

MCP SERVER

NO CODE

CLOUD HOSTED

Mistral AI MCP

Control Inference, Embeddings, and Agents from One Place

Mistral AI connects your agent to a full suite of state-of-the-art language model capabilities. You can run complex conversational tasks, generate dense text embeddings for search, or perform specialized code completions like Fill-in-the-Middle (FIM). It also allows you to audit available models and trigger custom multi-step AI workflows.

A+ Quality Score 100/100

llm

inference

rag

embeddings

natural-language-processing

model-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 — connect your AI agent in under 60 seconds.

Mistral AI (Frontier LLMs & Embeddings) MCP

7 tools available

Cloud-hosted on Vinkius

This MCP lets your agent interact with Mistral's advanced model suite without needing complex SDK setup. You get full control over running different types of inference, whether it's general chat or highly specialized tasks like code completion. Need to power a semantic search? Use the embedded tools to calculate vector representations from any text block. For building autonomous systems, you can trigger custom multi-step workflows and even check content against safety policies before deployment. If your current development stack uses various API keys for different providers, Vinkius brings all these advanced Mistral capabilities together into one place. You connect once through the Vinkius catalog and immediately gain access to this comprehensive set of tools.

Core Capabilities

01 — Run Conversational Inference

Execute high-fidelity chat completions using Mistral's various models, giving you detailed control over system instructions and message history.

03 — Complete Code Logic (FIM)

Fill in missing sections of code, bridging the logical gap between existing prefixes and required suffixes.

05 — Inspect Model Metadata

List all available Mistral AI models and retrieve detailed configuration settings to determine the best model for a job.

02 — Calculate Text Embeddings

Generate dense numerical vectors for any text. This powers semantic search engines and knowledge retrieval systems.

04 — Execute Agent Workflows

Trigger multi-step, autonomous agent processes that handle complex reasoning tasks on your behalf.

One Click on Vinkius — From Prompt to Execution

Available at vinkius.com/mcp/mistral-ai-frontier-llms-embeddings — connect your AI agent in three steps.

- 01 Subscribe to this MCP on Vinkius and enter your unique Mistral AI API Key.
- 02 Select your preferred connection point, like Claude or Cursor, and activate the Mistral tools within your agent's context.
- 03 Call a specific tool—for example, `generate_embeddings`—and pass the required text data to get immediate results.

The bottom line is you talk to your AI client normally, but it uses this MCP to handle all the complex model calls and data processing behind the scenes.

Built For

This connects with ML Engineers who are tired of managing dozens of API keys for different tasks. It's also perfect for AI Developers building prototypes who need to test multiple frontier models quickly, and Data Scientists needing consistent embedding generation across services.

ML Engineer

Tests model performance by calling `list_models` to compare metadata or runs `generate_embeddings` to verify vector distribution for new datasets.

AI Developer

Builds application features by using `chat_completion` and `fim_completion`, allowing the agent to perform both general conversation and specialized coding tasks in one go.

Data Scientist

Uses `generate_embeddings` to map large document sets into searchable vectors, then uses `get_model` to confirm the correct embedding model ID.

What Changes When You Connect

-
- 01** Deep control over model selection. Use `list_models` to compare different Mistral variants and get the exact metadata you need before running `chat_completion`.

 - 02** Build high-performance search features instantly. The `generate_embeddings` tool lets your agent convert any text into searchable vectors, making RAG pipelines easy.

 - 03** Improve code quality with FIM. Instead of generic auto-complete, `fim_completion` fills in logical gaps, requiring you to provide only the start and end points.

 - 04** Manage complexity through automation. You don't write multi-step API calls; you just call `agent_completion` to run a sophisticated workflow.

 - 05** Ensure safety compliance upfront. `moderate_content` checks inputs against toxicity policies, giving you confidence that the content is safe before it hits production.
-

Real-World Applications

Building a Document Search Portal

A data scientist needs to index 10,000 legal documents. They use `generate_embeddings` to convert all text into vectors and then pass the list of model IDs to `get_model`, confirming that the embedding process is using the correct, stable version.

Testing Agent Logic Flow

A researcher wants a multi-step agent to analyze market sentiment. They call `agent_completion`, which runs a sequence of reasoning steps and returns a final structured report without the developer needing to write orchestration code.

Creating a Code Copilot Feature

A developer wants an in-IDE assistant. They use `fim_completion` when they type `'def calculate_fib(n):'` and only need to write the closing bracket; the tool fills in all the complex loop logic.

Pre-deployment Content Scrubbing

A content team uploads user reviews that might contain prohibited material. Before storing them, they use `moderate_content` to run safety checks on every single entry, rejecting anything that fails the compliance filter.

Patterns to Avoid

Using a generic LLM endpoint for everything

✗ AVOID

Trying to generate code snippets using simple `chat_completion` prompts often results in incomplete or overly generalized functions, making them useless for production code.

✓ INSTEAD

Always use `fim_completion` when coding. By explicitly providing the prefix and suffix, you force the model to fill only the missing logical gap, resulting in precise, functional code.

Hardcoding Model IDs

✗ AVOID

Writing your application to only accept `'mistral-large-v1'` when a newer, better version is released, causing immediate breakage.

✓ INSTEAD

First, call `list_models`. This gives you the current inventory and metadata for all available Mistral AI models, ensuring your system always knows which IDs are valid.

