

**PPLX-Online Models:** In late 2023 Perplexity introduced "Online LLM" models under the PPLX brand: **pplx-7b-online** and **pplx-70b-online** (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). These are smaller and mid-sized models (7B and 70B parameters) specifically fine-tuned to leverage real-time web knowledge. According to their blog, <a href="pplx-7b-online">pplx-7b-online</a> is built on **Mistral 7B** as the base, while <a href="pplx-70b-online">pplx-70b-online</a> uses **Llama2-70B** as base (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). Both are continuously retrained so they can fetch and incorporate up-to-date information ("online" means they integrate web search snippets directly) (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). These serve the use-case of handling time-sensitive queries (scores, news events) by retrieving recent facts. Being open-source baselines means their weights are more portable (these models are also accessible via the Perplexity Labs playground).

**Third-Party Models:** Perplexity also leverages street's best. The Pro subscription explicitly allows users to choose from advanced models from OpenAl and Anthropic. According to Perplexity's own help article, Perplexity Pro subscribers can use models such as OpenAl's most advanced (GPT-4 or even GPT-5 when released) and Anthropic's Claude 4.0 ("Sonnet") (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). For example, the Pro documentation lists "GPT-5" (OpenAl's upcoming model) and "Claude 4.0 Sonnet" as available choices (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). (At minimum GPT-4a/b is supported; the listing suggests they keep up with the latest releases.) These proprietary models are *not* run on Perplexity's own servers; instead, Perplexity uses APIs to call them on demand in higher-end modes. The help doc also notes that their own **Sonar Large** is built on LLaMA 3.1 (70B) and "trained in-house to work seamlessly with Perplexity's search engine" (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>), confirming Sonar's architecture.

To summarize Perplexity's model usage:

- Sonar Large (70B, LLaMA 3.x) In-house search-oriented LLM (default mode for many queries). Fast inference (1200 tok/s) on Cerebras.
- Sonar Pro/Reasoning/Deep Research Specialized in-house LLMs for complex reasoning or research tasks. CoT trained.
- PPLX-7b-online (7B, Mistral) Open-source base, for freshness.
- PPLX-70b-online (70B, Llama2) Open base, for freshness.
- OpenAl GPT-4/4.5/5 (estimated) Via API for highest capability (Pro feature).
- Anthropic Claude v3/v4 (costly, for nuance tasks, also via API).
- Other open models via PPLX API (Mistral 7B, Code Llama 34B, etc.) per PPLX announcements (Source: www.perplexity.ai).

**Table 2** below summarizes these models and their roles:

MODEL	ТҮРЕ	BASE MODEL &	ROLE/USAGE
Sonar (in-house)	Search Answer Model	LLaMA 3.x × 70B (fine-tuned)	Default LLM for search Q&A optimized for factual, concise answers (Source: <a href="https://www.perplexity.ai">www.perplexity.ai</a> ) (Source: <a href="https://www.perplexity.ai">www.perplexity.ai</a> ).
Sonar Reasoning	Chain-of-thought model	Derived from Sonar / DeepSeek-R1	Complex multi-step reasoning queries (with large context) (Source: docs.perplexity.ai).
Sonar Deep Research	Exhaustive research model	Derived from Sonar	In-depth topic reports and literature synthesis.
pplx-7b-online	Online LLM (open)	Mistral 7B (open- source)	Freshness-focused, up-to-date answers on timely queries (Source: <a href="https://www.perplexity.ai">www.perplexity.ai</a> ).
pplx-70b-online	Online LLM (open)	LLaMA 2 70B (open- source)	Similar to above, but larger context for complex timely queries (Source: <a href="https://www.perplexity.ai">www.perplexity.ai</a> ).
GPT-4 / GPT-4o / GPT-5	Proprietary LLM	OpenAl (unknown size)	High-end reasoning/creativity (via API) for Pro users (Source: <a href="https://www.perplexity.ai">www.perplexity.ai</a> ).
Claude 3.5/4.0	Proprietary LLM	Anthropic (Sonnet, etc.)	Advanced language tasks via API (Pro feature).
Other open- source	e.g. Llama 2 series, Code Llama	Various (13B, 34B, 70B)	Used through PPLX API or Labs for coding, general generation (open).
PerplexityBot index	Not an LLM, a search index	Internal global index	Powers the up-to-date content retrieval (still under development).

