How to use Al Assistants with your Monitoring system

Alexander Marshalov



Agenda

What is MCP and how it works How to use **VictoriaMetrics** MCP Server Features overview Usage scenarios and examples Integration with Cloud Future plans About VictoriaLogs MCP Server



What is MCP and how it works?



- MCP Model Context Protocol
- Open standard that describes a way for AI applications and agents to connect to and work with your data sources and tools
- Introduced by Anthropic in November 2024



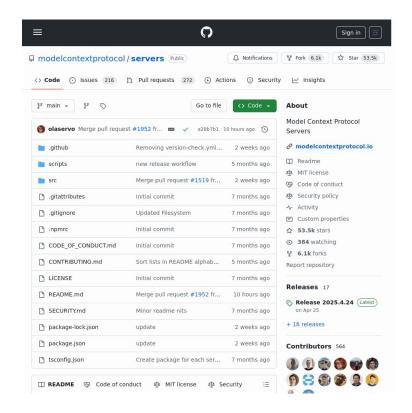
- MCP Model Context Protocol
- Open standard that describes a way for AI applications and agents to connect to and work with your data sources and tools
- Introduced by Anthropic in November 2024



- MCP Model Context Protocol
- Open standard that describes a way for AI applications and agents to connect to and work with your data sources and tools
- Introduced by Anthropic in November 2024

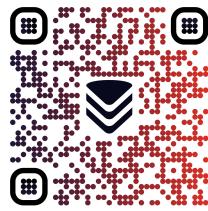


Since MCP was introduced





https://github.com/modelcontextprotocol/servers/



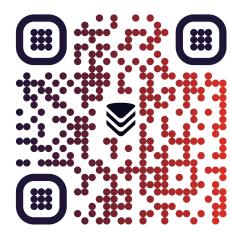
https://github.com/punkpeye/awesome-mcp-servers





Model Context Protocol

"Think of MCP like a USB-C port for Al applications. Just as USB-C provides a standardized way to connect your devices to various peripherals and accessories, MCP provides a standardized way to connect AI models to different data sources and tools."



https://modelcontextprotocol.io



Some software with chat interface

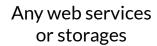


Some software with chat interface

Any web services or storages



Some software with chat interface



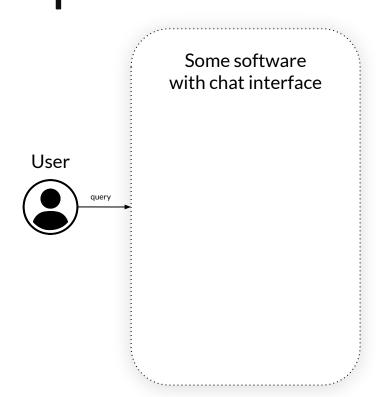






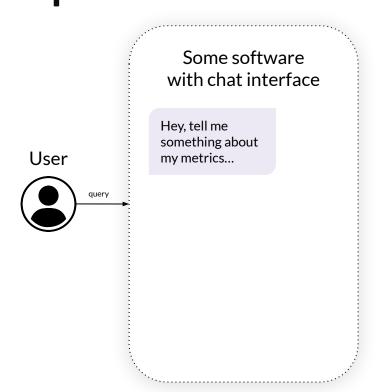


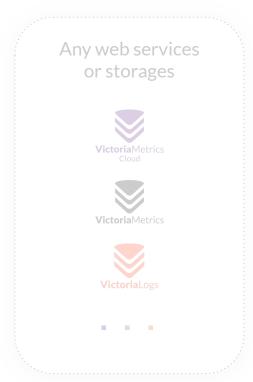




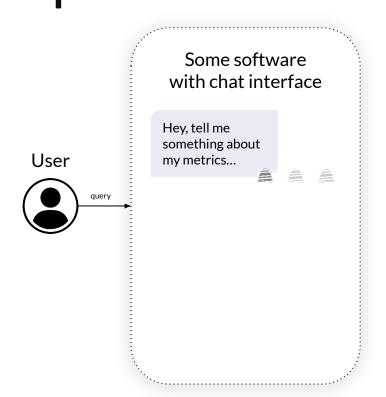


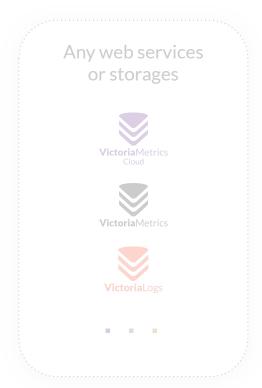




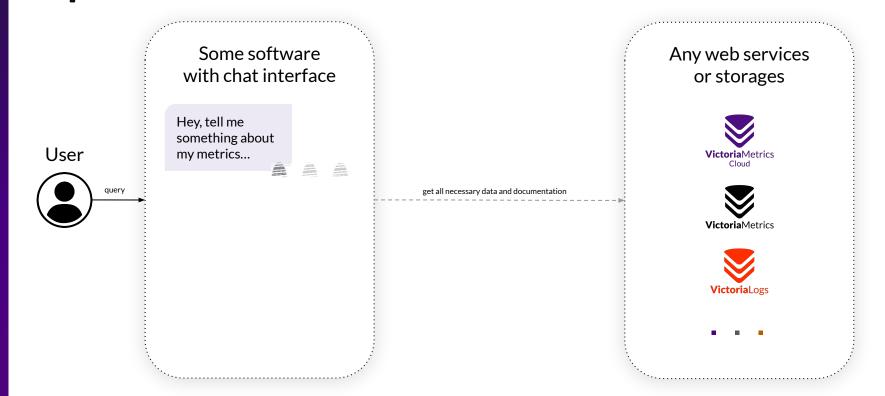




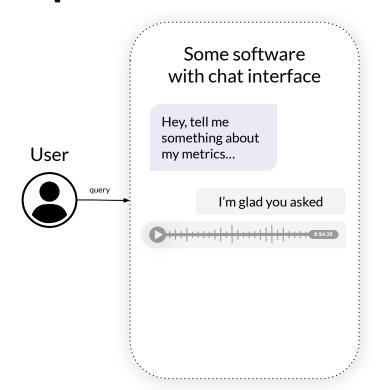


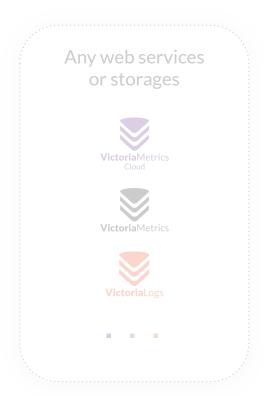




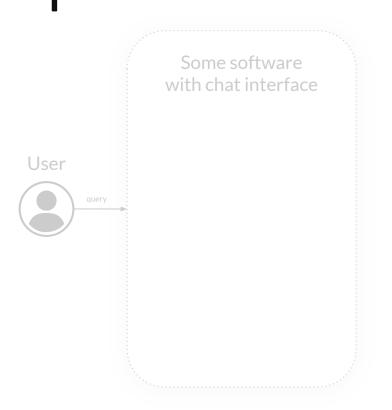








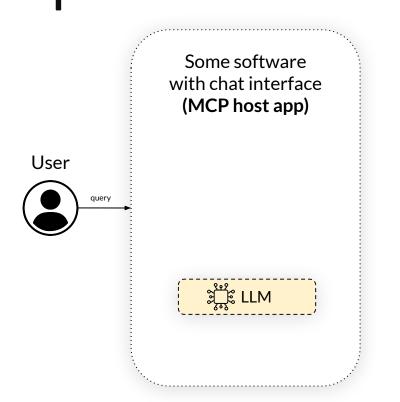
















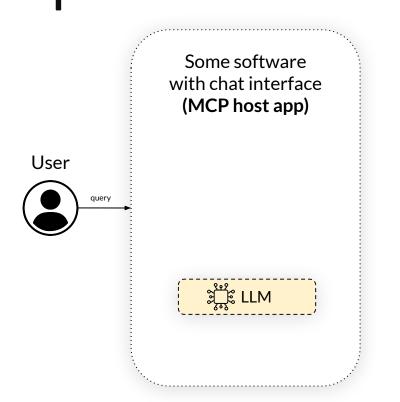


What is **LLM?**

$LLM \rightarrow Large Language Model$

A type of artificial intelligence model trained on vast amounts of text data to understand and generate human-like language. These models can perform various language tasks such as writing, translation, summarization, and question-answering by predicting the most likely next words or phrases based on patterns learned during training.

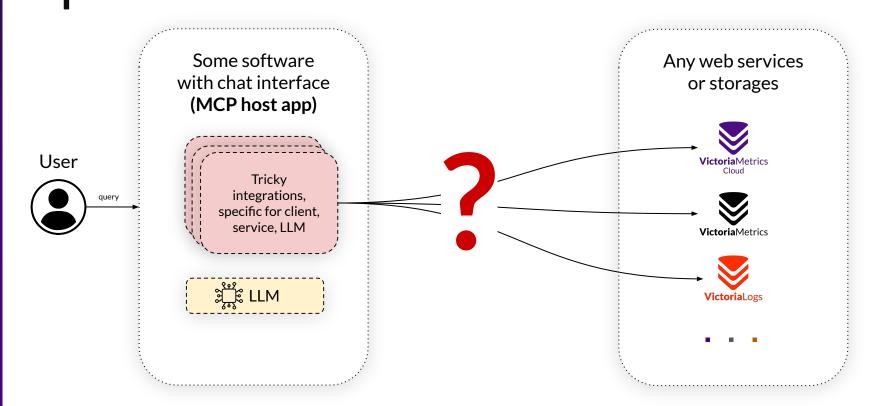




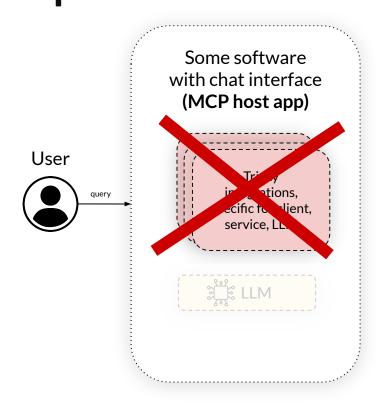








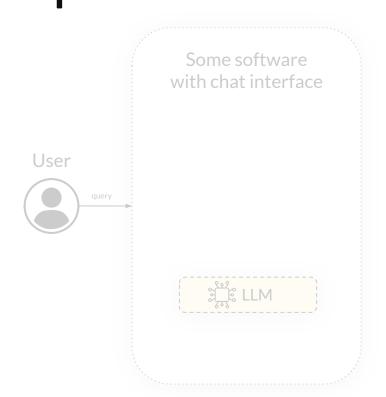


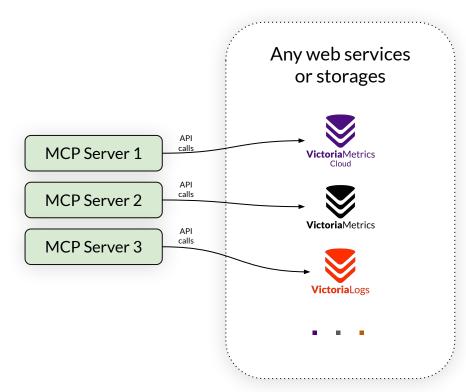


M x N problem

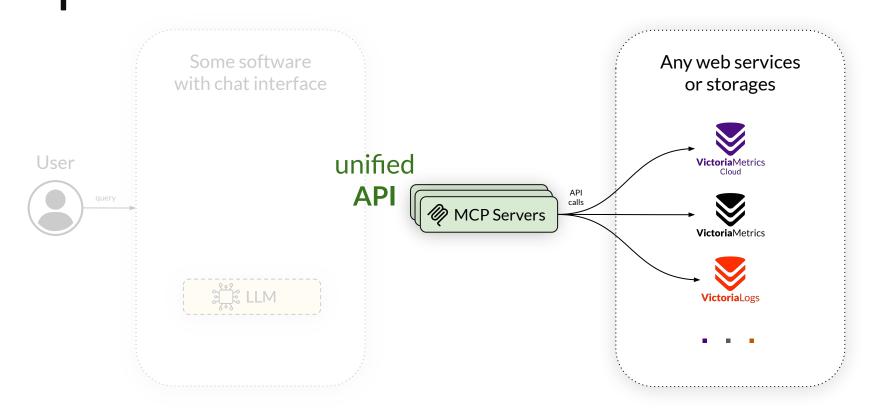




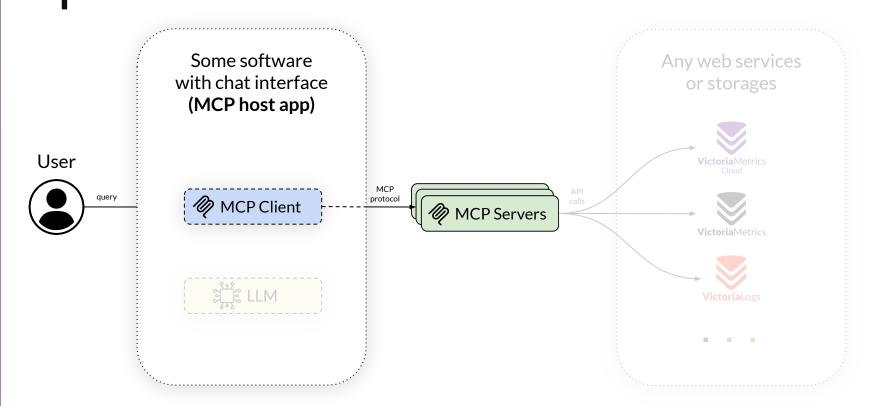




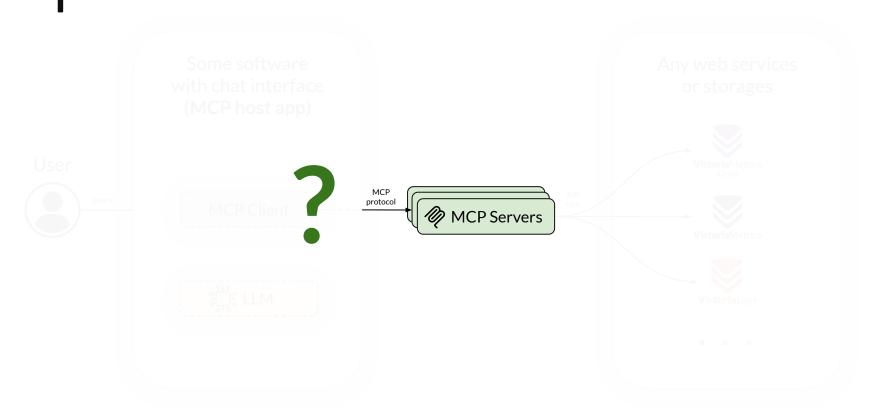




















MCP Servers declare and provide entities:



Tools - some executable functionality functionality for clients



Resources - some data to provide context for LLM







MCP Servers declare and provide entities:



Tools - some executable functionality for clients



Resources - some data to provide context for LLM







MCP Servers declare and provide entities:



Tools - some executable functionality for clients



Resources - some data to provide context for LLM







MCP Servers declare and provide entities:



Tools - some executable functionality for clients



Resources - some data to provide context for LLM







MCP Servers declare and provide entities:



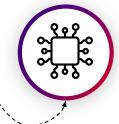
Tools - some executable functionality for clients



Resources - some data to provide context for LLM



Prompts - reusable prompt templates



Common recipes or scenarios that works with this MCP





MCP Servers declare and provide entities:



Tools - some executable functionality for clients



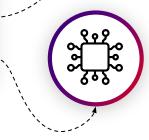
Resources - some data to provide context for LLM



Prompts - reusable prompt templates



Documentation, specifications, etc...