Forgetting content moderation

✗ AVOID

Sending user-generated data directly to a database or agent without checking it first, risking compliance violations or toxic output.

✓ INSTEAD

Always route sensitive inputs through `moderate_content`. This mandatory safety check verifies the text against toxicity policies before any other process runs.

The Right Fit

Use this MCP if your primary need is deep control over multiple, distinct AI tasks: chat conversation AND vector math AND code logic. You need a single source to manage everything from chat_completion to generate_embeddings and fim_completion.

Don't use it if you only need basic text generation or simple API calls for just one model type. If your needs are limited to just chatting, an alternative general-purpose LLM connector might suffice. But because this MCP handles specialized tasks like FIM completion and dedicated embedding calculation (generate_embeddings), it becomes indispensable when building complex, multi-faceted AI applications.

Handling model APIs used to be a messy process.

Today, integrating multiple specialized AI functions means writing boilerplate code for each one. You'll manage separate SDK calls just to chat with an agent, then another set of keys and logic just to generate embeddings, and yet a third block of code is needed for specialized tasks like code filling.

With this MCP, you connect your agent once, regardless of the task. Whether it's complex multi-step reasoning via agent_completion or simply calculating vectors with generate_embeddings, the whole process runs through one standardized connection point.

Mistral AI (Frontier LLMs & Embeddings) MCP gives you true model control.

You no longer have to guess which model is best. You can first call list_models, pull the metadata for specific variants with get_model, and then decide if a general chat_completion or specialized fim_completion is appropriate for the job at hand.

The result is an application that behaves like one cohesive system, not a collection of bolted-on API calls. You build reliability into your stack.

Mistral AI (Frontier LLMs & Embeddings) – 7 Tools

These tools allow your agent to perform specific tasks like calculating embeddings or running code completions directly through the Mistral AI model suite.

#	TOOL	DESCRIPTION
01	<code>chat_completion</code>	Runs conversational inference using Mistral AI's chat completion models for structured text output.
02	<code>generate_embeddings</code>	Calculates numerical vectors from provided text data using a dedicated embedding model.
03	<code>list_models</code>	Retrieves an inventory of all currently available Mistral AI models that the client can use.
04	<code>get_model</code>	Fetches specific details and metadata about one particular Mistral AI model ID.
05	<code>fim_completion</code>	Generates missing code logic by filling in the gap between a defined prefix and suffix of code.
06	<code>moderate_content</code>	Checks user-provided content against safety rules to ensure compliance before processing or deployment.
07	<code>agent_completion</code>	Initiates and manages a custom, multi-step autonomous agent workflow defined by Mistral AI.

See It in Action

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

U Run a chat completion using 'mistral-large-latest' to summarize this research paper: [text]



Executing conversational inference... I've processed the text using Mistral Large. Here is a 3-point summary: 1) The study identifies a 15% efficiency gain in transformers, 2) Data centers show reduced latency with the new protocol, and 3) Results are consistent across test sets. Would you like a more detailed breakdown?

U Generate code to complete this gap: Prefix 'def calculate_fib(n):', Suffix 'return sequence'



Executing Codestral FIM completion... I've generated the missing logic: 'sequence = [0, 1]; for i in range(2, n): sequence.append(sequence[i-1] + sequence[i-2])'. This perfectly bridges your prefix and suffix. Would you like me to generate unit tests for this function?

U List all available Mistral models and their IDs



I've retrieved the Mistral model inventory. Highlights include 'mistral-large-latest' (General purpose), 'mistral-small-latest' (Fast inference), 'codestral-latest' (Coding), 'pixtral-12b-2409' (Multimodal), and 'mistral-embed' (Embeddings). Which model would you like to inspect further?

Frequently Asked Questions

01 How do I use Mistral AI (Frontier LLMs & Embeddings) for semantic search?

You calculate dense numerical vectors using the generate_embeddings tool. This process converts raw text into a vector representation that powers your semantic search database.

02 Can I use Mistral AI (Frontier LLMs & Embeddings) for code filling?

Yes, you use `fim_completion`. You provide the existing code prefix and suffix, and the tool generates the missing logic in between.

03 What is the purpose of `list_models` with Mistral AI (Frontier LLMs & Embeddings)?

`list_models` provides an inventory of all active Mistral models. This helps you identify which model ID to use for a specific task, like choosing between 'mistral-large' and 'mistral-small'.

04 Does Mistral AI (Frontier LLMs & Embeddings) handle safety checks?

Yes, you can run `moderate_content`. This tool runs the content through rigorous toxicity policies to verify compliance before you deploy or store it.

05 Is `chat_completion` better than `agent_completion` for complex tasks?

No. Use `chat_completion` for single-turn conversations. If a task requires multiple steps of reasoning, calling `agent_completion` is the correct method for autonomous execution.

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 `"mistral-ai-frontier-llms-embeddings": { "url": "..."}`



Windsurf

MCP Settings → `mcp_settings.json` → 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

Mistral AI (Frontier LLMs & Embeddings) 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 Mistral AI (Frontier LLMs & Embeddings). 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	Mistral AI (Frontier LLMs & Embeddings) MCP
Server ID	019d75d5-9fe5-730c-88ed-3da746f21d8c
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/mistral-ai-frontier-llms-embeddings.