

MCP SERVER

NO CODE

CLOUD HOSTED

Anthropic MCP for AI Agents

Managing LLM Model Access and Token Counting for Prompt Engineering

Anthropic connects your AI agents to Claude models, letting you manage conversations and control costs without leaving your workflow. You can discover available models, count tokens before running a prompt, or submit large batches of requests for cost-effective processing.

A+ Quality Score 95.83/100

llm

model-discovery

token-counting

natural-language-processing

prompt-engineering

api-management



The infrastructure that powers AI agents in the real world.



Vinkius connects AI to the world's software through secure, enterprise-grade infrastructure — enabling real-world execution at scale, built on the Model Context Protocol (MCP).

Your AI Connections Run Through Vinkius Cloud

The world's largest
managed MCP catalog

Vinkius is the cloud infrastructure where AI agents connect 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 — connect your AI agent in under 60 seconds.

Anthropic MCP

6 tools available

Cloud-hosted on Vinkius

Need to use Claude's power but don't want to switch between different API interfaces? This MCP gives your AI agent direct access to Anthropic's entire model suite. You can send conversations and get responses using natural language, all managed through one place. It makes sense for developers or ML engineers who need reliable ways to test models, estimate costs, or process huge volumes of prompts efficiently.

For example, instead of running individual API calls for every prompt, you submit a batch job that your agent handles asynchronously. Plus, if you're worried about spending too much on context windows, you can use the token counting tool first to figure out exactly how big your messages are before hitting send. Finding and managing these different model options is simplified by connecting through Vinkius, giving all your AI clients a single catalog point of access.

Core Capabilities

01 — Send Conversations

Your agent sends natural language prompts to Claude models and receives the response text.

03 — Estimate Token Usage

Your agent counts the input tokens of a message before sending it to estimate costs or check context limits.

05 — Check Batch Status

Your agent monitors a batch job's progress and reports if the request succeeded or failed using ``get_batch_message``.

02 — Discover Available Models

You list every model Anthropic offers, getting their IDs and capabilities for use in your prompts.

04 — Process Batches of Prompts

You submit multiple, independent requests at once for cost-effective, asynchronous processing using ``create_batch_message``.

One Click on Vinkius — From Prompt to Execution

Available at vinkius.com/mcp/anthropic-alternative — connect your AI agent in three steps.

- 01 Subscribe to this MCP and paste your Anthropic API key.
- 02 Select this MCP in any compatible AI client, like Cursor or Windsurf.
- 03 Your agent now uses the integrated tools to manage model calls—whether sending a single message, checking tokens, or running a batch job.

The bottom line is that you treat Anthropic's entire suite of models as just another set of tools inside your AI client.

Built For

This MCP targets developers, data scientists, and product managers who build applications relying heavily on large language model APIs. If you spend time writing API wrappers or manually calculating prompt costs, this is for you.

ML Engineer

They use the tool to discover available models and run batch requests to compare performance across multiple prompts efficiently.

Software Developer

They connect this MCP to build features that require reliable message sending, often checking token counts first before integrating the final prompt logic.

Product Manager

They use it to review model output quality and track overall usage metrics for batch processing jobs via natural conversation.

What Changes When You Connect

- 01 Estimate costs before you send anything. Use the `count_tokens` tool to know exactly how many tokens your message will consume, preventing unexpected API overages.

-
- 02 Process massive volumes of data without manual coding. The `create_batch_message` and `get_batch_message` tools let you submit thousands of prompts asynchronously for cost-effective bulk processing.

 - 03 Never worry about model selection again. Use the `list_models` tool to see every available Claude version, their IDs, and specific capabilities in one spot.

 - 04 Build robust agents that handle failure gracefully. You can cancel a job using `cancel_batch_message` if you realize you started processing too many requests by accident.

 - 05 Maintain workflow simplicity. Your agent handles the complex API calls, letting you interact with Anthropic's models using plain conversation rather than raw code.
-

Real-World Applications

Analyzing large datasets for sentiment

A data scientist needs to analyze 5,000 customer reviews. Instead of writing a loop, they use the `create_batch_message` tool via their agent, submitting all prompts at once and tracking progress using `get_batch_message` until everything is complete.

Implementing cost guardrails

A developer integrating Claude needs to ensure prompts don't exceed a 400-token limit. They use the `count_tokens` tool first; if the count is too high, the agent automatically tells them to shorten the message before calling `send_message`.

Creating content for A/B testing

A product team needs to generate 10 variations of a marketing headline. They use the agent to first discover all relevant models via `list_models`, then send messages using `send_message` to test different tones and styles, reviewing output results in conversation.

Handling accidental large runs

An ML engineer accidentally triggers a batch request for 10,000 prompts. Realizing the cost implications immediately, they use the `cancel_batch_message` tool to stop processing before it wastes credits.

Patterns to Avoid

Calling APIs one by one

X AVOID

Trying to send 50 different prompts individually using only `send_message`. This is slow, expensive, and inefficient for bulk work.

✓ INSTEAD

For batch processing, always use `create_batch_message` and then monitor progress with `get_batch_message`. This handles the volume asynchronously.

Ignoring token limits

X AVOID

Sending a massive document to Claude without checking the input size first. The API might fail or cost more than necessary.

✓ INSTEAD

Always run `count_tokens` on your source material and prompt together. This ensures you stay within context window limits.

Assuming model availability

X AVOID

Hardcoding a specific model ID in your application without checking if it's still active or available.