Common recipes or scenarios that works with this MCP







MCP Servers declare and provide entities:



Tools - some executable functionality for clients



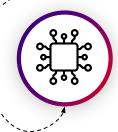
Resources - some data to provide context for LLM



Prompts - reusable prompt templates



Documentation, specifications, etc...



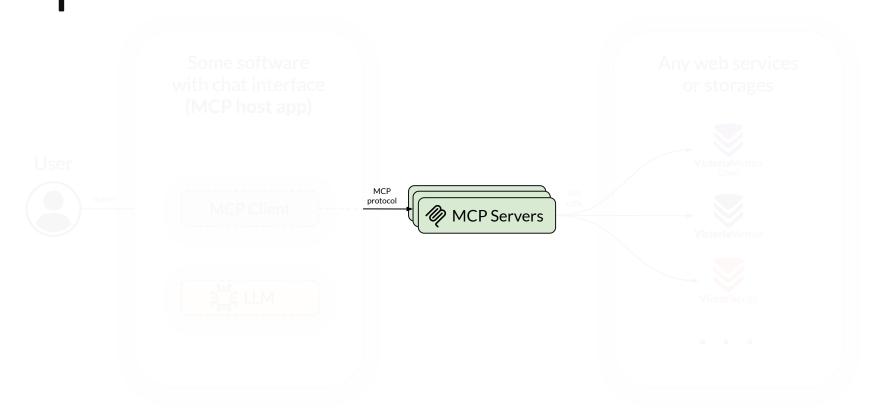
Common recipes or scenarios that works with this MCP



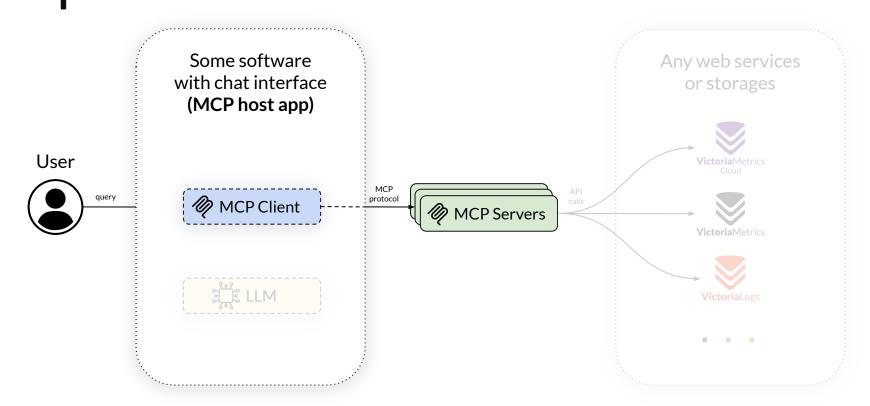
You can create **OpenAPI** based **MCP** Servers

\otimes	https://jedisct1.github.io/openapi-mcp/
\otimes	https://github.com/higress-group/openapi-to-mcpserver
\otimes	https://github.com/harsha-iiiv/openapi-mcp-generator
\otimes	https://github.com/janwilmake/openapi-mcp-server
\otimes	https://github.com/ckanthony/openapi-mcp
\otimes	https://mcp-link.vercel.app/
\otimes	https://github.com/hannesj/mcp-openapi-schema
\otimes	https://github.com/baryhuang/mcp-server-any-openapi

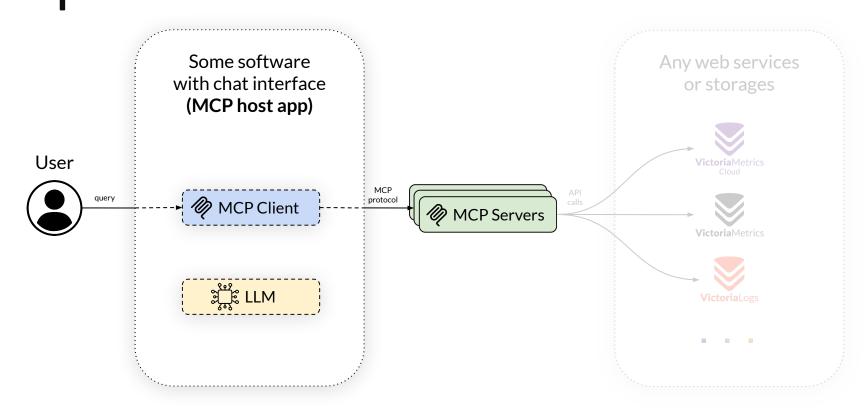




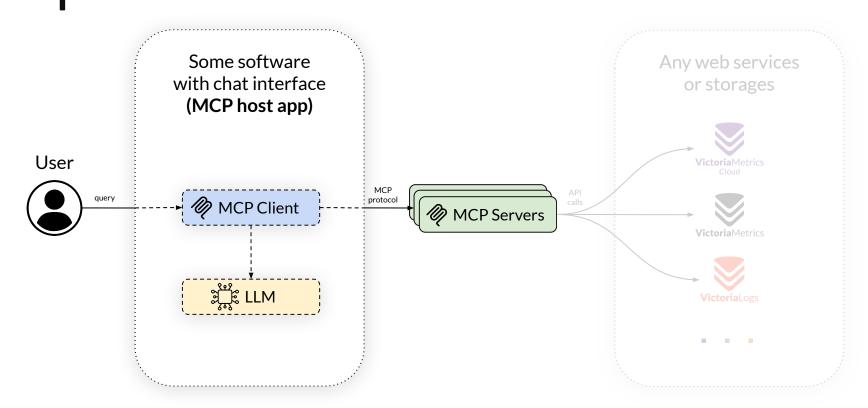




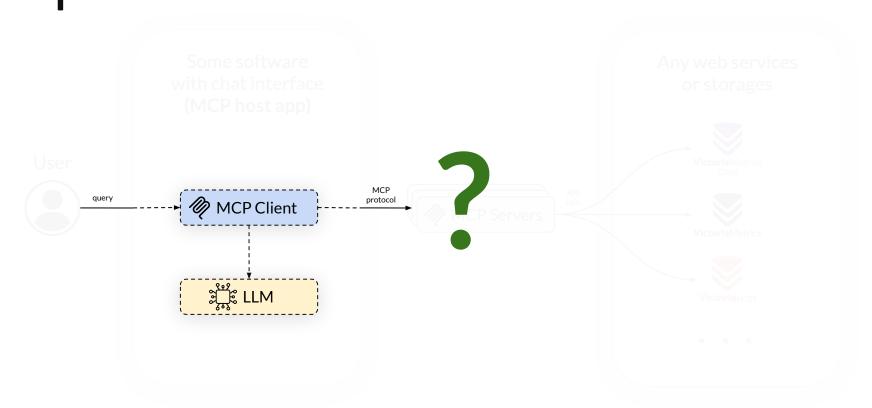




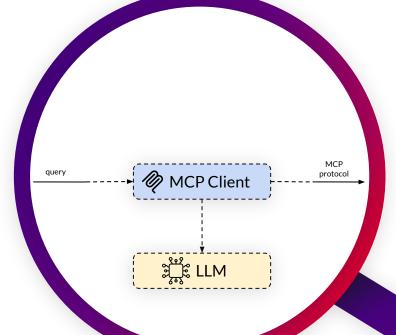




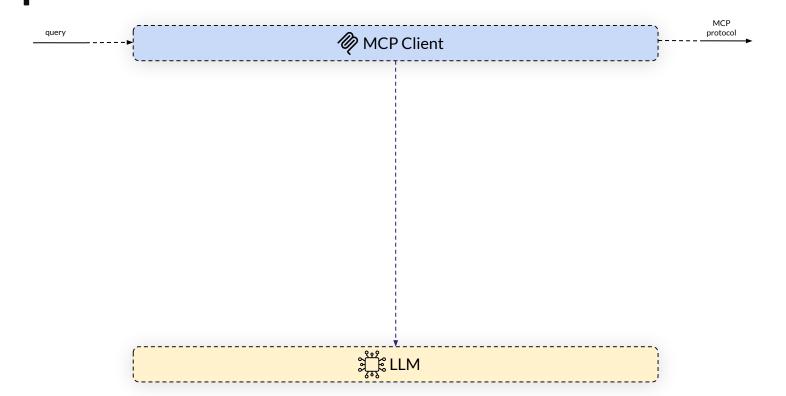




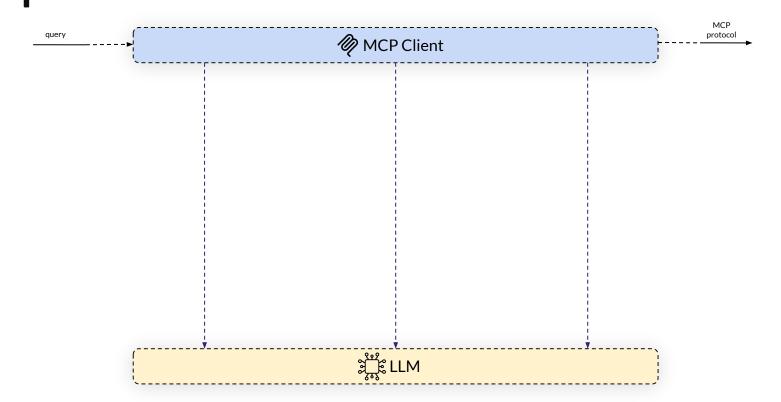




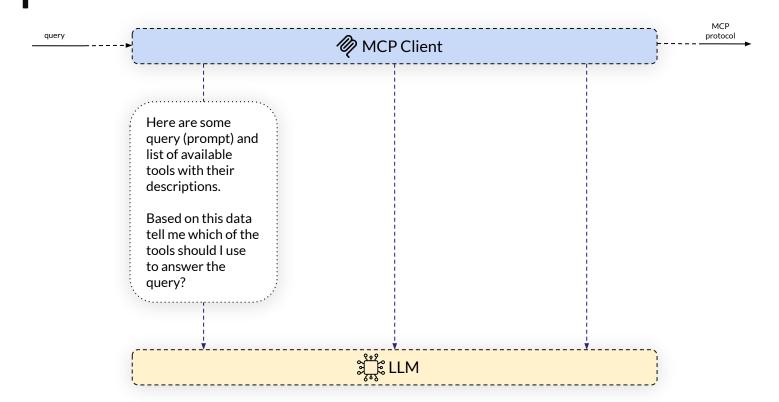




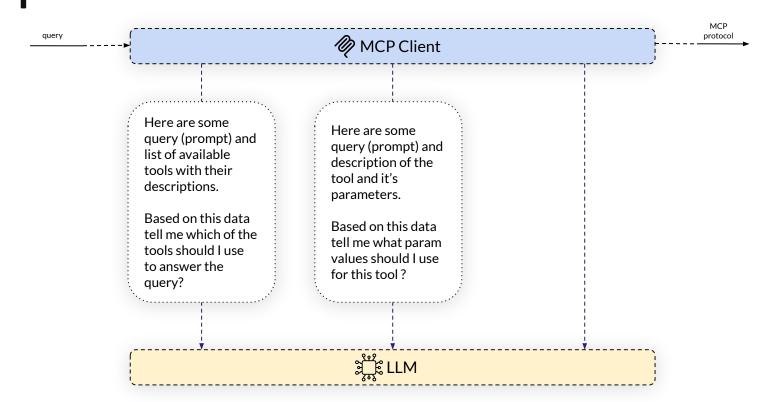




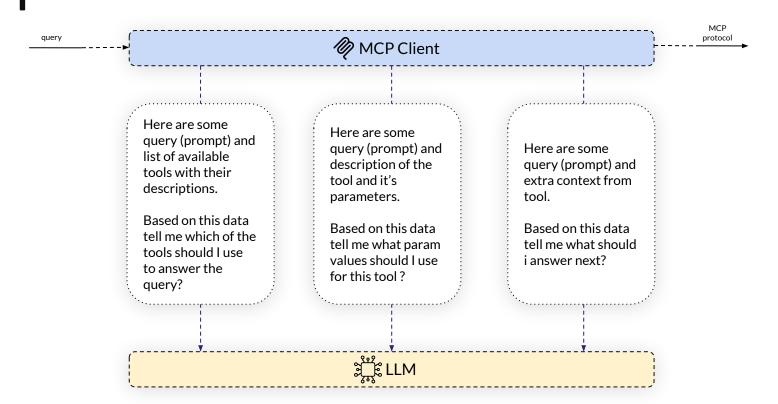




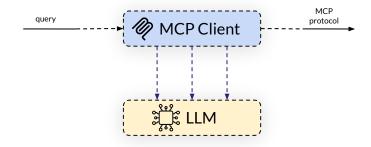




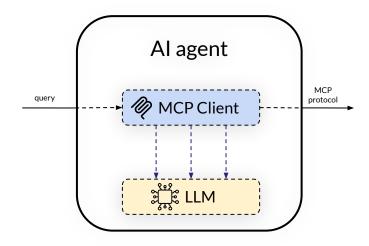




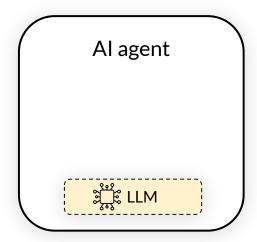




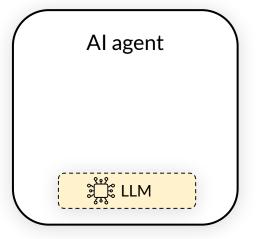


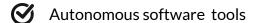


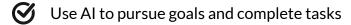






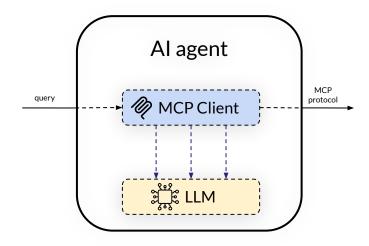




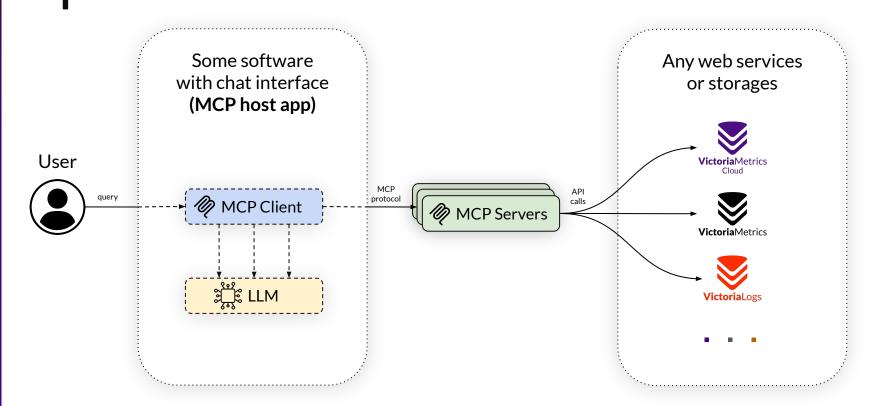


- Show reasoning, planning and have a level of autonomy to make decisions
- can work with other agents to coordinate and perform more complex workflows.