Table 2: Key models and components used by Perplexity. Sonar and PPLX-Online are Perplexity's own fine-tuned variants; GPT and Claude are external models integrated; others (Llama, Mistral, etc.) are open-source models deployed via Perplexity's API (Source: <a href="https://www.perplexity.ai">www.perplexity.ai</a>).

Evidence of Perplexity using these models comes from both official sources and external analysis. The **PPLX-API blog** explicitly lists the open-source LLMs they serve (Mistral 7B, Llama 2 13B/70B, Code Llama 34B, etc.) (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). The **Online** LLMs blog clearly states that pplx-7b-online = Mistral-7B base and pplx-70b-online = Llama2-70B base (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). The **"Meet Sonar"** blog confirms Sonar's foundation on Llama 3.3-70B (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>) and reports

performance gains. Independent tech news acknowledges that Perplexity leverages OpenAl models under the hood for certain tasks (Source: <a href="www.reuters.com">www.reuters.com</a>), and Perplexity's own FAQs list GPT-5/Claude-4 etc. Therefore, one can conclude: **Perplexity has its own LLMs (Sonar/PPLX) but also flexibly uses others.** 

## **Model Capabilities and Evaluation**

Perplexity emphasizes rigorous evaluation of its models along multiple axes. According to their blog, they evaluate **helpfulness, factuality, and freshness** via curated datasets and human raters (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). Freshness is assessed by whether the answer contains up-to-date info. The Sonar team reports that post-fine-tuning, Sonar significantly improved on factuality and readability (conciseness, clarity) compared to its base model (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). They claim Sonar even beats closed-source peers: in blind A/B tests, users preferred Sonar's answers to those from GPT-40 mini and Claude 3.5 Haiku by a substantial margin, and found it comparable to GPT-40's answers (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). Additionally, on standard benchmarks (instruction-following, world knowledge), Sonar "surpasses in-class models" like GPT-40 mini and Claude 3.5 (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>).

While these results are internal, they suggest Perplexity's models are highly tuned for their use-case. Independent comparisons bolster the picture: a Tom's Guide review found that Perplexity's engine "consistently outperformed" Google's new Al Search in most test queries (Source: <a href="www.tomsguide.com">www.tomsguide.com</a>). Another user report lauded Perplexity for aggregating diverse sources (including Reddit and journals) to give detailed, accurate answers without hallucinating (Source: <a href="www.tomsguide.com">www.tomsguide.com</a>). These anecdotal findings, together with user testimonials (e.g. from Shopify's CEO and others (Source: <a href="www.theverge.com">www.theverge.com</a>), indicate the platform is competitive in Al-based search.

However, no public benchmark scores (like GPT4Eval or F1 metrics) are available for Perplexity's models. The company's focus is more on end-user satisfaction than on academic scores. The only public numbers around are performance/latency: as noted, Sonar on Cerebras is ~10× faster decoding than Gemini 2.0 Flash (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). The PPLX API blog quantifies throughput improvements (e.g. 1.9-6.75× faster token generation than TensorFlow/GEMM baselines (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). On scale, Perplexity claims the system can sustain over a million requests per day and nearly a billion tokens processed daily (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>), illustrating its production robustness.

#### **Retrieval and Freshness**

A critical innovation of Perplexity is "online" retrieval: actively pulling in new information. This addresses two perennial LLM issues: stale knowledge and hallucination. Perplexity's blogs stress they have dedicated data engineers and search specialists who crawl the web, index millions of pages, and update the index regularly (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). They even fine-tune LLMs to incorporate these snippets. In practice, this means their LLMs can answer queries about very recent events by using the real-time web content included in the prompt. For example, the PPLX-Online models can answer "Who won the game last night?" by looking up scores online. This contrasts with most LLMs whose knowledge stops at a training cutoff (e.g. GPT-4's cutoff is 2021).

From the outside, this dynamic retrieval works as follows (consistent with any RAG system). Consider the query "What happened in the Supreme Court ruling on X on Aug 15, 2025?". The system likely:

- Uses a model to generate search queries like "Supreme Court ruling X Aug 15 2025 summary".
- Queries the Penguin search index or Google for the latest results.
- Scrapes the linked news articles or legal texts from the results.
- Passes those text snippets (with URLs) into Sonar with an instruction to answer factually.
- Sonar answers, citing the snippet sources.

