

MCP SERVER

NO CODE

CLOUD HOSTED

Anthropic MCP for AI Agents

Managing High-Volume Text Analysis and Prompt Engineering

The Anthropic MCP lets your AI client connect directly to Claude models. You can send prompts for complex reasoning or manage huge volumes of requests through asynchronous batch processing. It also keeps tabs on your account's rate limits and estimates costs before you hit send.

A Quality Score 94.17/100

llm-integration

natural-language-processing

batch-processing

prompt-engineering

api-access



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

8 tools available

Cloud-hosted on Vinkius

Working with advanced language models like Claude requires more than just sending a single prompt; it demands careful resource management, especially when you're running high-volume tasks. This MCP lets your AI agent talk directly to the Anthropic API, giving you granular control over every part of the process. Need to analyze thousands of documents for sentiment? You can set up large message batches and run them asynchronously, which drastically cuts down on token costs compared to live calls. Plus, it's built with monitoring in mind; your agent will tell you exactly what your current rate limits are and calculate how much a specific job is going to cost before you commit to running it. Connecting this MCP through Vinkius gives you access to this powerful Claude integration alongside thousands of other tools for your AI client.

Core Capabilities

01 — Send multi-turn prompts to Claude

You can send continuous messages and system instructions to any available Claude model.

03 — Estimate API costs

The MCP calculates the expected cost based on your prompt token counts and current Anthropic pricing structure.

05 — List available models

You can see a list of all Claude models currently supported by the API, including technical capabilities.

02 — Run high-volume message batches

Create and manage large groups of requests for non-realtime processing, which saves you money on tokens.

04 — Monitor usage limits

It tracks your account's Requests Per Minute (RPM) and Tokens Per Minute (TPM) to prevent unexpected throttling.

One Click on Vinkius — From Prompt to Execution

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

- 01 Subscribe to this MCP and plug in your Anthropic API Key.
- 02 Tell your AI client what you want to do—whether that's sending a quick message or initiating a large batch job.
- 03 Your agent interacts with the tools, handles the requests, and provides results like cost estimates or status updates.

The bottom line is you get direct, rate-managed access to Claude's full capabilities without having to worry about token limits or complex API setup.

Built For

Anyone doing serious work with LLMs needs this. If your job involves large amounts of text processing—anything from research analysis to content generation at scale—you need the cost control and throughput capacity this MCP provides.

ML Engineer

You use it to test prompt variations quickly or run large-scale evaluations by setting up message batches for significant cost reduction.

Content Manager

You send multi-turn messages to Claude models, getting consistent writing style and tone across multiple drafts or articles.

Data Scientist

You monitor API rate limits while running experiments and use the MCP's built-in cost estimation tools before committing compute resources.

What Changes When You Connect

- 01 Cut costs by 50% on large jobs. Setting up a message batch using `create_batch` slashes your token expense when you need to process thousands of items.

-
- 02 Stay operational without hitting limits. By checking rates with `check_rate_limits`, you know exactly how many requests and tokens are left before throttling hits.

 - 03 Manage complex projects efficiently. You can list all available models using `list_models` so your agent knows which Claude version is best for the job.

 - 04 Control everything after launch. Need to stop a runaway job? Use `cancel_batch`. If something breaks, you know how to get the status via `get_batch`.

 - 05 Speed up development cycles. The built-in cost estimation feature lets your agent give you an exact dollar figure for any prompt before running it.
-

Real-World Applications

Analyzing customer feedback at scale

A marketing team needs to analyze 10,000 pieces of customer survey text. Instead of sending them one by one, they use `create_batch` to process all the data overnight, saving money and getting results back via `get_batch_results` the next morning.

Automating content localization checks

A global team needs to check if 50 different articles are ready for publication. They first use `list_models` to pick the best Claude version and then send a batch job to ensure every article passes necessary formatting rules.

Building a multi-step character bot

A developer is building an interactive story tool. They use `create_message` repeatedly to send system prompts, managing context and ensuring the agent maintains a consistent personality across dozens of turns.

Stress testing an application's limits

An ML engineer needs to know how many API calls their app can handle per minute. They use `check_rate_limits` early in development and then monitor the usage data directly through the MCP.

Patterns to Avoid

Sending everything via live prompts

✗ AVOID

When you need to process 50,000 records, running them all through individual `create_message`` calls is slow and costs a ton of money.

✓ INSTEAD

Instead, use the batch features. Run a large job by calling `create_batch``, then check its progress with `get_batch``. This saves 50% on tokens.

Assuming you know your limits

✗ AVOID

Writing code that runs without ever checking the current API status means hitting rate limits and failing unexpectedly.

✓ INSTEAD

Always use `check_rate_limits`` first. This gives your agent real-time data on how many requests per minute are available right now.

Forgetting to clean up jobs

✗ AVOID

You start a massive batch job but forget about it, and the resources keep running until you manually stop them.