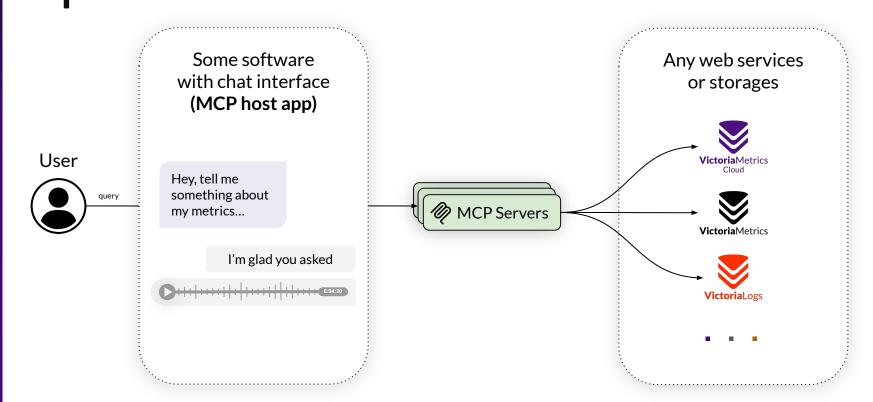










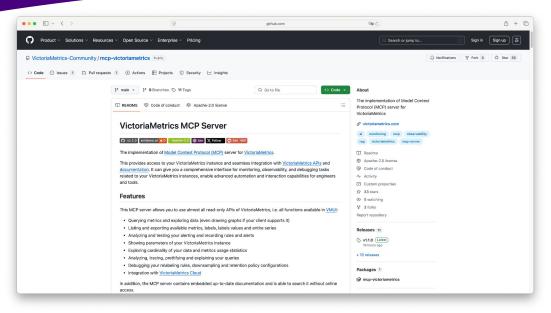




VictoriaMetrics MCP Server



mcp-victoriametrics





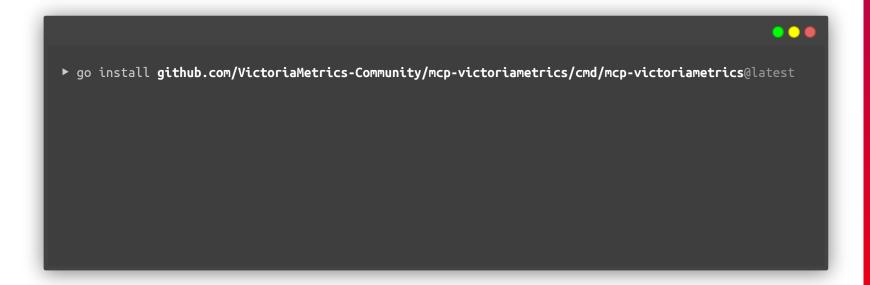


How to **install** with **docker**





How to **install** with **go**



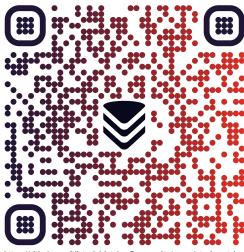


How to **install** from **source code**

```
    git clone https://github.com/VictoriaMetrics-Community/mcp-victoriametrics.git
    cd mcp-victoriametrics
    make build # if you want to build with local tooling
    docker build -t mcp-victoriametrics . # if you want to build with docker
```



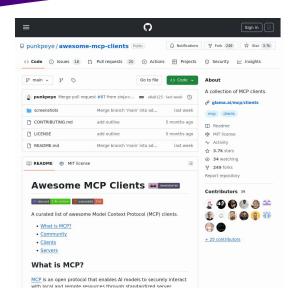
How to **install** using **binaries** or other methods

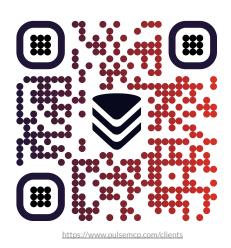


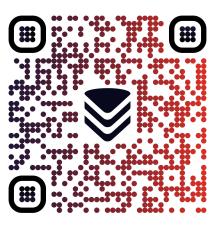




Awesome MCP clients list







https://github.com/punkpeve/awesome-mcp-clients



Which **clients** are supported?















Claude

VS Code

Cursor

JetBrains IDEs

Windsurf

LibreChat













Continue

Cline

Goose

Zed

RooCode

Warp







stdio mode



Name of mcp-server: it's just human readable name you'll see in UI



- Name of mcp-server: it's just human readable name you'll see in UI
- **Command to run**: can be binary file of server or tool for running it i.e. docker)



- Name of mcp-server: it's just human readable name you'll see in Ul
- Command to run: can be binary file of server or tool for running it i.e. docker)
- Arguments (args): command line arguments of binary file



- Name of mcp-server: it's just human readable name you'll see in UI
- Command to run: can be binary file of server or tool for running it i.e. docker)
- Arguments (args): command line arguments of binary file
- **Environment variables (env)**: parameters of MCP



- Name of mcp-server: it's just human readable name you'll see in UI
- Command to run: can be binary file of server or tool for running it i.e. docker)
- Arguments (args): command line arguments of binary file
- **Environment variables (env)**: parameters of MCP



```
Claude → claude_desktop_config.json
         "mcpServers": {
            "VictoriaMetrics": {
               "command": "/path/to/mcp-victoriametrics",
                  "VM_INSTANCE_ENTRYPOINT": "<YOUR_VM_INSTANCE>",
                  "VM_INSTANCE_TYPE": "<YOUR_VM_INSTANCE_TYPE>"
```



```
Claude → claude_desktop_config.json
            "VictoriaMetrics": { ← the name of MCP Server
```



```
Claude → claude_desktop_config.json
               "command": "/path/to/mcp-victoriametrics", 

the command to run
```



```
Claude → claude_desktop_config.json
               "env": { ← parameters in form of environment variables
                  "VM_INSTANCE_ENTRYPOINT": "<YOUR_VM_INSTANCE>",
                  "VM_INSTANCE_TYPE": "<YOUR_VM_INSTANCE_TYPE>"
```



```
Visual Studio Code → settings.json → "mcp"
        "servers": {
          "VictoriaMetrics": {
            "type": "stdio",
            "command": "/path/to/mcp-victoriametrics",
            "env": {
              "VM_INSTANCE_ENTRYPOINT": "<YOUR_VM_INSTANCE>",
              "VM_INSTANCE_TYPE": "<YOUR_VM_INSTANCE_TYPE>"
```



```
Cursor → mcp.json
         "mcpServers": {
            "VictoriaMetrics": {
               "command": "/path/to/mcp-victoriametrics",
                  "VM_INSTANCE_ENTRYPOINT": "<YOUR_VM_INSTANCE>",
                  "VM_INSTANCE_TYPE": "<YOUR_VM_INSTANCE_TYPE>"
```



```
Zed → settings.json
        "context_servers": {
              "path": "/path/to/mcp-victoriametrics",
                "VM_INSTANCE_ENTRYPOINT": "<YOUR_VM_INSTANCE>",
                "VM_INSTANCE_TYPE": "<YOUR_VM_INSTANCE_TYPE>"
```



```
JetBrains IDE → settings
         "mcpServers": {
            "VictoriaMetrics": {
               "command": "/path/to/mcp-victoriametrics",
                  "VM_INSTANCE_ENTRYPOINT": "<YOUR_VM_INSTANCE>",
                  "VM_INSTANCE_TYPE": "<YOUR_VM_INSTANCE_TYPE>"
```



stdio mode

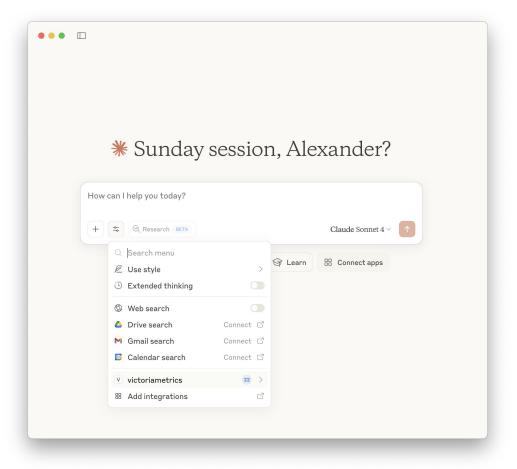
- Name of mcp-server: it's just human readable name you'll see in UI
- Command to run: can be binary file of server or tool for running it i.e. docker)
- Arguments (args): command line arguments of binary file
- **Environment variables (env)**: parameters of MCP



```
Claude → claude_desktop_config.json
             "command": "docker",
               "args":
                 "run", "-i", "--rm",
                 "-e", "VM INSTANCE ENTRYPOINT",
                 "-e", "VM INSTANCE TYPE",
                 "ghcr.io/victoriametrics-community/mcp-victoriametrics:latest",
               ],
               "VM_INSTANCE_ENTRYPOINT": "https://play.victoriametrics.com",
               "VM INSTANCE TYPE": "cluster"
```

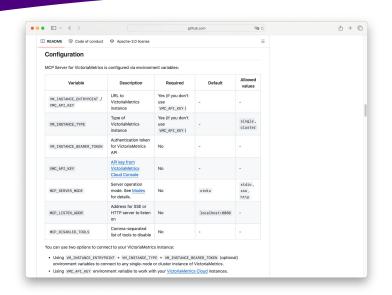


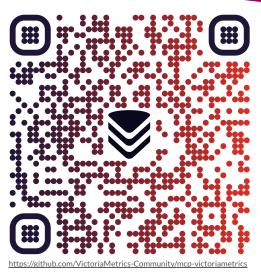
Tools in your **client**



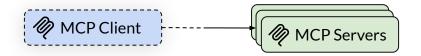


mcp-victoriametrics configuration

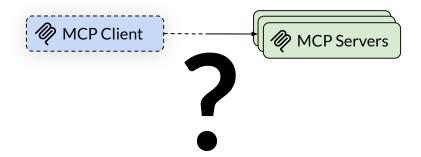














The stdio transport enables communication through standard input and output streams. This is particularly useful for local integrations and command-line tools



The most common transport



Simple, secure, no network



Only local usage, no scalability



The fastest (100µs)



Perfect for CLIs, desktop apps and local integrations



The stdio transport enables communication through standard input and output streams. This is particularly useful for local integrations and command-line tools



The most common transport



Simple, secure, no network



Only local usage, no scalability



The fastest (100µs)



Perfect for CLIs, desktop apps and local integrations

streamable

HTTP

The Streamable HTTP transport uses HTTP POST requests for client-to-server communication and optional Server-Sent Events (SSE) streams for server-to-client communication.



Good for web services, APIs



Multi-tenancy, concurrency



Scalability



Still a small adoption



Slower that others: ~1-10ms (network + HTTP overhead)



The stdio transport enables communication through standard input and output streams. This is particularly useful for local integrations and command-line tools



The most common transport



Simple, secure, no network



Only local usage, no scalability



The fastest (100µs)



Perfect for CLIs, desktop apps and local integrations

streamable

HTTP

The Streamable HTTP transport uses HTTP POST requests for client-to-server communication and optional Server-Sent Events (SSE) streams for server-to-client communication.



Good for web services, APIs



Multi-tenancy, concurrency



Scalability



Still a small adoption



Slower that others: ~1-10ms (network + HTTP overhead)

SSE

The legacy SSE transport enabled server-to-client streaming with HTTP POST requests for client-to-server communication.



Deprecated



 $Replaced\ with\ Streamable\ HTTP$



But still support as only way to connect some of the clients



The stdio transport enables communication through standard input and output streams. This is particularly useful for local integrations and command-line tools



Simple, secure, no network



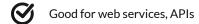
The fastest (100µs)

Perfect for CLIs, desktop apps and local integrations

streamable

HTTP

The Streamable HTTP transport uses HTTP POST requests for client-to-server communication and optional Server-Sent Events (SSE) streams for server-to-client communication.



Multi-tenancy, concurrency

Scalability

Still a small adoption

Slower that others: ~1-10ms (network + HTTP overhead)

SSE

The legacy SSE transport enabled server-to-client streaming with HTTP POST requests for client-to-server communication.