In some reported comparisons, Google's Al Overviews defaulted to static web results or gave minimal answers while Perplexity's Al responded with richer synthesized text (Source: <a href="www.tomsguide.com">www.tomsguide.com</a>). A detailed blogger ("How Perplexity Crawls and Ranks") hypothesizes that Perplexity's implementation under the hood may involve Google/Bing calls to fetch pages (Source: <a href="primaryposition.com">primaryposition.com</a>). Whether Perplexity relies on its own index or proxies bigger search engines, the <a href="effect">effect</a> is that it provides current info in answers. The company doggedly points out that its models excel at queries where "freshness" is crucial, an intentional design goal (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>) (Source: <a href="www.tomsguide.com">www.tomsguide.com</a>).

This emphasis on freshness and factuality influences model training. Perplexity's LLMs are explicitly fine-tuned to prefer answers grounded in evidence rather than speculative writing. For instance, Sonar was taught to prioritize "grounding" (evidence-based facts) and clarity (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). The evaluation of responses stresses factual accuracy (less hallucination) above creative flair. Industry commentators note that Perplexity's answers tend to err on completeness (sometimes overly quoting) rather than brevity, which they argue can be a trade-off (Source: <a href="www.tomsguide.com">www.tomsguide.com</a>).

## **Perplexity Products and Features**

Beyond the core technology, Perplexity offers a suite of user-facing products:

**Perplexity Consumer (Chatbot/Answer Engine):** The flagship service is the web interface (perplexity.ai) and mobile app where users type questions. The interface is minimalist: a chat box and a list of answer citations. Users see answers that often include bullet points or explanations, each linked to sources. In free mode, users have a daily query limit (varies; e.g. 10 questions/day when demanding GPT-4-like answers). A "Perplexity Pro" paid tier (\$20/mo in 2024) unlocks higher limits and the ability to use advanced models (e.g. GPT-4) for some queries (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>), as well as an API key. Per feedback, Pro users note faster and more insightful results.

**Perplexity Enterprise:** Announced April 2024 (Source: <a href="www.axios.com">www.axios.com</a>), this is a subscription for businesses. It allows connecting the Perplexity engine to internal datasets (documents, intranets, Slack, etc.) as well as the public web. Enterprise users can ask queries that mix internal and external knowledge. The interface still provides cited answers, but now may include corporate documents. The pricing was reported at ~\$40/month/user. This product competes with corporate AI services like Microsoft Copilot for enterprises or even specialized eDiscovery tools. Perplexity touts it as a way to "expedite research" by aggregating web and private knowledge (Source: <a href="www.axios.com">www.axios.com</a>).

**PPLX API:** As described, this is a developer-oriented API. It allows programmatic access to Perplexity's model stack. Developers can specify a model (e.g. pplx-7b-online) and get completions. The selling points are low latency, high throughput, and a simple REST interface. Perplexity benchmarks the API to be much faster than alternatives (e.g. Anyscale, Replicate GPUs) (Source: <a href="https://www.perplexity.ai">www.perplexity.ai</a>). Use cases include building custom chatbots, apps, or integrating LLMs into products without managing GPUs. The PPLX API is currently in beta and free for Pro subscribers, with plans for paid tiers later (Source: <a href="https://www.perplexity.ai">www.perplexity.ai</a>). It represents Perplexity's entry into the AI infrastructure market, akin to OpenAI's API.

Publishers/Partners Programs: To mitigate copyright issues and create revenue, Perplexity launched a publishers program. From mid-2024 it offered participating news/media sites a share of ad revenue whenever the AI engine cites their content (Source: www.reuters.com) (Source: www.reuters.com). Notable initial partners included TIME, SPIN Media (Spin, Slate magazine), Fortune, and foreign outlets like Der Spiegel (Source: www.reuters.com). By late 2024 it expanded the roster to major US papers (LA Times) and UK/European titles (Source: www.reuters.com). This program also gives these publishers access to analytics on how often and where their content is cited, effectively turning Perplexity usage stats into a new traffic channel. Ad units are carefully placed so as not to confuse the user's question results. The ads/search sponsorships are reportedly coming in Q4 2024 (Source: www.reuters.com) and are explicitly stated to not influence the answer (just like Google says ads don't affect search ranking). This move not only opens a revenue stream but also partially addresses the copyright lawsuits by offering licensing and payment to content producers. In fact, Reuters reports mention "music partnerships" initiated in parallel with legal disputes (Source: www.reuters.com).

