

MCP SERVER

NO CODE

CLOUD HOSTED

# Perplexity AI MCP

## Get Web-Grounded Answers with Citations

Perplexity AI MCP connects any AI agent to Perplexity's advanced search and chat models, guaranteeing web-grounded answers with source citations. Stop relying on general LLMs that hallucinate facts; use our tools to run live searches and deep conversations using Sonar and other specialized models directly from your favorite client.

**A+** Quality Score 98.33/100

web-search

ai-chat

citations

real-time-data

natural-language-processing

search-api



# The connectivity layer between AI and the world's software.



Vinkius sits between AI and every application. All communication passes through Vinkius Cloud via the Model Context Protocol (MCP) — with governance, observability, and security at every layer.

# Your AI Connections Run Through Vinkius Cloud

The world's largest  
managed MCP catalog

Vinkius is the connectivity layer where AI connects to the software your business already runs. We handle the hosting, the security, the credentials, the uptime — you get agents that actually do things.

We operate the world's largest managed MCP catalog. Major SaaS platforms, CRMs, databases, and cloud providers — running, monitored, production-ready. This MCP server is hosted and maintained by the Vinkius Cloud for AI Agents.

*The agent doesn't manage credentials, doesn't manage uptime, doesn't manage security. Vinkius does.*

— Architecture principle

---

## Four Pillars of the Vinkius Runtime

### 01 — Security by design

Credentials stay encrypted at rest via AES-256. The AI agent never touches raw keys — they're injected into a sandboxed V8 isolate at runtime. Actions are logged, and connections have an emergency kill switch.

### 03 — Deterministic observability

Eight immutable metrics per endpoint: request volume, p95 latency, error rate, active connections, cost attribution. A live payload feed logs every tool call with mutation detection.

### 02 — Built on MCP Fusion

This MCP server was built with **MCP Fusion**, the open-source framework (Apache 2.0) that powers the entire Vinkius catalog. Schema-as-firewall strips undeclared fields, compiled PII redaction runs at zero overhead, and cryptographic lockfiles produce git-diffable audit trails.

### 04 — Autonomous operations

Servers are deployed, monitored, and patched autonomously. New capabilities and security patches ship weekly. Zero-downtime deployments ensure continuous availability across all managed MCP servers.

**AES-256**

Encryption at rest

**Ed25519**

PKI vault signatures

**24h TTL**

Ephemeral session keys

**V8 Isolate**

Sandboxed execution

---

## One Token. Instant Access.

Every MCP server on Vinkius is accessed through a **Connection Token**. Tokens are generated in the cloud dashboard and produce a unique MCP endpoint URL. Paste this URL into any MCP-compatible client — no SDK required.

A single token can serve **multiple AI clients simultaneously**, or you can issue separate tokens per client for granular access control. Each token tracks its own request count, last activity timestamp, and can be individually enabled or revoked.

MCP ENDPOINT

`https://edge.vinkius.com/{token}/mcp`

Claude



Cursor



VS Code



Windsurf



Grok



Gemini

---

## Security Is the Architecture

Security in Vinkius is not a feature — it's the foundation of the runtime. The gateway enforces multiple independent protection layers between AI agents and third-party APIs.

### 01 — Ed25519 PKI Vault

Every workspace has an Ed25519 Master Key. Session keys are generated ephemerally (24h TTL) and signed by the Master Key. Credentials never leave the vault boundary.

### 02 — V8 Isolate Sandboxing

Tool code runs inside isolated-vm V8 isolates with 64 MB memory caps and per-request timeouts. No filesystem access, no network access except through the SSRF-guarded fetch bridge.

### 03 — SSRF Guard

All outbound HTTP requests are DNS-resolved and validated before execution. Private IP ranges (10.x, 172.16-31.x, 192.168.x, AWS metadata 169.254.x) are blocked at the network layer.

### 05 — Cryptographic Audit Trail

Every request is signed into a SHA-256 hash chain with Ed25519 signatures. Events form a tamper-proof, SIEM-exportable forensic record.

### 04 — DLP & PII Redaction

A ResponseGuard pipeline intercepts every tool response. Configurable redaction patterns strip sensitive fields (emails, SSNs, card numbers) before data reaches the AI agent.

### 06 — Honeypot Trap System

Phantom credentials are injected into isolated environments. If a honeypot is used outside Vinkius infrastructure, the server is quarantined instantly.