O Deprecated

Replaced with Streamable HTTP

But still support as only way to connect some of the clients

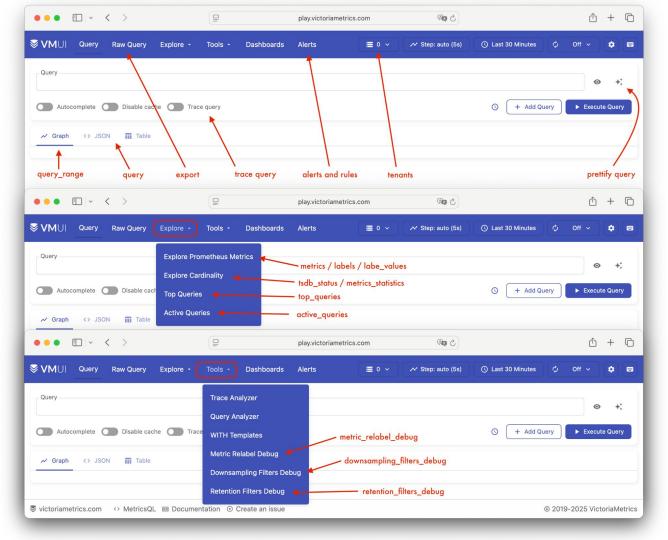


Features overview



\odot	tenants (for cluster)
8	metrics / labels / label_values / series
\otimes	rules / alerts / test_rules
\odot	active_queries / top_queries
\otimes	query / query_range
\otimes	tsdb_status (cardinality explorer) / metric_statistic
\otimes	metrics_relabel_debug / downsampling_filters_debug / retention_filters_debug
8	explain_query / prettify_query
8	documentation





tenants (for cluster)
metrics / labels / label_values / series
rules / alerts / test_rules
active_queries / top_queries
query / query_range
tsdb_status (cardinality explorer) / metric_statistic
metrics_relabel_debug / downsampling_filters_debug / retention_filters_debug
explain_query / prettify_query
documentation



\odot	tenants (for cluster)
8	metrics / labels / label_values / series
\otimes	rules / alerts / test_rules
\odot	active_queries / top_queries
\otimes	query / query_range
\otimes	tsdb_status (cardinality explorer) / metric_statistic
\otimes	metrics_relabel_debug / downsampling_filters_debug / retention_filters_debug
8	explain_query / prettify_query
8	documentation

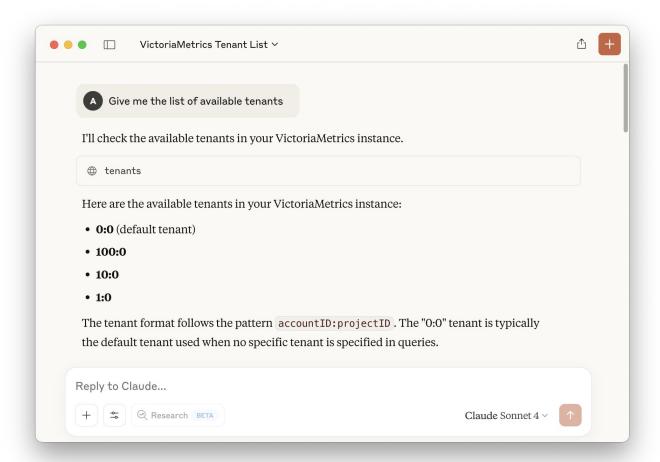


Usage scenarios and examples

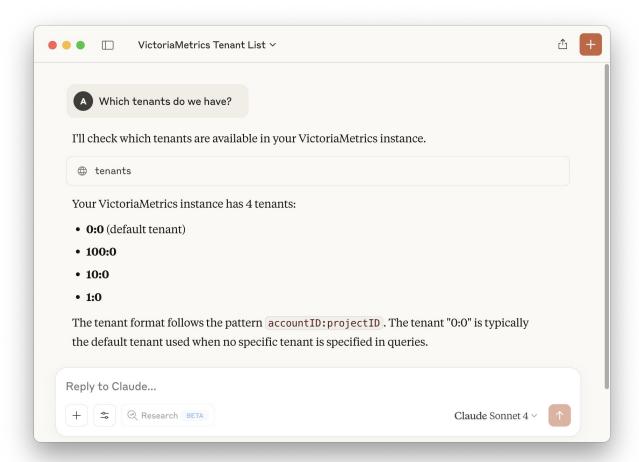


\otimes	tenants (for cluster)
	metrics / labels / label_values / series
	rules / alerts / test_rules
	active_queries / top_queries
	query / query_range
	tsdb_status (cardinality explorer) / metric_statistic
	metrics_relabel_debug / downsampling_filters_debug / retention_filters_debug
	explain_query / prettify_query
	documentation



















 $I'II\ check\ which\ tenants\ are\ available\ in\ your\ Victoria Metrics\ instance.$



tenants

AAA



I'll check which tenants are available in your VictoriaMetrics instance.



Your VictoriaMetrics instance has 4 tenants:

- 0:0 (default tenant)
- 100:0
- 10:0
- 1:0

The tenant format follows the pattern accountID:projectID. The tenant "0:0" is typically the default tenant used when no specific tenant is specified in queries.



I'll check which tenants are available in your VictoriaMetrics instance.



tenants

Your Victoria Metrics instance has 4 tenants:

- 0:0 (default tenant)
- 100:0
- 10:0
- 1:0

The tenant format follows the pattern account ID: project ID. The tenant "0:0" is typically the default tenant used when no specific tenant is specified in queries.

from now on, let's work with default tenant









I'll check which tenants are available in your VictoriaMetrics instance.



tenants

Your Victoria Metrics instance has 4 tenants:

- 0:0 (default tenant)
- 100:0
- 10:0
- 1:0

The tenant format follows the pattern account ID: project ID. The tenant "0:0" is typically the default tenant used when no specific tenant is specified in queries.

from now on, let's work with default tenant

Understood! I'll use the default tenant "0" (which corresponds to "0:0") for all subsequent queries and operations. This is actually the default behavior, so all Victoria Metrics API calls will automatically use the default tenant unless you specify otherwise.



\otimes	tenants (for cluster)
0	metrics / labels / label_values / series
	rules / alerts / test_rules
	active_queries / top_queries
	query / query_range
	tsdb_status (cardinality explorer) / metric_statistic
	metrics_relabel_debug / downsampling_filters_debug / retention_filters_debug
	explain_query / prettify_query
	documentation



	tenants (for cluster)
\otimes	metrics / labels / label_values / series
0	rules / alerts / test_rules
	active_queries / top_queries
	query / query_range
	tsdb_status (cardinality explorer) / metric_statistic
	metrics_relabel_debug / downsampling_filters_debug / retention_filters_debug
	explain_query / prettify_query
	documentation



from now on, let's work with default tenant

Understood! I'll use the default tenant "0" (which corresponds to "0:0") for all subsequent queries and operations. This is actually the default behavior, so all VictoriaMetrics API calls will automatically use the default tenant unless you specify otherwise.



Understood! I'll use the default tenant "0" (which corresponds to "0:0") for all subsequent queries and operations. This is actually the default behavior, so all VictoriaMetrics API calls will automatically use the default tenant unless you specify otherwise.

What metrics do we have?







from now on, let's work with default tenant

Understood! I'll use the default tenant "0" (which corresponds to "0:0") for all subsequent queries and operations. This is actually the default behavior, so all VictoriaMetrics API calls will automatically use the default tenant unless you specify otherwise.

What metrics do we have?

I'll get the list of available metrics from your VictoriaMetrics instance.

AAA





I'll get the list of available metrics from your VictoriaMetrics instance.



Your VictoriaMetrics instance has a comprehensive set of metrics across various categories:

- Infrastructure & System Metrics:
 - Kubernetes components: Extensive metrics from API server (apiserver_*), kubelet (kubelet_*), controller runtime (controller runtime *)
 - Node metrics: CPU, memory, disk, network statistics (node_*)
 - Container metrics: Resource usage, filesystem, network (container_*)
- Monitoring & Observability:
 - Prometheus ecosystem: Alertmanager (alertmanager_*), kube-state-metrics (kube_*)
 - VictoriaMetrics components: VM storage (vm_*), VMAgent (vmagent_*), VMAlert (vmalert_*), VMAuth (vmauth_*)
 - VictoriaLogs: (vl_*) for log management
 - Vector log collector: (vector_*)
 - ArgoCD: Application and GitOps metrics (argocd_*)
- Application & Custom Metrics:
 - Go runtime: Garbage collection, memory, goroutines (go_*)
 - HTTP & networking: Request durations, connection stats (grpc_*, rest_client_*)
 - Custom applications: Various business metrics like all_visits_count, unique_visits_count, github_*, docker_*

You have 808 unique metrics covering infrastructure monitoring, application performance, and business analytics. This suggests a robust monitoring setup with full observability across your stack.

Would you like me to analyze any specific metric category or help you explore particular aspects of your monitoring data?



Your VictoriaMetrics instance has a comprehensive set of metrics across various categories:

- Infrastructure & System Metrics:
 - Kubernetes components: Extensive metrics from API server (apiserver *), kubelet (kubelet *), controller runtime (controller runtime *)
 - Node metrics: CPU, memory, disk, network statistics (node *)
 - Container metrics: Resource usage, filesystem, network (container *)
- Monitoring & Observability:
 - Prometheus ecosystem: Alertmanager (alertmanager *), kube-state-metrics (kube *)
 - VictoriaMetrics components: VM storage (vm *), VMAgent (vmagent *), VMAlert (vmalert *), VMAuth (vmauth *)
 - VictoriaLogs: (vl *) for log management
 - Vector log collector: (vector *)
 - ArgoCD: Application and GitOps metrics (argocd *)
- Application & Custom Metrics:
 - Go runtime: Garbage collection, memory, goroutines (go *)
 - HTTP & networking: Request durations, connection stats (qrpc *, rest client *)
 - Custom applications: Various business metrics like all visits count, unique visits count, github *, docker *

You have 808 unique metrics covering infrastructure monitoring, application performance, and business analytics. This suggests a robust monitoring setup with full observability across your stack.

Would you like me to analyze any specific metric category or help you explore particular aspects of your monitoring data?

Show me argcd metrics







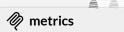
- Container metrics: Resource usage, filesystem, network (container_*)
- Monitoring & Observability:
 - Prometheus ecosystem: Alertmanager (alertmanager *), kube-state-metrics (kube *)
 - VictoriaMetrics components: VM storage (vm_*), VMAgent (vmagent_*), VMAlert (vmalert_*), VMAuth (vmauth *)
 - VictoriaLogs: (vl_*) for log management
 - Vector log collector: (vector_*)
 - ArgoCD: Application and GitOps metrics (argocd_*)
- Application & Custom Metrics:
 - Go runtime: Garbage collection, memory, goroutines (qo *)
 - HTTP & networking: Request durations, connection stats (grpc *, rest client *)
 - Custom applications: Various business metrics like all_visits_count, unique_visits_count, github_*, docker_*

You have 808 unique metrics covering infrastructure monitoring, application performance, and business analytics. This suggests a robust monitoring setup with full observability across your stack.

Would you like me to analyze any specific metric category or help you explore particular aspects of your monitoring data?

Show me argcd metrics

I'll get the ArgoCD metrics from your Victoria Metrics instance.





I'll get the ArgoCD metrics from your VictoriaMetrics instance.



Perfect! Here are all the ArgoCD metrics available in your VictoriaMetrics instance:

ArgoCD Metrics Categories

Application Management

- argocd_app_info Information about ArgoCD applications
- argocd app orphaned resources count Count of orphaned resources per application
- argocd_app_reconcile_bucket/count/sum Application reconciliation performance metrics

Cluster Management

- argocd_cluster_api_resource_objects Number of API resource objects in clusters
- argood cluster api resources Count of API resources in clusters
- argocd_cluster_cache_age_seconds Age of cluster cache data
- argocd_cluster_connection_status Connection status to managed clusters
- argocd_cluster_events_total Total cluster events processed
- argocd_cluster_info Information about managed clusters

Git Operations

- argocd_git_request_duration_seconds_* Git request latency metrics (histogram)
- argocd_git_request_total Total Git requests made

• • •



Kubectl Operations

- argood kubectl rate limiter duration seconds *-Rate limiting duration metrics
- argord kubectl request duration seconds *-Kubectl request performance
- argocd_kubectl_request_size_bytes_* Size of kubectl requests
- argocd_kubectl_requests_total Total kubectl requests
- argocd_kubectl_response_size_bytes_* Size of kubectl responses

Redis Operations

- argood redis request duration * Redis request latency metrics
- argocd_redis_request_total Total Redis requests

Repository Management

• argocd_repo_pending_request_total - Pending repository requests

Resource Processing

- argocd_resource_events_processed_in_batch Git request latency metrics (histogram)
- argocd_resource_events_processing_* Resource event processing performance

General Info

• argocd_info - General ArgoCD instance information

Would you like me to query specific ArgoCD metrics to see their current values or analyze particular aspects of your ArgoCD deployment performance?

What labels does the argocd_app_info have?





• argord kubecti response size bytes *-Size of kubecti responses

Redis Operations

- argood redis request duration * Redis request latency metrics
- argood redis request total Total Redis requests

Repository Management

• argood repo pending request total - Pending repository requests

Resource Processing

- argocd_resource_events_processed_in_batch Git request latency metrics (histogram)
- argord resource events processing * Resource event processing performance

General Info

• argocd info - General ArgoCD instance information

Would you like me to query specific ArgoCD metrics to see their current values or analyze particular aspects of your ArgoCD deployment performance?

What labels does the argocd_app_info have?

I'll guery the series for argood app info to see what labels are available.



O series

A A A



I'll query the series for argocd_app_info to see what labels are available.



Based on the series data, the argocd_app_info metric has the following labels:

ArgoCD-Specific Labels

- name The name of the ArgoCD application (e.g., "grafana", "monitoring", "argocd")
- autosync enabled Whether auto-sync is enabled ("true"/"false")
- dest_namespace The destination namespace where the app is deployed
- health_status Application health status (e.g., "Healthy")
- project ArgoCD project name (e.g., "default")
- repo Git repository URL or Helm chart repository
- sync status Synchronization status (e.g., "Synced")

Standard Kubernetes/Prometheus Labels

- namespace Kubernetes namespace where ArgoCD is running ("argocd")
- cluster Cluster identifier ("sandbox")
- service Kubernetes service name ("argood-metrics")
- job Prometheus job name ("argocd-metrics")
- instance Instance endpoint ("10.71.10.3:8082")
- pod Pod name ("argocd-application-controller-0")
- container Container name ("argocd-application-controller")
- endpoint Metrics endpoint ("metrics")





Sync Status - Synchronization status (e.g., Synced)

Standard Kubernetes/Prometheus Labels

- namespace Kubernetes namespace where ArgoCD is running ("argocd")
- cluster Cluster identifier ("sandbox")
- service Kubernetes service name ("argood-metrics")
- job Prometheus job name ("argood-metrics")
- instance Instance endpoint ("10.71.10.3:8082")
- pod Pod name ("argocd-application-controller-0")
- container Container name ("argood-application-controller")
- endpoint Metrics endpoint ("metrics")
- exported namespace Original namespace label ("argood")
- prometheus Prometheus instance ("monitoring/monitoring-vmagent")

Your Current Applications

I can see you have 24 ArgoCD applications deployed, all with:

- Auto-sync enabled
- Healthy status
- Synced state

Some notable applications include monitoring tools (Grafana, Prometheus), Victoria Metrics components, and various operational tools.