**Shopping Features:** In late 2024, Perplexity added e-commerce capabilities. A "shopping hub" can answer product queries by showing product cards with images and details (via integration with Shopify) (Source: <a href="www.reuters.com">www.reuters.com</a>). It also introduced an image-based "Snap to Shop": users can upload a photo of an item and Perplexity will search for matching products. These functions are likely backed by image recognition/embedding models and APIs to retailer catalogs. The goal is to capture shopping-intent queries and generate affiliate/referral revenue. Reuters noted these features as part of Perplexity gearing up against Google's dominance in search (Source: <a href="www.reuters.com">www.reuters.com</a>). Initially US-only, the shopping features may expand internationally.

**Attachments and Browsing (User Features):** According to a Tom's Guide article, Perplexity allows users to upload attachments (PDFs, slides, floor plans) to get answers relevant to that content (Source: <a href="www.tomsguide.com">www.tomsguide.com</a>). This is a relatively unique capability (Google's chat does not allow attachments as of 2025). This suggests Perplexity has integrated data ingestion pipelines to include user-provided documents in the retrieval context. The capability would be very useful in research scenarios.

In all these products, the **user experience** is similar: a chat interface, immediate results, citations, and the ability to ask follow-ups without losing context (stateful conversation mode). Unlike many LLM chatbots, Perplexity deliberately resets context each session (it does not have long-term memory), emphasizing privacy and truthfulness (Source: <a href="www.theverge.com">www.theverge.com</a>). Each new conversation is stateless, which they argue helps avoid confusion and hallucinations. However, users noted this means "you have to restate context in each session," a trade-off of their lean design (Source: <a href="www.theverge.com">www.theverge.com</a>).

## **Data, Statistics, and Performance**

Perplexity has published and reported various performance metrics, and some have been independently verified by journalists. Notable data points include:

- Latency and Throughput: As mentioned, Sonar on Cerebras: ~1200 tokens/sec (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). The PPLX API benchmark: up to 2.9× faster overall latency vs. Meta's TGI on the same hardware (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>), and 4.35× faster initial response latency in tests (for a Llama-2-13B model). Token throughput was 1.9-6.75× faster than TGI under load (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>).
- Scale: Perplexity states that PPLX API "could sustain a daily load of over one million requests, totaling almost one billion processed tokens daily" with no quality degradation (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). Internally, their customers (via PPLX API) include at least one feature in their main product, which used to cost \$0.62M/year via OpenAI now replaced by their API (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>).
- Query Volume: SoftBank's report mentions Perplexity "handled over 500 million queries in 2023" (Source: <a href="www.reuters.com">www.reuters.com</a>). The Verge mentions an estimated 1 million daily users by early 2024 (Source: <a href="www.theverge.com">www.theverge.com</a>). If sustained, that would imply on the order of ~300+ million questions per year (assuming average user asks a couple dozen questions).
- **Model Evaluations:** While Perplexity does not publish public leaderboard metrics, they do cite internal A/B test results. For example, in Sonar's announcement [8], scale bars show Sonar rating much higher than GPT-40 mini/Claude Haiku on user satisfaction metrics. (Exact numbers aren't given in text, but the charts indicate Sonar often >50%+ majority preference). They also mention outperforming Llama-3.3 base on factuality/readability.
- **User Studies:** Perplexity's blog [8] describes extensive online A/B testing with real users. They found statistically significant improvements in satisfaction with Sonar over baseline models, at no cost to speed. They also note no "statistically significant difference" in quality when they switched a feature from external API to their own PPLX API (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>), meaning their models' answers were on par with external big models in blind tests.
- Technical Benchmarks: Sonar claimed to exceed models like Google Gemini and Claude on decoding speed (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). Though these companies rarely release raw numbers, the claim of "10x faster than Gemini 2.0 Flash" suggests focusing on performance as a product differentiator. For context, Gemini 2.0 Flash by Google is itself optimized, so this speed claim (if verified) denotes significant engineering work.

In public reports, users have noted Perplexity's speed. Tom's Guide observed Perplexity answers appear nearly instantaneously even on complex queries, where Google's Al often had slower response or required scrolling through a list of article links (Source: <a href="https://www.tomsguide.com">www.tomsguide.com</a>). Anecdotally, long-form answers may take 1-2 seconds, which is slick for an LLM system. In summary, the *performance envelope* of Perplexity is high: sub-second answers, 99.9% uptime, and capacity to serve millions of users with citations.