## Emergency Kill Switch

EU AI Act Art. 14(1)  
Compliant

The kill switch is an **emergency halt** mechanism — not a simple toggle. When triggered, it executes three actions atomically:

#### 01 — Server deactivated

The MCP server is immediately taken offline across the entire cluster.

#### 02 — All tokens revoked

Every connection token is invalidated. Total lockout — reconnection blocked until new tokens are issued.

#### 03 — WebSocket connections killed

Active connections terminated via Redis pubsub broadcast. Propagates to every runtime node in the cluster.

## Full Visibility. Zero Guesswork.

The Vinkius cloud dashboard includes a full MCP Governance suite — real-time analytics and security controls for production AI operations.

**Control Plane**

KPI dashboard with request volume, latency, success rate, token consumption, and AI-generated operational briefings.

**FinOps**

Cost tracking per tool, payload compression savings, budget optimization signals, and consumption trends.

**Firewall & DLP**

PII redaction activity, sensitive data protection counters, and security event timeline.

**Agent Activity**

Which AI clients are connecting, how often, and what they're doing — real-time session tracking.

**Tool Health**

Slowest and most error-prone tools, with actionable root-cause insights and performance baselines.

**Incident Log**

Error trends, failure rates, status-code breakdowns, and forensic audit trail access.

Get started at [cloud.vinkius.com](https://cloud.vinkius.com) — connect your AI agent in under 60 seconds.

# Perplexity AI MCP

8 tools available

Cloud-hosted on Vinkius

This MCP connects your AI agent to Perplexity's powerful search engine, letting you get web-grounded answers through natural conversation. Instead of relying on a general chat model that might make up facts, this setup forces the AI to cite its sources, giving you real-time context and verifiable data.

Using this MCP means your AI becomes a true research assistant. You can run basic searches using the search tool, or dive deep with specialized models like Sonar Pro for enhanced responses. The entire catalog is managed by Vinkius, so once you connect through any compatible client, all these advanced capabilities are available.

It's built for accuracy. Whether you need a quick check of current market data or a detailed analysis with step-by-step reasoning, the models provide citations and source URLs right in the response. No more switching between Google and ChatGPT; your AI handles it all.

---

## Core Capabilities

### 01 — Search the Live Web

Runs dedicated searches against the web to pull current articles, snippets, and links.

### 03 — Deep Reasoning and Analysis

Uses specialized Sonar reasoning models to break down complex topics into detailed, step-by-step explanations.

### 02 — Generate Cited Responses

Engages chat models (sonar/sonar-pro) to answer questions using real-time data found online, including citations for every claim.

### 04 — Discover Model Options

Lists all available Perplexity AI models so you know exactly which depth or type of model you're running.

# One Click on Vinkius — From Prompt to Execution

Available at [vinkius.com/mcp/perplexity-ai-alternative](https://vinkius.com/mcp/perplexity-ai-alternative) — connect your AI agent in three steps.

- 01 Subscribe to this MCP and input your unique Perplexity API Key.
- 02 Connect the key to any compatible client (like Cursor or Claude).
- 03 Run a search query, and the agent will return web-grounded answers complete with source citations.

The bottom line is that you get reliable, sourced information without having to leave your AI workflow.

---

## Built For

This MCP is essential for researchers and analysts who can't afford to rely on unverified LLM outputs. If your job requires citing sources or knowing what happened yesterday, this tool saves you hours of manual verification.

### Academic Researcher

Needs web-grounded answers with citations for literature reviews and fact-checking research hypotheses.

### Financial Analyst

Searches the current web using domain filters to pull real-time data on market trends or company news.

### Technical Developer

Integrates live search capabilities directly into applications, needing reliable context for code generation tasks.

---

## What Changes When You Connect

- 01 Accuracy is built in. When you use the chat tool, every factual claim comes paired with a citation and source URL, eliminating hallucination.
- 02 Go beyond simple questions. The chat\_with\_reasoning model forces your agent to show its work, providing step-by-step reasoning for complex topics.

- 
- 03 Fine-tune your search results using the search tool's domain filter, letting you limit outcomes to specific types of websites (e.g., only academic journals).

---

  - 04 Access specialized depth with `chat_pro` and `chat_with_reasoning_pro`; these modes give you higher performance models for complex data extraction.

---

  - 05 Monitor your usage via the `get_usage` tool. You always know where you stand so you don't hit unexpected API limits mid-project.
- 

---

## Real-World Applications

### Validating a Research Thesis

A PhD candidate needs to prove three competing theories about climate change. Instead of relying on general chat, they use the search tool and then feed the results into the sonar model. The resulting output provides citations from Nature and Science Daily for every claim, allowing them to build their argument with verifiable data.