Would you like me to query specific applications or analyze the health/sync status across your ArgoCD deployment?

What values i have for label name of this metric?









- pod Pod name ("argocd-application-controller-0")
- container Container name ("argood-application-controller")
- endpoint Metrics endpoint ("metrics")
- exported namespace Original namespace label ("argood")
- prometheus Prometheus instance ("monitoring/monitoring-vmagent")

Your Current Applications

I can see you have 24 ArgoCD applications deployed, all with:

- Auto-sync enabled
- Healthy status
- Synced state

Some notable applications include monitoring tools (Grafana, Prometheus), Victoria Metrics components, and various operational tools.

Would you like me to query specific applications or analyze the health/sync status across your ArgoCD deployment?

What values i have for label name of this metric?

I'll get the distinct values for the name label of the argood app info metric.

AAA

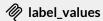


label_values



What values i have for label name of this metric?

I'll get the distinct values for the name label of the argocd_app_info metric.



Here are all the values for the name label of the argocd_app_info metric, representing your deployed ArgoCD applications:

. . .

Your ArgoCD Applications (24 total)

Monitoring & Observability

- analytic-exporter
- grafana
- grafana-dashboards
- kube-state-metrics
- monitoring
- node-exporter
- play-grafana
- play-grafana-dashboards

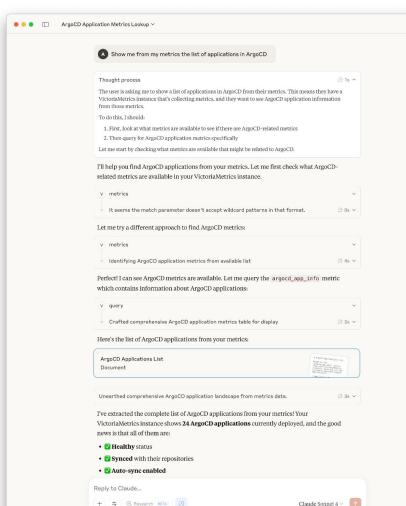
VictoriaMetrics Ecosystem

- victoria-metrics-operator
- vmlicense
- vmlogs-single
- vmlogs-slack2logs
- benchmark-vm-cluster

MCP (Model Context Protocol) Services

mcp-victorialogs

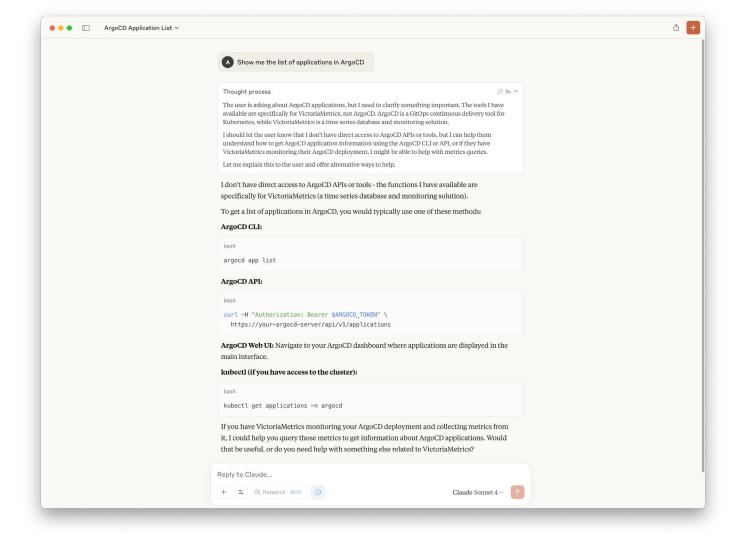




ArgoCD Applications List

Based on the argocd_app_info metric from your VictoriaMetrics instance, here are all the currently deployed ArgoCD applications:

Application Name	Destination Namespace	Health Status	Sync Status	Auto- Sync	Project	Repository
analytic-	analytic-					
exporter	exporter	Healthy	Synced		default	https://github.com/vio
app-of-apps	argood	Healthy	Synced		default	https://github.com/vio
argocd	argood	Healthy	Synced		default	https://github.com/vio
benchmark- vm-cluster	benchmark	Healthy	Synced	W	default	https://github.com/vio
external-dns	external-dns	Healthy	Synced		default	https://kubernetes-sig
grafana	grafana	Healthy	Synced	2	default	https://grafana.githul
grafana- dashboards	grafana	Healthy	Synced	V	default	https://github.com/vio
kube-state- metrics	kube-state- metrics	Healthy	Synced		default	https://prometheus-c
load-balancer	-	Healthy	Synced	☑	default	https://github.com/vio
mcp- victorialogs	mcp- victorialogs	Healthy	Synced		default	https://github.com/vio
mcp- victoriametrics	mcp- victoriametrics	Healthy	Synced		default	https://github.com/vio
monitoring	monitoring	Healthy	Synced	V	default	https://github.com/vio
node-exporter	node-exporter	Healthy	Synced		default	https://prometheus-c
play-grafana	play-grafana	Healthy	Synced	☑	default	https://grafana.githul
play-grafana- dashboards	play-grafana	Healthy	Synced		default	https://github.com/vio
prometheus- benchmark-Its	prometheus- benchmark	Healthy	Synced	2	default	https://github.com/vio
prometheus- benchmark- stable	prometheus- benchmark	Healthy	Synced	2	default	https://github.com/vio
prometheus- benchmark- staging	prometheus- benchmark	Healthy	Synced	☑	default	https://github.com/vio benchmark
slack2logs	slack2logs	Healthy	Synced		default	https://github.com/vio
tg-spam	tg-spam	Healthy	Synced		default	https://github.com/vio
victoria- metrics- operator	vm-operator	Healthy	Synced	2	default	https://github.com/vio
vmlicense	vmlicense	Healthy	Synced	2	default	https://github.com/videnterprise/
vmlogs-single	vmlogs	Healthy	Synced		default	https://victoriametric
vmlogs-						



	tenants (for cluster)
\odot	metrics / labels / label_values / series
\odot	rules / alerts / test_rules
	active_queries / top_queries
	query / query_range
	tsdb_status (cardinality explorer) / metric_statistic
	metrics_relabel_debug / downsampling_filters_debug / retention_filters_debug
	explain_query / prettify_query
	documentation



	tenants (for cluster)
	metrics / labels / label_values / series
\odot	rules / alerts / test_rules
\odot	active_queries / top_queries
	query / query_range
	tsdb_status (cardinality explorer) / metric_statistic
	metrics_relabel_debug / downsampling_filters_debug / retention_filters_debug
	explain_query / prettify_query
	documentation











show me current alerts

I'll check the current alerts in your VictoriaMetrics instance.





I'll check the current alerts in your VictoriaMetrics instance.



Alert Summary

You have many firing alerts and several pending alerts currently active in your sandbox cluster.

Critical Alerts (Firing)

HostFilesystemDeviceError - Multiple instances

- Severity: Critical
- Status: Firing (since May 26, 2025)
- Issue: Permission denied errors on filesystem devices
- Affected instances: Multiple nodes (10.142.0.x:9100)
- Devices affected:
 - o tmpfs filesystems (Kubernetes projected volumes, secrets, config maps)
 - o /dev/sdb, /dev/sdc persistent disk mount points
 - Various Kubernetes volume mount paths

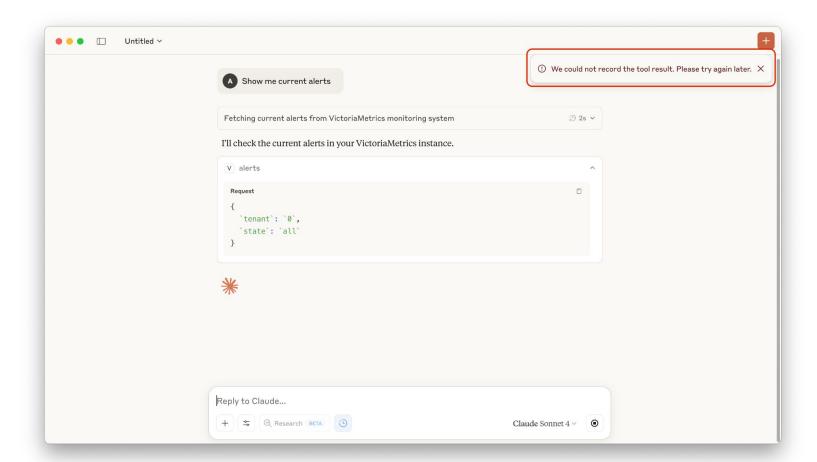
HostCpulsUnderutilized - Some instances

- Severity: Info
- Status: Firing
- Issue: CPU utilization < 20% for extended periods
- Recommendation: Consider reducing CPU allocation

Pending Alerts 🗾

HostCpulsUnderutilized - Multiple instances







Various Kubernetes volume mount paths

HostCpulsUnderutilized - Some instances

- Severity: Info
- Status: Firing
- Issue: CPU utilization < 20% for extended periods
- Recommendation: Consider reducing CPU allocation

Pending Alerts \(\big| \)

HostCpulsUnderutilized - Multiple instances

- Severity: Info • Status: Firing
- Issue: CPU utilization < 20% for 1 week
- Recommendation: Various nodes showing underutilization

Key Issues to Address

- 1. Permission Issues: The filesystem device errors suggest there are permission problems accessing mounted volumes on your Kubernetes nodes. This could impact application functionality.
- 2. Resource Optimization: Multiple CPU underutilization alerts suggest you may be able to optimize your cluster resource allocation.
- 3. Cluster Health: The large number of filesystem permission errors across multiple nodes indicates a systemic issue that should be investigated.

Most of these alerts are from your sandbox-gke environment in the sandbox cluster. The filesystem permission errors in particular warrant immediate attention as they could affect application performance and stability.

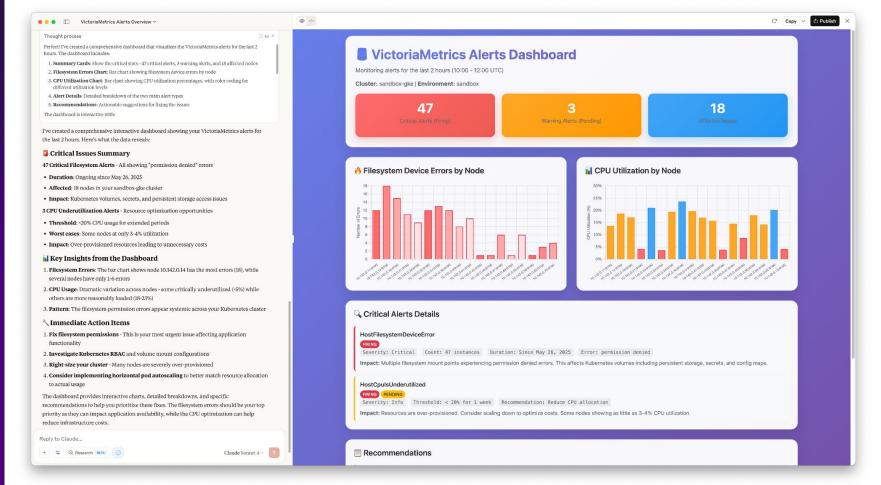
show me the graph for the queries for these alerts



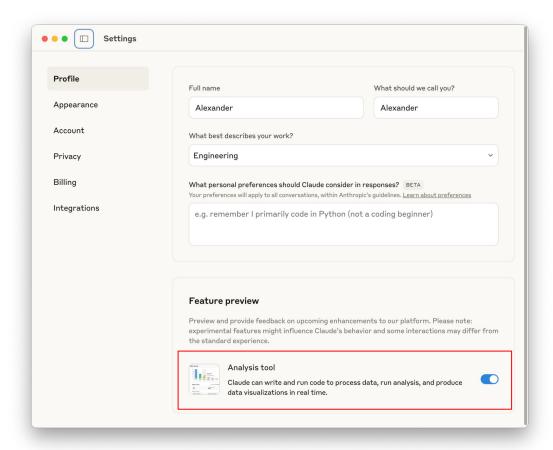






















what alerting rules i have?

 $I'II\ check\ what\ alerting\ rules\ are\ configured\ in\ your\ Victoria\ Metrics\ instance.$





I'll check what alerting rules are configured in your VictoriaMetrics instance.