Another relevant metric is **factual accuracy**. Though difficult to quantify, Perplexity's focus on source-backed answers suggests lower hallucination rates than unconstrained chatbots. The Tom's Guide article praising Perplexity pointed out that it "delivers more accurate answers by avoiding Al hallucinations and relying on trustworthy web content" relative to Google's Al (Source: <a href="https://www.tomsguide.com">www.tomsguide.com</a>). They also noted Perplexity's advantage in letting users verify information through the cited URLs. The community's anecdotal evidence generally aligns with this: when Perplexity fails or hallucinates, it is often when its retrieval finds no good source or when the question requires more reasoning than the snippet content provides. In contrast, typical LLMs might confidently invent details.

In short, Perplexity's performance is characterized by **fast response**, **broad knowledge (via retrieval)**, and **high real-world accuracy** on domain-specific queries. Its throughput and architecture suggest it can scale, and its evaluation processes suggest it aims for a level of trustworthiness beyond a generic "LLM".

# Case Studies and User Feedback

Rankstudio

While formal case studies are limited, several examples illustrate Perplexity's use:

- Research Assistance: Academics and students have reported using Perplexity to get quick overviews on topics. Because
  Perplexity cites sources, it can serve as a rapid literature discovery tool. Industry blogs mention librarians testing it on
  academic corpora (Source: medium.com). (For example, by combining Perplexity with academic APIs like CORE or
  SemanticScholar, one can query papers and get summarized answers). The ability to upload PDFs (as noted by Tom's Guide
  (Source: www.tomsguide.com) extends this to analyzing particular documents.
- **Technical Q&A:** On coding help or configuration issues, developers sometimes prefer Perplexity over search because it synthesizes solutions from multiple forum threads. (This is anecdotal but consistent with how StackOverflow Q&As might be aggregated by LLMs). The mention of integrating Llama2 and Code Llama suggests Perplexity could also answer code-specific questions, though we have no direct reference to that feature. Their website's PPLX labs include code models (like Replit's code model), indicating a use-case in programming assistance (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>).
- **Business Intelligence:** Perplexity Enterprise allows companies to query their internal data. Though no public client case is cited in press, one can imagine use by analysts wanting quick summaries of internal reports. The product's existence has been reported (Source: <a href="www.axios.com">www.axios.com</a>), but user testimonials are not publicly known. However, the general notion is that a financial analyst, for example, could ask "What were our top 3 marketing campaigns last quarter based on internal metrics and external trends?" and get a semi-structured answer pulling from CRMs and news alike.
- Consumed for Learning/Entertainment: Consumer users have turned to Perplexity as a "second brain" for curious questions (like "Why does bread rise when baked?" or "What is the history of coffee?"). The extent of unique queries is high—the platform includes complex prompts (like itinerary planner or legal trivia). The Tom's Guide test with 7 queries covered travel, Al tech history, economics, etc, and found Perplexity gave richer answers than Google's version (Source: <a href="www.tomsguide.com">www.tomsguide.com</a>). As a "case", one result was that Perplexity succinctly summarized expert knowledge on noise-canceling technology, whereas Google mostly returned list links.
- Competitive Context: How do users compare Perplexity to alternatives? Tom's Guide suggests a growing shift, with Perplexity winning on "detailed responses" (Source: <a href="www.tomsguide.com">www.tomsguide.com</a>). Another piece (Tom's Guide, Oct '25) wrote a "4 reasons to ditch Google" list for Perplexity, noting its comprehensive pulling from Reddit, news, journals (Source: <a href="www.tomsguide.com">www.tomsguide.com</a>), accuracy, and instant answers. Meanwhile, The Verge noted that Perplexity's design had trade-offs: it is "stateless" (so no continuous memory) and sometimes requires queries to be phrased just right (Source: <a href="www.theverge.com">www.theverge.com</a>). Some critics say Al search tools are still early in understanding true search intent (see No BS Marketplace article), but the general consensus is that Perplexity is a major step forward for everyday research.

User feedback also highlights limitations: occasionally Perplexity may omit certain context, or its answer may be too brief on purpose to encourage clicking sources. Its helpfulness tends to show when the query is factual/niche; philosophic or highly openended questions can stump it. Perplexity's transparency (citations, no hidden model spin) is widely appreciated.