### Debugging Complex Code Logic

A developer is stuck on an obscure API integration error. They use the `chat_with_reasoning` model, asking it to explain the protocol failure step-by-step and citing documentation pages from known tech sources.

### Checking Today's Market Data

A financial analyst needs the current stock price of a niche company. They run a targeted search query using the domain filter set to `'financial-news.com'`. The agent pulls up real-time market data and source links, which they can immediately incorporate into their report.

### Preparing a Quarterly Business Review

A marketing manager needs to summarize competitor movements. They run multiple focused searches (e.g., `'competitor X product launch 2025'`) using the search tool, gathering snippets and links from several sources in one pass.

---

---

# Patterns to Avoid

---

## Asking for current facts via general chat

### ✗ AVOID

Pasting 'What is the latest price of Bitcoin?' into a generic AI client. The response will be based on its training data and likely outdated or wrong.

### ✓ INSTEAD

Use the sonar tool or the search tool instead. This forces the agent to perform a live web lookup, guaranteeing that the answer reflects current market prices with direct source attribution.

---

## Overloading the prompt with vague concepts

### ✗ AVOID

Prompting: 'Tell me about AI.' The response is generic marketing fluff and provides no actionable data or sources.

### ✓ INSTEAD

Use chat\_with\_reasoning, specifying a narrow focus like 'Compare BERT vs GPT architecture'. This forces deep analysis and structured reasoning instead of broad concepts.

---

## Trying to filter results manually

### ✗ AVOID

Running a search and then having to copy-paste 15 different links into a spreadsheet for verification.

### ✓ INSTEAD

Use the search tool with domain filtering. This narrows down the initial pool of links, giving you highly relevant snippets from specific sites right out of the gate.

---

## The Right Fit

Use this MCP if your primary concern is verifiability and real-time context. If your task requires knowing what happened on the internet in the last 24 hours, or citing a source for every factoid, this is your tool. Don't use it if you need pure creative brainstorming (use a general chat model) or if you are writing fiction (the web search results will ruin the tone). If you just want to summarize an article you already have, a simple text-based LLM might suffice. But when the source data matters—like in finance, law, or research—you need the citations and live context this MCP provides.

---

---

## The Problem with Outdated AI Answers

Today, if you ask a standard LLM about a piece of technology or market data, it gives you an answer based on everything it was trained on. That means the information is stale; it's history, not fact. You end up wasting time cross-referencing those answers with actual news sites and databases just to confirm if they're right.

With this MCP, your AI doesn't guess. It executes a live web search using the dedicated search tool. It pulls real snippets from the current internet—and it shows you exactly which websites provided that information. You get data-backed answers instantly.

---

## Get Cited Answers with sonar

Manual verification means opening tabs, copying key quotes, and pasting them into a citation tracker. It's slow, tedious, and you inevitably miss something important when switching between windows.

The sonar model changes that completely. You ask the question once, and it returns the complete answer with all sources linked directly in the output. The process is one prompt, one verifiable result.

---

# Perplexity AI - 8 Tools Available

These tools give you granular control over the entire Perplexity search stack, from simple queries to deep reasoning and model discovery.

#	TOOL	DESCRIPTION
01	chat	Sends a chat message to various Perplexity models and gets an answer with web citations.
02	chat_pro	Sends a chat message specifically to the Sonar Pro model for enhanced, highly detailed responses with citations.
03	chat_with_reasoning	Sends a message to the Sonar Reasoning model, forcing it to return a detailed step-by-step chain of thought and citations.
04	chat_with_reasoning_pro	Sends a complex message to the Sonar Reasoning Pro model for deep analysis, providing both reasoning steps and citations.
05	get_usage	Checks your API usage statistics so you can monitor consumption limits.
06	list_models	Lists all available Perplexity models, showing their IDs and specific capabilities for use in other tools.
07	search	Performs a web search using the Perplexity Search API and returns snippets with source URLs and citations.
08	sonar	Sends a message to the core Sonar model, which provides answers grounded in real-time web search results.

---

## See It in Action

Real prompts you can use once this MCP is connected to your AI agent through Vinkius Cloud.