Great! Now I can see your alerting rule groups. You have quite a comprehensive alerting setup with 16 different rule groups. Here's what I found:

Your Alerting Rule Groups

You have 16 alerting rule groups configured across several files:

Rule Groups by Category:

Infrastructure & Host Monitoring:

- node-exporter Host-level monitoring alerts
- container-events Container event monitoring
- kubernetes-apps Kubernetes application alerts
- kubernetes-storage Kubernetes storage alerts

VictoriaMetrics Monitoring:

- vm-health VictoriaMetrics health monitoring
- vm-license VictoriaMetrics license monitoring
- vmagent VMAgent monitoring alerts
- vmalert VMAlert monitoring alerts
- vmcluster VMCluster monitoring alerts

Application Monitoring:

- ArgoCD ArgoCD application alerts
- analitycs-exporter Analytics exporter monitoring



Simple, Reliable, Efficient Monitoring

	tenants (for cluster)
	metrics / labels / label_values / series
\otimes	rules / alerts / test_rules
0	active_queries / top_queries
	query / query_range
	tsdb_status (cardinality explorer) / metric_statistic
	metrics_relabel_debug / downsampling_filters_debug / retention_filters_debug
	explain_query / prettify_query
	documentation

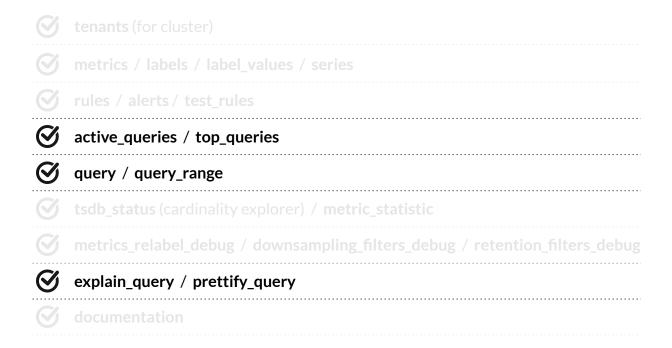


	tenants (for cluster)
	metrics / labels / label_values / series
\odot	rules / alerts / test_rules
\odot	active_queries / top_queries
0	query / query_range
	tsdb_status (cardinality explorer) / metric_statistic
	metrics_relabel_debug / downsampling_filters_debug / retention_filters_debug
	explain_query / prettify_query
	documentation



	tenants (for cluster)
	metrics / labels / label_values / series
	rules / alerts / test_rules
\odot	active_queries / top_queries
\odot	query / query_range
\odot	tsdb_status (cardinality explorer) / metric_statistic
	metrics_relabel_debug / downsampling_filters_debug / retention_filters_debug
	explain_query / prettify_query
	documentation















Show me 3 slowest queries

I'll help you find the 3 slowest queries by checking the top queries with the biggest average execution duration.



top_queries



I'll help you find the 3 slowest queries by checking the top queries with the biggest average execution duration.



Here are the 3 slowest queries by average execution duration in your VictoriaMetrics instance:

Top 3 Slowest Queries

1. Filesystem Space Prediction (4h) - 0.249 seconds

```
( node_filesystem_avail_bytes{job="node-exporter",fstype!=""} /
node_filesystem_size_bytes{job="node-exporter",fstype!=""} * 100 < 15 and
predict_linear(node_filesystem_avail_bytes{job="node-exporter",fstype!=""}[6h], 4*60*60) < 0 and
node_filesystem_readonly{job="node-exporter",fstype!=""} == 0 )</pre>
```

2. Filesystem Inodes Prediction (24h) - 0.206 seconds

```
( node_filesystem_files_free{job="node-exporter",fstype!=""} /
node_filesystem_files{job="node-exporter",fstype!=""} * 100 < 40 and
predict_linear(node_filesystem_files_free{job="node-exporter",fstype!=""}[6h], 24*60*60) < 0 and
node_filesystem_readonly{job="node-exporter",fstype!=""} == 0 )</pre>
```

3. Filesystem Inodes Prediction (4h) - 0.194 seconds

```
( node_filesystem_files_free{job="node-exporter",fstype!=""} /
node_filesystem_files{job="node-exporter",fstype!=""} * 100 < 20 and
predict_linear(node_filesystem_files_free{job="node-exporter",fstype!=""}[6h], 4*60*60) < 0 and
node_filesystem_readonly{job="node-exporter",fstype!=""} == 0 )</pre>
```



• •

```
node_filesystem_size_bytes{job="node-exporter",fstype!=""} * 100 < 15 and
predict_linear(node_filesystem_avail_bytes{job="node-exporter",fstype!=""}[6h], 4*60*60) < 0 and
node_filesystem_readonly{job="node-exporter",fstype!=""} == 0 )</pre>
```

2. Filesystem Inodes Prediction (24h) - 0.206 seconds

```
( node_filesystem_files_free{job="node-exporter",fstype!=""} /
node_filesystem_files{job="node-exporter",fstype!=""} * 100 < 40 and
predict_linear(node_filesystem_files_free{job="node-exporter",fstype!=""}[6h], 24*60*60) < 0 and
node_filesystem_readonly{job="node-exporter",fstype!=""} == 0 )</pre>
```

3. Filesystem Inodes Prediction (4h) - 0.194 seconds

```
( node_filesystem_files_free{job="node-exporter",fstype!=""} /
node_filesystem_files{job="node-exporter",fstype!=""} * 100 < 20 and
predict_linear(node_filesystem_files_free{job="node-exporter",fstype!=""}[6h], 4*60*60) < 0 and
node_filesystem_readonly{job="node-exporter",fstype!=""} == 0 )</pre>
```

Key Observations

All of your slowest queries are related to filesystem monitoring and use predict_linear() functions, which are computationally expensive as they perform linear regression calculations over time series data. These queries are likely part of alerting rules for predicting when filesystems will run out of space or inodes.

The slowest query takes about 249ms on average, which suggests your VictoriaMetrics instance is performing well overall with relatively fast query execution times.

Let's work with the first one. Please prettify it





```
node_fitesystem_fites_free{job= node-exporter ',fstype!=""} * 100 < 40 and
predict_linear(node_filesystem_files_free{job="node-exporter",fstype!=""}[6h], 24*60*60) < 0 and
node_filesystem_readonly{job="node-exporter",fstype!=""} == 0 )</pre>
```

3. Filesystem Inodes Prediction (4h) - 0.194 seconds

```
( node_filesystem_files_free{job="node-exporter",fstype!=""} /
node_filesystem_files{job="node-exporter",fstype!=""} * 100 < 20 and
predict_linear(node_filesystem_files_free{job="node-exporter",fstype!=""}[6h], 4*60*60) < 0 and
node_filesystem_readonly{job="node-exporter",fstype!=""} == 0 )</pre>
```

Key Observations

All of your slowest queries are related to filesystem monitoring and use predict_linear() functions, which are computationally expensive as they perform linear regression calculations over time series data. These queries are likely part of alerting rules for predicting when filesystems will run out of space or inodes.

The slowest query takes about 249ms on average, which suggests your VictoriaMetrics instance is performing well overall with relatively fast query execution times.

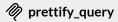
Let's work with the first one. Please prettify it

I'll prettify the slowest query for you using VictoriaMetrics' query formatting endpoint.





I'll prettify the slowest query for you using VictoriaMetrics' query formatting endpoint.



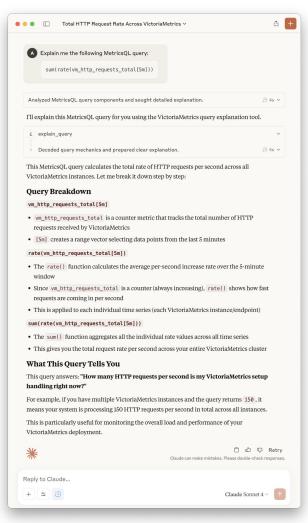
Here's the prettified version of your slowest query:

```
node_filesystem_avail_bytes{job="node-exporter",fstype!=""}
       node filesystem size bytes{job="node-exporter",fstype!=""}
      100
   and
   predict linear(
     node filesystem avail bytes{job="node-exporter",fstype!=""}[6h],
(node filesystem readonly{job="node-exporter",fstype!=""} == 0)
```

The prettified version makes the query structure much clearer:

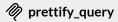


• • •





I'll prettify the slowest query for you using VictoriaMetrics' query formatting endpoint.



Here's the prettified version of your slowest query:

```
node_filesystem_avail_bytes{job="node-exporter",fstype!=""}
       node filesystem size bytes{job="node-exporter",fstype!=""}
      100
   and
   predict linear(
     node filesystem avail bytes{job="node-exporter",fstype!=""}[6h],
(node filesystem readonly{job="node-exporter",fstype!=""} == 0)
```

The prettified version makes the query structure much clearer:



• • •

```
node_filesystem_avail_bytes{job="node-exporter",fstype!=""}

node_filesystem_size_bytes{job="node-exporter",fstype!=""}

node_filesystem_size_bytes{job="node-exporter",fstype!=""}

promq!

promq!

/

15

and

(node_filesystem_avail_bytes{job="node-exporter",fstype!=""}[6h],

(4 * 60) * 60

)

and

(node_filesystem_readonly{job="node-exporter",fstype!=""} == 0)
```

The prettified version makes the query structure much clearer:

- 1. Main condition: Filesystem usage > 85% (less than 15% available)
- 2. **Prediction condition**: Linear prediction suggests the filesystem will be full within 4 hours (14,400 seconds)
- 3. **Safety condition**: Filesystem is not read-only

This query is designed to alert when a filesystem is both critically low on space AND predicted to run out completely within 4 hours based on the current trend.

now explain me this query









The prettified version makes the query structure much clearer:

- 1. Main condition: Filesystem usage > 85% (less than 15% available)
- 2. **Prediction condition**: Linear prediction suggests the filesystem will be full within 4 hours (14,400 seconds)
- 3. **Safety condition**: Filesystem is not read-only

This query is designed to alert when a filesystem is both critically low on space AND predicted to run out completely within 4 hours based on the current trend.

now explain me this query







explain_query

What This Query Does

This query identifies filesystems that are **critically low on disk space AND predicted to run out completely within 4 hours**. It's a proactive alerting rule for preventing disk space exhaustion.

Step-by-Step Breakdown

1. Current Usage Check (First Condition)

```
(node_filesystem_avail_bytes{job="node-exporter",fstype!=""} / node_filesystem_size_bytes{job="node-exporter",fstype!=""}) * 100 < 15
```

What it does:

- node filesystem avail bytes: Available space in bytes (space non-root users can use)
- node_filesystem_size_bytes: Total filesystem size in bytes
- **Division**: Calculates the ratio of available space to total space
- Multiply by 100: Converts to percentage
- < 15: Finds filesystems with less than 15% free space

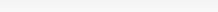
Filters applied:

- job="node-exporter": Only data from Node Exporter
- fstype!="": Excludes entries without filesystem type (filters out pseudo-filesystems)

2. Future Prediction Check (Second Condition)

```
predict_linear(node_filesystem_avail_bytes{job="node-exporter",fstype!=""}[6h], (4 * 60) * 60) < 0</pre>
```

What it does:





- **Division**: Calculates the ratio of available space to total space
- Multiply by 100: Converts to percentage
- < 15: Finds filesystems with less than 15% free space

Filters applied:

- job="node-exporter": Only data from Node Exporter
- fstype!="": Excludes entries without filesystem type (filters out pseudo-filesystems)

2. Future Prediction Check (Second Condition)

```
predict_linear(node_filesystem_avail_bytes{job="node-exporter",fstype!=""}[6h], (4 * 60) * 60) < 0</pre>
```

What it does:

- predict_linear(): Uses linear regression to predict future values
- [6h]: Analyzes the last 6 hours of data to establish the trend
- (4 * 60) * 60 = 14,400 seconds: Predicts the value 4 hours into the future
- < 0: Checks if the predicted available space will be negative (completely full)

How prediction works:

- Looks at how disk space has been decreasing over the past 6 hours
- Extrapolates this trend 4 hours into the future
- If the trend suggests 0 or negative bytes available, it triggers

3. Read-Only Filter (Third Condition)

```
node_filesystem_readonly{job="node-exporter",fstype!=""} == 0
```

. . .

What it does:

- node_filesystem_readonly: Indicates if filesystem is mounted read-only (1) or read-write (0)
- == 0: Only includes read-write filesystems
- Why this matters: No point alerting on read-only filesystems since you can't free up space on them



3. Read-Only Filter (Third Condition

```
node_filesystem_readonly{job="node-exporter",fstype!=""} == 0
```

What it does:

- node_filesystem_readonly: Indicates if filesystem is mounted read-only (1) or read-write (0)
- == 0: Only includes read-write filesystems
- Why this matters: No point alerting on read-only filesystems since you can't free up space on them

Combined Logic

The query uses **AND** operators, so **ALL** three conditions must be true:

```
(Current free space < 15%)

AND
(Predicted to be full in 4 hours)

AND
(Filesystem is writable)
```

Why This Query is Slow

- 1. predict_linear() function: Computationally expensive as it performs linear regression analysis
- 2. 6-hour lookback window: Must analyze 6 hours of historical data for each filesystem
- 3. Multiple metric queries: Three different metrics with filtering
- 4. Label matching: Complex filtering across multiple label dimensions

Practical Use Case

This would typically fire alerts like:

• "Filesystem /var/log on server-01 is 12% full and predicted to be completely full in 4 hours based on current usage trends"

This gives administrators time to take preventive action before the disk actually fills up, which is much better than alerting only when the disk is already 95% full.

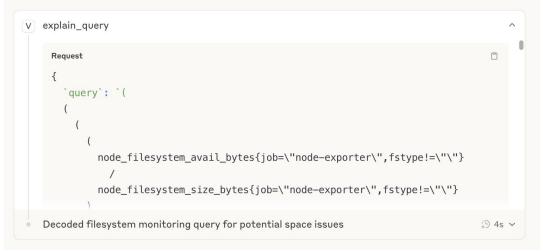


I'll analyze this MetricsQL query for you. Let me first get a detailed explanation using VictoriaMetrics' query explanation tool.

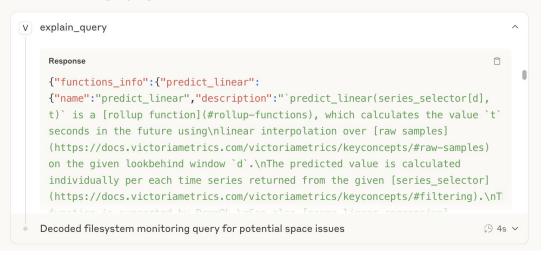
V explain_query

Decoded filesystem monitoring query for potential space issues

I'll analyze this MetricsQL query for you. Let me first get a detailed explanation using VictoriaMetrics' query explanation tool.