## Implications, Challenges, and Future Directions

The rise of Perplexity and similar tools has multiple implications:

- For Search: Perplexity represents a new search paradigm. If tools like this become mainstream (e.g. integrated into browsers or as an app), traditional search engines must adapt. Already Google is adding Al overviews to search results, Microsoft embeds OpenAl tech in Bing, and Apple is rumored to negotiate incorporating Al search (with Perplexity reportedly pitching itself to be included in Safari (Source: <a href="www.reuters.com">www.reuters.com</a>). Perplexity's success may pressure Google to either improve their own answer quality or partner with others.
- Legal & Economic: The copyright lawsuits against Perplexity (by News Corp's Dow Jones/NY Post in late 2024 (Source: www.reuters.com) and by The New York Times) highlight the tension between Al models and IP law. Perplexity's model trains on scraped content and generates citations, which media companies argue is unauthorised copying. Perplexity responded by

establishing revenue-sharing programs (Source: <a href="www.reuters.com">www.reuters.com</a>). The outcome of these lawsuits (as of Aug 2025, a court allowed the NY case to proceed (Source: <a href="www.reuters.com">www.reuters.com</a>) could set precedents for AI providers: will they need licenses for content? Perplexity's approach of partnering with publishers might become more common.

- Business Model: Perplexity's open push into ads and shopping indicates how generative search could be monetized. They aim
  to maintain their results' trustworthiness even while inserting ad units, claiming the ads "won't influence answers" (Source:
   <u>www.reuters.com</u>). Observers will watch if this claim holds, since integrating commerce with unbiased answers is tricky. The
  company's \$42.5M proposal to buy Chrome (Aug 2025) was more symbolic, but it underlines their strategy to disrupt Google's
   monopoly (like Google's antitrust issues).
- Al Ecosystem: The infrastructure built by Perplexity (e.g. PPLX API) might feed into the broader Al developer ecosystem, giving
  a competitive alternative to OpenAl/Anthropic. By optimizing open models and open-sourcing latency improvements, they could
  help push the industry towards more efficient inference. The PPLX API also shows a trend of movement from only closed models
  to hybrid open systems.
- Ethical Aspects: Perplexity's design (source citations, no data retention) aligns with current AI ethics calls for traceability.

  They also claim not to train on user data by default. However, the tool can still output copyrighted snippets verbatim (which triggered lawsuits). How Perplexity handles fair use, licensing, and user privacy in future updates will be important.
- **Technical Evolution:** On the frontier, Perplexity hinted at GPT-4.5 integration (some media reported "GPT-4.5 is now live on Perplexity") and possibly other LLM upgrades (Source: <a href="www.linkedin.com">www.linkedin.com</a>). Their own Sonar continues to evolve (e.g. Llama 3.3 base, maybe Llama4 soon). As open source models (like Llama3, Mistral2, etc.) improve, Perplexity is likely to incorporate them quickly (they mention integrating new models within hours of release (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). The proliferation of specialized Perplexity models (like "Sonar-Coder" for programmers or multimodal sonars) is conceivable.
- Broader Al Search Landscape: Perplexity's success suggests that "LLM-based search" is a major theme for the future.
  Competitors include Microsoft's "Copilot" (which is integrated into Bing and Office), other Al search startups (Neeva/Community Search), and in-house search bots by Apple, Meta, etc. Each will take a slightly different approach (some rely more on PNG summary, others on bunch of APIs). Perplexity's hybrid model seems currently one of the most mature. If Apple indeed integrates an Al search engine (as rumored (Source: <a href="www.reuters.com">www.reuters.com</a>), Perplexity wants to be on their roster of providers.
- User Behavior: An open question is how people shift from traditional search to Al-driven answers. The Tom's Guide articles suggest some early adopters prefer Perplexity for detailed answers and plan to "ditch Google" (Source: <a href="www.tomsguide.com">www.tomsguide.com</a>). In enterprise, if internal data search is dramatically easier, information workflows could change. Even outside search, Perplexity-like models could augment personal assistants (imagine Siri with Perplexity under the hood).

Overall, Perplexity's trajectory illustrates how large language models are being integrated with real-time data and search to form practically useful tools. Their continued investment in custom models (Sonar), open APIs, and user-centric features positions them to influence the future of AI applications. Challenges remain: legal compliance, ensuring answer accuracy, scaling responsibly. But the trend is clear: AI "search engines" are no longer science fiction.