### **U** Search the web for 'latest advances in quantum computing 2025'.



Found recent articles on quantum error correction breakthroughs, Google's Willow chip achieving below-threshold error rates, and IBM's 1000+ qubit processors. Key sources include Nature, MIT Technology Review and Science Daily.

### **U** Ask Sonar: What is the current price of Bitcoin?



Bitcoin is currently trading at \$87,234 USD according to CoinMarketCap data from today. The price has increased 3.2% in the last 24 hours with a market cap of \$1.73 trillion. Sources: CoinMarketCap, CoinGecko.

### **U** Send a chat to sonar-pro asking 'Explain how transformers work in NLP' with `return_related_questions` enabled.



Sonar-pro responded with a detailed explanation of transformer architecture including self-attention, multi-head attention, positional encoding and the encoder-decoder structure. Related questions returned include: 'What is the difference between BERT and GPT?', 'How does self-attention work?' and 'What are transformer limitations?'

---

## Frequently Asked Questions

### **01** How does the Perplexity AI MCP handle real-time data?

The MCP uses dedicated search tools (like 'search' or 'sonar') that execute live queries against the web, ensuring answers are based on current information rather than old training data.

---

**02 Is the chat tool suitable for complex research?**

Yes. For deeper analysis, use `chat_with_reasoning` or `chat_with_reasoning_pro`. These models force the AI to show its reasoning steps before giving the final answer.

---

**03 What is the difference between 'sonar' and general chatting?**

General chat relies on trained knowledge; sonar uses a specialized process that forces web-grounding. It guarantees citations, which you won't get from basic conversation tools.

---

**04 How can I restrict the search results?**

The MCP allows you to use the search tool with domain filtering. You can limit your results to specific types of sites (e.g., only academic domains) for better accuracy.

---

**05 Do I need a separate API key for each model?**

No, you connect one API Key through the MCP, and it gives you access to all available models listed via `list_models`. You just select which tool or mode you want to use.







---

# Go Live in 60 Seconds

Get your connection token from [cloud.vinkius.com](https://cloud.vinkius.com), then paste the endpoint URL into any MCP-compatible client.

YOUR MCP ENDPOINT

```
https://edge.vinkius.com/[TOKEN]/mcp
```

CLIENT	WHERE TO CONFIGURE
 <b>Claude AI</b>	Profile → Customize → Connectors → "+" → Add custom connector → Paste endpoint
 <b>Cursor</b>	Settings → Features → MCP Servers → "+ Add New MCP Server" → Type: SSE → Paste endpoint
 <b>VS Code</b>	Ctrl/Cmd+Shift+P → "MCP: Add Server" → add <code>"perplexity-ai-alternative": { "url": "..." }</code>
 <b>Windsurf</b>	MCP Settings → <code>mcp_settings.json</code> → Add endpoint URL
 <b>ChatGPT</b>	Settings → Tools & plugins → Add MCP server → Paste endpoint
 <b>Gemini</b>	Extensions → Add MCP Server → Paste endpoint URL

## ASK AN AI ABOUT THIS

Let your preferred AI explain this MCP server

-  **Ask ChatGPT** 
-  **Ask Claude** 
-  **Ask Perplexity** 
-  **Ask Gemini** 
-  **Ask Grok** 

READY TO CONNECT

# Perplexity AI is live on Vinkius Cloud.

Get your connection token, paste it into your AI agent, and  
start building. No SDK. No deployment. Just results.

[Start at cloud.vinkius.com](https://cloud.vinkius.com) →

[vinkius.com](https://vinkius.com) · [support@vinkius.com](mailto:support@vinkius.com)

### INDEPENDENT PLATFORM DISCLAIMER

Vinkius is an independent platform and is not affiliated with, endorsed by, sponsored by, verified by, or otherwise authorized by Perplexity AI. All third-party trademarks, logos, and brand names are the property of their respective owners. Their use in this document is strictly for informational purposes to identify service compatibility and interoperability.

### DOCUMENT INFORMATION

Generated	June 2026
MCP Server	Perplexity AI MCP
Server ID	019d846b-3d68-71c0-a084-c532f9568d49
Platform	Vinkius Cloud for AI Agents
Endpoint	<code>https://edge.vinkius.com/{token}/mcp</code>

### LICENSE & USAGE

This document is generated automatically by the Vinkius PDF Engine. Content reflects the MCP server configuration at the time of generation and may change as updates are deployed. For the most current information, visit [vinkius.com/mcp/perplexity-ai-alternative](https://vinkius.com/mcp/perplexity-ai-alternative).