I'll analyze this MetricsQL query for you. Let me first get a detailed explanation using VictoriaMetrics' query explanation tool.



```
explain query → response.json
      "name": "predict linear",
      "description": "`predict_linear(series_selector[d], t)` is a [rollup function](#rollup-functions), which calculates the value `t` seconds
in the future using\nlinear interpolation over [raw samples](https://docs.victoriametrics.com/victoriametrics/keyconcepts/#raw-samples) on the
given lookbehind window `d`.\nThe predicted value is calculated individually per each time series returned from the given
[series selector](https://docs.victoriametrics.com/victoriametrics/keyconcepts/#filtering).\nThis function is supported by PromOL.\nSee also
[range linear regression](#range linear regression).",
      "category": "Rollup functions"
      "group": "node-exporter",
     "name": "node_filesystem_avail_bytes",
      "help": "Filesystem space available to non-root users in bytes.",
      "type": "gauge".
      "labels": ["device", "device error", "fstype", "mountpoint"]
   "node filesystem readonly": { "group": "node-exporter", "name": "node filesystem readonly", "help": "Filesystem read-only status.", "type":
"gauge", "labels": ["device", "device_error", "fstype", "mountpoint"] },
    "node_filesystem_size_bytes": {    "group": "node-exporter", "name": "node_filesystem_size_bytes", "help": "Filesystem size in bytes.",
'type": "gauge", "labels": ["device", "device_error", "fstype", "mountpoint"] }
  "syntax_tree": {
    'group_modifier": { "args": null, "op": "", "type": "ModifierExpr" },
    "join_modifier": { "args": null, "op": "", "type": "ModifierExpr" },
      "group_modifier": { "args": null, "op": "", "type": "ModifierExpr" },
      "join_modifier": {        "args": null, "op": "", "type": "ModifierExpr"    },
```

```
explain query → response.json
     "name": "predict_linear",
     "description": "`predict_linear(series_selector[d], t)` is a [rollup function](#rollup-functions), which calculates the value `t` seconds
in the future using\nlinear interpolation over [raw samples](https://docs.victoriametrics.com/victoriametrics/keyconcepts/#raw-samples) on the
given lookbehind window `d`.\nThe predicted value is calculated individually per each time series returned from the given
[series_selector](https://docs.victoriametrics.com/victoriametrics/keyconcepts/#filtering).\nThis function is supported by PromQL.\nSee also
[range linear regression](#range linear regression).",
     "category": "Rollup functions"
```

```
explain query → response.json
     "group": "node-exporter",
     "name": "node_filesystem_avail_bytes",
     "help": "Filesystem space available to non-root users in bytes.",
     "type": "gauge",
     "labels": ["device", "device error", "fstype", "mountpoint"]
   "node_filesystem_readonly": { "group": "node-exporter", "name": "node_filesystem_readonly", "help": "Filesystem read-only status.", "type":
"gauge", "labels": ["device", "device_error", "fstype", "mountpoint"] },
   "node_filesystem_size_bytes": {  "group": "node-exporter", "name": "node_filesystem_size_bytes", "help": "Filesystem size in bytes.",
```

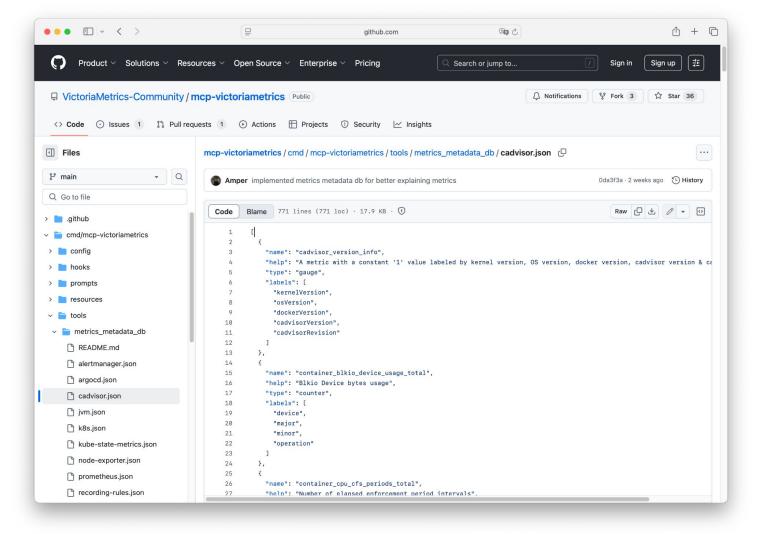
```
explain_query → response.json
```

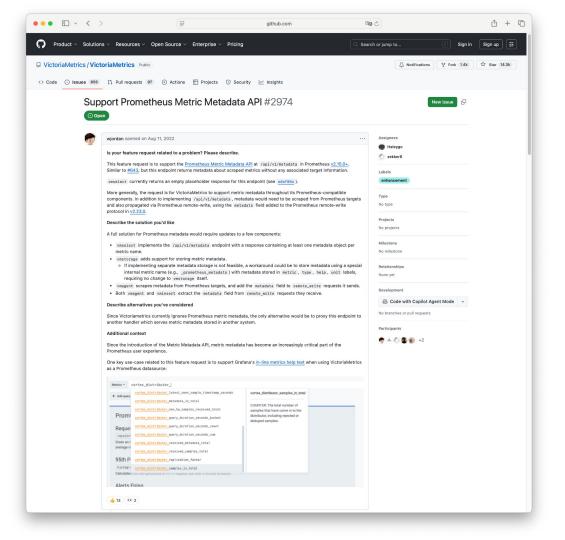
```
"group_modifier": {        "args": null, "op": "", "type": "ModifierExpr" },
"join_modifier": { "args": null, "op": "", "type": "ModifierExpr" },
  "group_modifier": { "args": null, "op": "", "type": "ModifierExpr" },
  "join_modifier": { "args": null, "op": "", "type": "ModifierExpr" },
```

```
explain query → response.json
      "name": "Binary operation",
      "description": "BinaryOpExpr represents binary operation such as `+`, `-`, `*`, `/`, etc.",
        "bool": { "name": "", "description": "Bool indicates whether `bool` modifier is present. For example, `foo \u003ebool bar`",
"data_type": "bool" },
        "group modifier": { "name": "", "description": "GroupModifier contains modifier such as \"on\" or \"ignoring\".", "data type":
"ModifierExpr" },
        "join modifier": { "name": "", "description": "JoinModifier contains modifier such as \"group left\" or \"group right\".", "data type":
"ModifierExpr" },
          "description": "JoinModifierPrefix is an optional prefix to add to labels specified inside group left() or group right() lists. The
syntax is `group_left(foo,bar) prefix \"abc\"`",
         "data_type": "StringExpr"
        "keep metric names": { "name": "", "description": "If KeepMetricNames is set to true, then the operation should keep metric names.",
"data type": "bool" },
        "left": { "name": "", "description": "Left contains left arg for the `left op right` expression.", "data type": "Expr" },
        "op": { "name": "", "description": "Op is the operation itself, i.e. `+`, `-`, `*`, etc.", "data_type": "string" },
        "right": { "name": "", "description": "Right contains right arg for the `left op right` expression.", "data type": "Expr" }
      "name": "Duration",
      "description": "DurationExpr represents a duration, e.g. `5m`, `1h`. Supported suffixes are `s` (seconds), `m` (minutes), `h` (hours),
`d` (days), `w` (weeks), and `y` (years).",
      "fields": { "value": { "name": "", "description": "Value is the duration value as a string, e.g. `5m`, `1h`.", "data type": "string" } }
      "name": "Function",
      "description": "uncExpr represents MetricsQL function such as `foo(...)`",
        "args": { "name": "", "description": "Arguments of the function, which can be other expressions.", "data type": "array of Expr" },
        "keep_metric_names": {    "name": "",    "description": "If KeepMetricNames is set to true, then the function should keep metric names.",
'data_type": "bool" },
```

```
explain query → response.json
      "name": "predict linear",
      "description": "`predict_linear(series_selector[d], t)` is a [rollup function](#rollup-functions), which calculates the value `t` seconds
in the future using\nlinear interpolation over [raw samples](https://docs.victoriametrics.com/victoriametrics/keyconcepts/#raw-samples) on the
given lookbehind window `d`.\nThe predicted value is calculated individually per each time series returned from the given
[series selector](https://docs.victoriametrics.com/victoriametrics/keyconcepts/#filtering).\nThis function is supported by PromOL.\nSee also
[range linear regression](#range linear regression).",
      "category": "Rollup functions"
      "group": "node-exporter",
     "name": "node_filesystem_avail_bytes",
      "help": "Filesystem space available to non-root users in bytes.",
      "type": "gauge".
      "labels": ["device", "device error", "fstype", "mountpoint"]
   "node filesystem readonly": { "group": "node-exporter", "name": "node filesystem readonly", "help": "Filesystem read-only status.", "type":
"gauge", "labels": ["device", "device_error", "fstype", "mountpoint"] },
    "node_filesystem_size_bytes": {    "group": "node-exporter", "name": "node_filesystem_size_bytes", "help": "Filesystem size in bytes.",
'type": "gauge", "labels": ["device", "device_error", "fstype", "mountpoint"] }
  "syntax_tree": {
    'group_modifier": { "args": null, "op": "", "type": "ModifierExpr" },
    "join_modifier": { "args": null, "op": "", "type": "ModifierExpr" },
      "group_modifier": { "args": null, "op": "", "type": "ModifierExpr" },
      "join_modifier": {        "args": null, "op": "", "type": "ModifierExpr"    },
```

```
explain query → response.json
     "group": "node-exporter",
     "name": "node_filesystem_avail_bytes",
     "help": "Filesystem space available to non-root users in bytes.",
     "type": "gauge",
     "labels": ["device", "device error", "fstype", "mountpoint"]
   "node_filesystem_readonly": { "group": "node-exporter", "name": "node_filesystem_readonly", "help": "Filesystem read-only status.", "type":
"gauge", "labels": ["device", "device_error", "fstype", "mountpoint"] },
   "node_filesystem_size_bytes": {  "group": "node-exporter", "name": "node_filesystem_size_bytes", "help": "Filesystem size in bytes.",
```





The guery uses **AND** operators, so **ALL** three conditions must be true:

```
(Current free space < 15%)
(Predicted to be full in 4 hours)
(Filesystem is writable)
```

Why This Query is Slow

- 1. predict_linear() function: Computationally expensive as it performs linear regression analysis
- 2. 6-hour lookback window: Must analyze 6 hours of historical data for each filesystem
- 3. Multiple metric queries: Three different metrics with filtering
- 4. Label matching: Complex filtering across multiple label dimensions

Practical Use Case

This would typically fire alerts like:

• "Filesystem /var/log on server-01 is 12% full and predicted to be completely full in 4 hours based on current usage trends"

This gives administrators time to take preventive action before the disk actually fills up, which is much better than alerting only when the disk is already 95% full.

please trace this query









(Filesystem is writable)

Why This Query is Slow

- 1. predict_linear() function: Computationally expensive as it performs linear regression analysis
- 2. 6-hour lookback window: Must analyze 6 hours of historical data for each filesystem
- 3. Multiple metric queries: Three different metrics with filtering
- 4. Label matching: Complex filtering across multiple label dimensions

Practical Use Case

This would typically fire alerts like:

AAA

• "Filesystem /var/log on server-01 is 12% full and predicted to be completely full in 4 hours based on current usage trends"

This gives administrators time to take preventive action before the disk actually fills up, which is much better than alerting only when the disk is already 95% full.

please trace this query

I'll trace the execution of this guery to show you exactly how VictoriaMetrics processes it step by step.



query



I'll trace the execution of this guery to show you exactly how VictoriaMetrics processes it step by step.



Here's a detailed trace analysis of your query execution:

Query Execution Summary

- Total Time: 217.9ms
- Series Fetched: 3,309 time series
- Result: 0 matching filesystems (none currently meet all criteria)
- Data Scanned: Over 2.1 million samples

Execution Breakdown

1. Read-Only Filesystem Check (21ms)

```
node_filesystem_readonly{job="node-exporter",fstype!=""} == 0
```

.- , -- (

- Time: 21ms
- Series Found: 1,110 filesystems
- Purpose: Filter out read-only filesystems
- **Result**: Found 1,110 writable filesystems

2. Current Usage Percentage (75ms)

```
(node\_filesystem\_avail\_bytes \ / \ node\_filesystem\_size\_bytes) \ * \ 100 \ < \ 15
```

. . .



2. Current Usage Percentage (75ms)

(node_filesystem_avail_bytes / node_filesystem_size_bytes) * 100 < 15</pre>

promgl

- **Time**: 75ms
- Series Processed: 733 filesystems
- Data Fetched: 67,816 samples
- Storage Calls: Parallel calls to 2 storage nodes
- Result: Found 733 filesystems with <15% free space

3. Prediction Analysis (113ms) - The Bottleneck

predict_linear(node_filesystem_avail_bytes[6h], 14400) < 0</pre>

omal

- Time: 113ms (52% of total execution time)
- Series Processed: 733 filesystems
- Lookback Window: 6 hours (21,600 seconds)
- Data Analyzed: 1,244,406 samples over 6 hours
- Storage Impact: Fetched 9,913 blocks from storage
- Linear Regression: Performed on each of 733 time series
- Result: None predicted to be full in 4 hours

Why This Query is Slow

Primary Bottleneck: predict_linear() Function

- 94ms spent on linear regression calculations
- 1.2M+ samples analyzed for trend calculation
- 6-hour lookback requires significant historical data
- No optimization available for predict_linear() function





• Parallel processing - storage nodes worked simultaneously

• Memory efficient - only 744KB RAM needed for calculations

Performance Metrics

• Samples Scanned: 2,111,040 total • Blocks Fetched: 25.635 blocks

• Network Efficiency: Good - parallel storage access • Cache Hit Rate: High - used tag filters cache effectively

Final Result

Currently, no filesystems meet all three criteria:

- 1,110 filesystems are writable
- **733 filesystems** have <15% free space
- X 0 filesystems are predicted to be full in 4 hours

The prediction algorithm determined that none of the low-space filesystems have a consumption trend steep enough to exhaust space within 4 hours, which is actually good news for your infrastructure!

show me the graph for this query for last 24 hours









Performance Metrics

- Samples Scanned: 2,111,040 total
- Blocks Fetched: 25,635 blocks
- Network Efficiency: Good parallel storage access
- Cache Hit Rate: High used tag filters cache effectively

Final Result

Currently, no filesystems meet all three criteria:

- **1,110 filesystems** are writable
- **733 filesystems** have <15% free space
- X 0 filesystems are predicted to be full in 4 hours

The prediction algorithm determined that none of the low-space filesystems have a consumption trend steep enough to exhaust space within 4 hours, which is actually good news for your infrastructure!

show me the graph for this query for last 24 hours

I'll fetch the data for this query over the last 24 hours and create a graph for you.



query_range



inal Resul

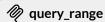
Currently, no filesystems meet all three criteria:

- 1,110 filesystems are writable
- **733 filesystems** have < 15% free space
- X 0 filesystems are predicted to be full in 4 hours

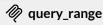
The prediction algorithm determined that none of the low-space filesystems have a consumption trend steep enough to exhaust space within 4 hours, which is actually good news for your infrastructure!

show me the graph for this query for last 24 hours

I'll fetch the data for this guery over the last 24 hours and create a graph for you.