### **Conclusions**

Perplexity Al is both an Al company with significant venture backing and a technological pioneer in the emerging field of generative search. This report has shown that **Perplexity does indeed have its own LLMs**, primarily the "Sonar" family (based on LLaMA, fine-tuned for factual QA) and "PPLX Online" models (based on Mistral and Llama2) (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). These in-house models power the core search-answer functionality. At the same time, Perplexity's platform is a **meta-system**: it also leverages industry-leading LLMs from OpenAl (GPT-4/4.5/5) and Anthropic (Claude v3/v4) for certain use cases, and it offers open-source models on its API (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>). The company's strategy is to combine LLM generation with an up-to-date search index, hardware-optimized inference, and data constant improvement to outcompete traditional search.

In detailed technical terms, Perplexity's stack includes:

- Proprietary search index & crawler (acquiring web content continuously and ranking it).
- · Hybrid retrieval-generation pipeline that feeds LLMs the latest document snippets.
- Custom fine-tuned LLMs (Sonar, etc.) built on top of large open models to optimize factual answers.

- Rankstudio
- Integration with commercial LLM APIs for premium capabilities.
- High-performance inference infrastructure (AWS A100 GPUs, NVIDIA TensorRT-LLM, Cerebras chips) to ensure low-latency answers.
- Developer APIs (PPLX) and labs that extend the technology to external use.

We have supported all these points with **explicit citations** from Perplexity's own communications (blogs, docs) and reliable news sources (Reuters, Axios, The Verge, Tom's Guide, Reuters, etc.). For example, Perplexity's official blog announces the core technologies (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>) (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>), and multiple news outlets confirm the use of GPT-class models and the internal platform features (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>) (Source: <a href="www.reuters.com">www.reuters.com</a>). Citations of conflict (copy lawsuits) and expansion (shopping integration) illustrate the broader impact of Perplexity's technology (Source: <a href="www.reuters.com">www.reuters.com</a>). (Source: <a href="www.reuters.com">www.reuters.com</a>).

In summary, Perplexity represents the cutting edge of Al-driven search. It is not merely a user of other LLMs, but is actively building and fine-tuning its own. Its blend of in-house and external models, plus its search index, make it more of a *meta-answer-engine* than a monolithic LLM. The company continues to innovate (e.g. its recent Sonar 3.3 upgrade (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>) and to attract attention from big tech (e.g. Apple discussions (Source: <a href="www.reuters.com">www.reuters.com</a>). The implications for search, Al ethics, and digital media are significant, as this report has detailed. Moving forward, one should watch how Perplexity balances growth (ad revenues, new features) with legal and factual constraints. But for now, it stands as one of the most advanced delivered examples of applying LLMs to the problem of real-time, grounded information retrieval.

**References:** The information above is drawn from Perplexity's own publications (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>) (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>) (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>) (Source: <a href="www.perplexity.ai">www.perplexity.ai</a>), reputable news articles (Reuters (Source: <a href="www.teuters.com">www.teuters.com</a>), Axios (Source: <a href="www.teuters.com">www.teuters.com</a>), The Verge (Source: <a href="www.teuters.com">www.teuters.com</a>), Tom's Guide (Source: <a href="www.teuters.com">www.teuters.com</a>), Axios (Source: <a href="www.teuters.com">www.teuters.com</a>), Tom's Guide (Source: <a href="www.teuters.com">www.teuters.com</a>), Source: <a href="www.teuters.com">www.teuters.com</a>), Axios (Source: <a href="www.teuters.com">www.teuters.com</a>), Tom's Guide (Source: <a href="www.teuters.com">www.teuters.com</a>), and technical analyses (Source: <a href="primary.position.com">primary.position.com</a>). Each claim in this report is backed by specific citations as indicated.

Tags: perplexity ai, large language model, pplx, sonar llm, ai answer engine, llm architecture, open-source llm

#### DISCLAIMER

This document is provided for informational purposes only. No representations or warranties are made regarding the accuracy, completeness, or reliability of its contents. Any use of this information is at your own risk. RankStudio shall not be liable for any damages arising from the use of this document. This content may include material generated with assistance from artificial intelligence tools, which may contain errors or inaccuracies. Readers should verify critical information independently. All product names, trademarks, and registered trademarks mentioned are property of their respective owners and are used for identification purposes only. Use of these names does not imply endorsement. This document does not constitute professional or legal advice. For specific guidance related to your needs, please consult qualified professionals.