✓ INSTEAD

Use the `list_models` tool first. This guarantees you know which models are currently accessible and what their IDs are for reliable code.

The Right Fit

You should use this MCP when your workflow requires controlled, high-volume interaction with Claude's APIs. If your primary need is simply to chat or ask a single question, `send_message` works fine. However, if you are building an application that needs enterprise reliability—like processing thousands of records, managing budget constraints, or needing model version discovery—you must use the batch tools (`create_batch_message`, etc.). Don't use this if your goal is just to test a single prompt once; use it when scale and cost management matter. This MCP makes Anthropic reliable for production pipelines.

Managing Anthropic Claude Model Access with the Anthropic MCP

Today, using Claude often means jumping between multiple interfaces or writing complex API scripts just to handle basic tasks like counting tokens or running a batch job. You have to manually manage model versions and track costs across different endpoints, which is tedious and prone to failure.

With this MCP, your agent handles all that complexity for you. Instead of dealing with raw HTTP requests, you simply ask the tool to count tokens or list models. The punchline? You get reliable access and cost visibility without writing any boilerplate API code.

Anthropic Claude Batch Processing via the Anthropic MCP

Manually processing hundreds of prompts involves creating massive, unwieldy scripts that run sequentially. If one prompt fails, the whole process often halts, and you have no easy way to track which ones succeeded or failed.

Now, you use `create_batch_message`. You submit all your independent requests in a single job, letting the MCP handle the queuing and tracking. This means reliable throughput and full visibility into every result.

6 Tools for Anthropic LLM Batch Management and Token Counting

Use these tools to send single messages, list models, check token counts, or process massive batches of requests with Claude.

#	TOOL	DESCRIPTION
01	<code>cancel_batch_message</code>	Stops a large, ongoing message batch request if you submitted it by mistake, saving costs.
02	<code>count_tokens</code>	Calculates the total input tokens for a given message array, useful for estimating API costs and context limits.
03	<code>create_batch_message</code>	Submits multiple independent prompts to Claude in one go, which is more cost-effective than running them individually.
04	<code>get_batch_message</code>	Checks the current status of a batch job using its ID, reporting success counts and results.
05	<code>list_models</code>	Retrieves a list of all Claude models available, including their IDs and specific capabilities.
06	<code>send_message</code>	Sends a single message prompt to Claude with customizable parameters like system prompts and temperature.

See It in Action

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

U What's the difference between Sonnet and Opus for summarizing documents?



Model Comparison: Claude Versions

MODEL	BEST FOR	CONTEXT WINDOW	COST/TOKEN ESTIMATE
claude-sonnet	General use, balance of speed and capability.	200K tokens	Medium
claude-opus	Complex reasoning, intricate data analysis.	200K tokens	High
claude-haiku	Simple classification, rapid initial triage.	200K tokens	Low

Use the model discovery tool to confirm current IDs.

U I have three articles I need summarized and categorized by industry. How do I do this?



Batch Processing Plan:

You can submit all three requests simultaneously using a batch job.

- Action:** `create_batch_message` (Submit the 3 prompts)
- Tracking:** Use `get_batch_message` with the returned ID to check status.
- Result:** You'll get a list of three structured summaries and categories, allowing you to process them all at once.

U Before I send this 10-page report for analysis, how many tokens will it use?



Token Count Estimate

Please provide the model ID and the text. For instance:

- **Input Text:** [Paste your text here]
- **Model Used:** cLaude-sonnet ...
- **Estimated Tokens:** 895 input tokens.
- **Warning:** This is above the recommended context limit for this type of analysis. Consider summarizing the report first.

Frequently Asked Questions

01 How do I manage model costs when using Anthropic through the Anthropic MCP?

You control costs by proactively checking token usage before sending anything. The `count_tokens` tool lets you estimate input size, and the batch tools make large-scale processing much more efficient than calling APIs individually.

02 Can this Anthropic MCP handle thousands of prompts at once?

Yes. By using the batch creation tools, your agent can submit massive jobs asynchronously. You simply monitor the status with `get_batch_message` until all requests are complete.

03 What if a large batch job fails or runs too long?

You've got options to manage that. If you run into an issue, you can use the tool to check the status of your batch and even stop processing early with `cancel_batch_message` to save credits.

04 Do I need to know all my model IDs beforehand?

No. You can use the dedicated function within this MCP to list every available Claude model ID, making sure your agent is always pointing to a current and working version.

05 Is using this MCP better than writing custom API calls for Anthropic?







Most times, yes. This MCP wraps the complexity into simple actions within your agent, letting you focus on what the AI does with the data instead of how to connect to the API.

Go Live in 60 Seconds

Get your connection token from 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
 Claude AI	Profile → Customize → Connectors → "+" → Add custom connector → Paste endpoint
 Cursor	Settings → Features → MCP Servers → "+ Add New MCP Server" → Type: SSE → Paste endpoint
 VS Code	Ctrl/Cmd+Shift+P → "MCP: Add Server" → add <code>"anthropic-alternative": { "url": "..."} </code>
 Windsurf	MCP Settings → <code>mcp_settings.json</code> → Add endpoint URL
 ChatGPT	Settings → Tools & plugins → Add MCP server → Paste endpoint
 Gemini	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

Anthropic 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 · 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 Anthropic. 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	Anthropic MCP
Server ID	019d8416-47d9-732a-983f-276099624a35
Platform	Vinkius Cloud for AI Agents
Endpoint	https://edge.vinkius.com/{token}/mcp

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/anthropic-alternative.