✓ INSTEAD

Use `list_batches`` to see what's running, and if necessary, use `cancel_batch`` to shut down any unnecessary processes immediately.

The Right Fit

Use this MCP if your primary need is high-throughput processing or cost control when working with Claude. If you are doing a few single-prompt interactions occasionally, the base API might suffice. But if you're dealing with thousands of records, complex multi-turn conversations, or just need to know exactly how much it will cost before you start, this MCP is essential. Don't use it if your goal is simply connecting Claude via basic chat; use a general messaging tool for that. You *must* use the batch tools like `create_batch`` and its related functions (`get_batch`` , `list_batches``) when scaling up to manage costs effectively.

Anthropic MCP: Scaling Text Analysis with Claude Models

Right now, processing large datasets in a language model is a nightmare. You have to copy data into the prompt, hit send, wait for the response, and then repeat that process manually hundreds or thousands of times. This isn't just tedious; it's expensive because every single API call counts toward your rate limit.

With this MCP, you let your agent handle the heavy lifting. Instead of sending prompts one by one, your agent initiates a batch job via `create_batch`. You send the entire payload once, and Claude processes everything asynchronously in the background. The result? A massive drop in cost and a stable workflow that handles volume without breaking.

Anthropic MCP: Controlling API Usage with Anthropic Models

The biggest headache is always uncertainty. You don't know if you have enough tokens for the job, or if your rate limit will drop mid-process, leaving your whole pipeline hanging. Manually tracking usage across different tools and documentation is a full-time job.

This MCP gives you visibility. Your agent can check your limits using `check_rate_limits` before starting work, and it can estimate the total cost for any complex prompt run. You know exactly what you're spending and when to slow down.

8 Tools for Anthropic Message Batching and Rate Management

Use these tools to manage Claude message batches, check API rates, estimate token costs, and interact with various Claude models programmatically.

#	TOOL	DESCRIPTION
01	<code>cancel_batch</code>	Stops a message batch job that you started but no longer need running.
02	<code>check_rate_limits</code>	Retrieves your current limits for requests and tokens from the Anthropic account.
03	<code>create_batch</code>	Starts a message batch process, which saves you 50% on token costs compared to live API calls.
04	<code>create_message</code>	Sends a direct prompt to Claude and gets the resulting text response back.
05	<code>get_batch_results</code>	Pulls the completed results from a message batch that has already finished processing.
06	<code>get_batch</code>	Checks the current status of a specific, ongoing message batch job.
07	<code>list_batches</code>	Shows you a list of all the message batches that have been created on your account.
08	<code>list_models</code>	Retrieves an accurate list of every Claude model currently available for use.

See It in Action

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

U What models do I have access to right now?



Anthropic API Models Available

The following Claude versions are currently active in your account, optimized for different use cases:

- **Claude 3.5 Sonnet:** Best balance of speed and capability.
- **Claude 3 Opus:** Reserved for the most complex reasoning tasks (e.g., legal review).
- **Claude 3 Haiku:** Ideal for simple classification or quick summarization.

Always check the documentation to see which model best fits your budget and complexity needs.

U Help me figure out how much it costs to process 1 million tokens.



Cost Estimate Report

Based on current Anthropic pricing, processing a hypothetical payload of 500,000 input tokens and 20,000 output tokens would cost approximately **\$X.XX**.

Note: This estimate does not include usage fees for the message batch management itself.

U I need to run a sentiment analysis on these 50 articles.



I've queued up your request. A new Message Batch (ID: msgbatch_d1c2f) has been created with 50 documents for processing.

- **Status:** Pending
- **Estimated Completion:** Within 4 hours.
- You can track progress at any time using the `get_batch` tool.

Frequently Asked Questions

01 How do I handle processing thousands of prompts with Anthropic MCP?

Use batching. Instead of running single messages, use the message batch functionality to process huge volumes of data asynchronously. This saves you money and keeps your workflow stable.

02 Does the Anthropic MCP track my rate limits?

Yes, it does. You can ask the MCP to check your current Requests Per Minute (RPM) and Tokens Per Minute (TPM) usage instantly, preventing unexpected service slowdowns.

03 Is this better than just using the Anthropic website?

Absolutely. This MCP gives you programmatic control over everything—from starting batches to monitoring cost. It moves beyond a simple chat interface and into scalable engineering workflows.

04 What if I want to stop a big job before it finishes?

You can list all your running jobs using the MCP, find the ID of the batch you want gone, and then use the cancellation tool. It stops processing immediately.

05 Does Anthropic MCP calculate my token usage cost?







Yes, it includes a built-in cost estimator that calculates your expected spending based on input and output tokens before you run the job. This is crucial for budgeting large projects.

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": { "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	019d754e-55d5-702e-b5e1-12b68627b1ba
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.