The query returned no results over the past 24 hours, which means no filesystems met all three criteria simultaneously. This is actually good news! Let me create a more informative dashboard by breaking down the query components to show you what's happening with your filesystem monitoring.





o mesystems are predicted to be full in 4 hours

The prediction algorithm determined that none of the low-space filesystems have a consumption trend steep enough to exhaust space within 4 hours, which is actually good news for your infrastructure!

show me the graph for this query for last 24 hours

I'll fetch the data for this query over the last 24 hours and create a graph for you.



The query returned no results over the past 24 hours, which means no filesystems met all three criteria simultaneously. This is actually good news! Let me create a more informative dashboard by breaking down the query components to show you what's happening with your filesystem monitoring.





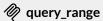
AAA



exhaust space within 4 hours, which is actually good news for your infrastructure!

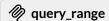
show me the graph for this query for last 24 hours

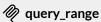
I'll fetch the data for this query over the last 24 hours and create a graph for you.



The query returned no results over the past 24 hours, which means no filesystems met all three criteria simultaneously. This is actually good news! Let me create a more informative dashboard by breaking down the query components to show you what's happening with your filesystem monitoring.







AAA

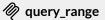


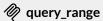
I'll fetch the data for this query over the last 24 hours and create a graph for you.



The query returned no results over the past 24 hours, which means no filesystems met all three criteria simultaneously. This is actually good news! Let me create a more informative dashboard by breaking down the query components to show you what's happening with your filesystem monitoring.







Perfect! Now I have filesystem data to create a meaningful visualization. I can see the filesystem usage patterns over the past few hours. Let me create an interactive dashboard showing the filesystem monitoring data.

Filesystem Monitoring Dashboard

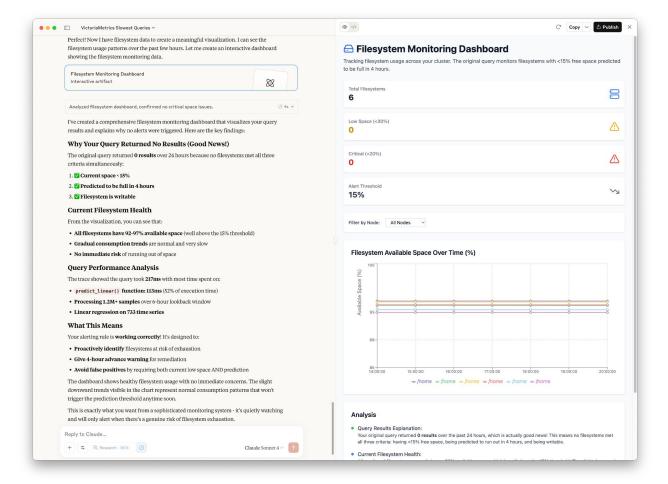
Interactive artifact



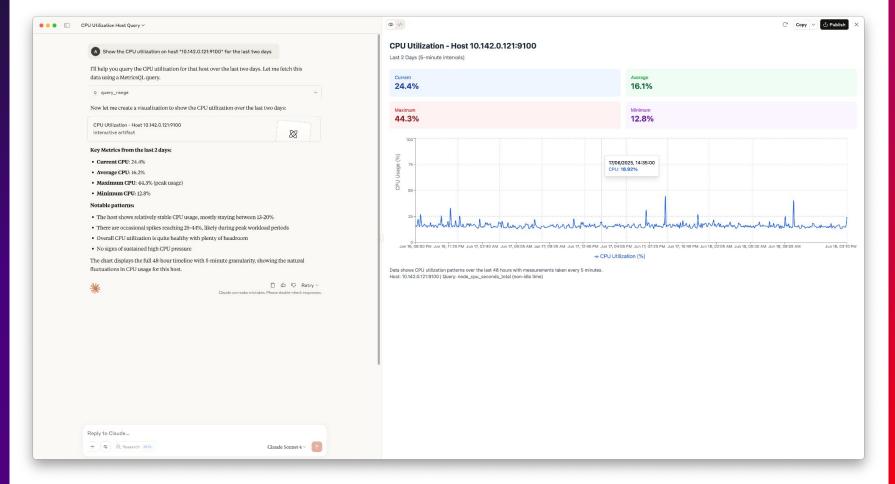
I've created a comprehensive filesystem monitoring dashboard that visualizes your query results and explains why no alerts were triggered. Here are the key findings:

• • •



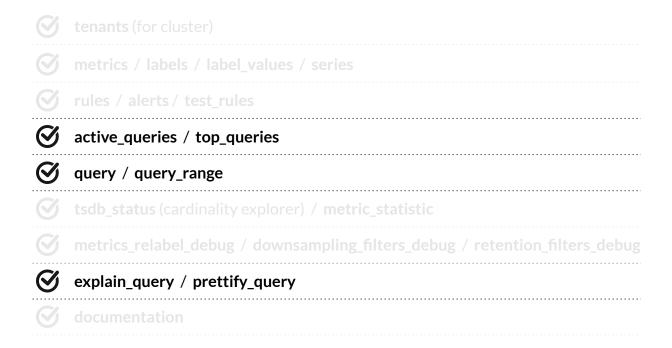








Available tools





Available tools

	tenants (for cluster)
	metrics / labels / label_values / series
	rules / alerts / test_rules
	active_queries / top_queries
0	query / query_range
\otimes	tsdb_status (cardinality explorer) / metric_statistic
0	metrics_relabel_debug / downsampling_filters_debug / retention_filters_debug
	explain_query / prettify_query
	documentation



Available tools

	tenants (for cluster)
	metrics / labels / label_values / series
	rules / alerts / test_rules
	active_queries / top_queries
	query / query_range
\otimes	tsdb_status (cardinality explorer) / metric_statistic
\otimes	metrics_relabel_debug / downsampling_filters_debug / retention_filters_debug
0	explain_query / prettify_query
	documentation











Give me the top metric names by cardinality

 $I'll\ get\ the\ top\ metric\ names\ by\ cardinality\ from\ your\ Victoria\ Metrics\ instance.$

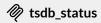


tsdb_status



Give me the top metric names by cardinality

I'll get the top metric names by cardinality from your Victoria Metrics instance.



Here are the top 20 metric names by cardinality in your VictoriaMetrics instance:

Rank	Metric Name	Series Count	Requests Count
1	etcd_request_duration_seconds_bucket	11,304	2,658,477
2	apiserver_request_duration_seconds_bucket	10,920	3,606,927
3	github_downloads_total	8,387	4,181,762
4	apiserver_request_sli_duration_seconds_bucket	6,908	1,751,128
20	kube_node_status_condition	755	590,471

Key Observations:





17		1,517	446,691
18	go_sched_latencies_seconds_bucket	1,517	381,902,240
19	csi_operations_seconds_bucket	1,498	109,372
20	kube_node_status_condition	755	590,471

Key Observations:

- Total series across all metrics: 249,452
- **Top contributors**: The etcd and apiserver duration bucket metrics have the highest cardinality, likely due to multiple label combinations (different endpoints, methods, status codes, etc.)
- **High query activity**: go_sched_latencies_seconds_bucket has extremely high request count (381M+) despite moderate cardinality
- Bucket metrics dominate: Most high-cardinality metrics are histogram buckets (_bucket suffix), which is typical as they create multiple series per metric

The histogram bucket metrics typically have high cardinality because they create separate time series for each bucket boundary (le label) combined with other dimensional labels.

Which of them are never queried?







19	csi_operations_seconds_bucket	1,498	109,372
20	kube_node_status_condition	755	590,471

Key Observations:

• Total series across all metrics: 249,452

A A A

- **Top contributors**: The etcd and apiserver duration bucket metrics have the highest cardinality, likely due to multiple label combinations (different endpoints, methods, status codes, etc.)
- **High query activity**: go_sched_latencies_seconds_bucket has extremely high request count (381M+) despite moderate cardinality
- Bucket metrics dominate: Most high-cardinality metrics are histogram buckets (_bucket suffix), which is typical as they create multiple series per metric

The histogram bucket metrics typically have high cardinality because they create separate time series for each bucket boundary (le label) combined with other dimensional labels.

Which of them are never queried?

I can see from the data that all of the top 20 metrics by cardinality have been used in queries (they all have non-zero request counts). However, let me check for metrics that are rarely or never queried to give you a better picture.





Which of them are never queried?

I can see from the data that all of the top 20 metrics by cardinality have been used in queries (they all have non-zero request counts). However, let me check for metrics that are rarely or never queried to give you a better picture.



metric_stastistics

Based on the query statistics, all of the top 20 metrics by cardinality are actively being used in queries. None of them have zero query requests.

However, I did find some completely unused metrics (0 queries) in your system, though these are not among the high-cardinality ones:

Completely unused metrics (0 queries):

- apiserver_nodeport_repair_reconcile_errors_total
- kube_pod_init_container_status_last_terminated_reason
- vmanomaly reader response parsing seconds created

Key findings about your top cardinality metrics:

- Most actively queried: go_sched_latencies_seconds_bucket with 381M+ queries
- Least queried among top 20: volume_operation_total_seconds_bucket with 96,335 queries
- All have substantial usage: Even the least queried high-cardinality metric has nearly 100k queries

This actually indicates a **healthy situation** - your high-cardinality metrics are being actively used, meaning they're providing value and not just consuming storage space unnecessarily.

If you're concerned about storage costs, you might want to focus on:

- Retention policies for less frequently queried metrics
- Downsampling strategies for histogram buckets that don't need full granularity
- Label optimization to reduce unnecessary cardinality combinations



Based on the query statistics, all of the top 20 metrics by cardinality are actively being used in queries. None of them have zero query requests.

However, I did find some completely unused metrics (O queries) in your system, though these are not among the high-cardinality ones:

Completely unused metrics (0 queries):

- apiserver_nodeport_repair_reconcile_errors_total
- kube_pod_init_container_status_last_terminated_reason
- vmanomaly reader response parsing seconds created

Key findings about your top cardinality metrics:

- Most actively queried: go sched latencies seconds bucket with 381M+ queries
- Least queried among top 20: volume_operation_total_seconds_bucket with 96,335 queries
- All have substantial usage: Even the least queried high-cardinality metric has nearly 100k queries

This actually indicates a healthy situation - your high-cardinality metrics are being actively used, meaning they're providing value and not just consuming storage space unnecessarily.

- If you're concerned about storage costs, you might want to focus on: • Retention policies for less frequently gueried metrics
 - Downsampling strategies for histogram buckets that don't need full granularity
 - Label optimization to reduce unnecessary cardinality combinations

let's create relabeling rule to remove these completely unused metrics









high-cardinality ones:

Completely unused metrics (0 queries):

- apiserver_nodeport_repair_reconcile_errors_total
- kube_pod_init_container_status_last_terminated_reason
- vmanomaly_reader_response_parsing_seconds_created

Key findings about your top cardinality metrics:

- Most actively queried: go_sched_latencies_seconds_bucket with 381M+ queries
- Least queried among top 20: volume_operation_total_seconds_bucket with 96,335 queries
- All have substantial usage: Even the least queried high-cardinality metric has nearly 100k queries

This actually indicates a **healthy situation** - your high-cardinality metrics are being actively used, meaning they're providing value and not just consuming storage space unnecessarily.

If you're concerned about storage costs, you might want to focus on:

- Retention policies for less frequently queried metrics
- **Downsampling** strategies for histogram buckets that don't need full granularity
- Label optimization to reduce unnecessary cardinality combinations

ok, let's create relabeling rule to drop all completely unused metrics

Great idea! Let me first get a complete list of completely unused metrics (0 queries), then we'll create relabeling rules to drop them.



metric_stastistics



- Least queried among top 20: volume operation total seconds bucket with 96,335 queries
- All have substantial usage: Even the least queried high-cardinality metric has nearly 100k queries

This actually indicates a healthy situation - your high-cardinality metrics are being actively used, meaning they're providing value and not just consuming storage space unnecessarily.

If you're concerned about storage costs, you might want to focus on:

- Retention policies for less frequently gueried metrics
- Downsampling strategies for histogram buckets that don't need full granularity
- Label optimization to reduce unnecessary cardinality combinations

ok, let's create relabeling rule to drop all completely unused metrics

> Great idea! Let me first get a complete list of completely unused metrics (0 queries), then we'll create relabeling rules to drop them.



metric_stastistics

Perfect! I found several completely unused metrics (O queries). Let me create relabeling rules to drop them. Here are the metrics with exactly **0 queries**:

Relabeling Rules To Drop Unused Metrics Code







providing value and not just consuming storage space unnecessarily.

If you're concerned about storage costs, you might want to focus on:

- Retention policies for less frequently queried metrics
- Downsampling strategies for histogram buckets that don't need full granularity
- Label optimization to reduce unnecessary cardinality combinations

ok, let's create relabeling rule to drop all completely unused metrics

Great idea! Let me first get a complete list of completely unused metrics (0 queries), then we'll create relabeling rules to drop them.



metric_stastistics

Perfect! I found several completely unused metrics (0 queries). Let me create relabeling rules to drop them. Here are the metrics with exactly **0 queries**:

Relabeling Rules To Drop Unused Metrics





Now let me test these relabeling rules using the debug tool to ensure they work correctly:



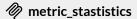
metric_relabel_debug



- Downsampling strategies for histogram buckets that don't need full granularity
- Label optimization to reduce unnecessary cardinality combinations

ok, let's create relabeling rule to drop all completely unused metrics

Great idea! Let me first get a complete list of completely unused metrics (0 queries), then we'll create relabeling rules to drop them.



Perfect! I found several completely unused metrics (0 queries). Let me create relabeling rules to drop them. Here are the metrics with exactly **0 queries**:

Relabeling Rules To Drop Unused Metrics





Now let me test these relabeling rules using the debug tool to ensure they work correctly:



Let me test another unused metric to verify the pattern works:





ok, let's create relabeling rule to drop all completely unused metrics

Great idea! Let me first get a complete list of completely unused metrics (0 queries), then we'll create relabeling rules to drop them.



metric_stastistics

Perfect! I found several completely unused metrics (0 queries). Let me create relabeling rules to drop them. Here are the metrics with exactly **0 queries**:

Relabeling Rules To Drop Unused Metrics





Now let me test these relabeling rules using the debug tool to ensure they work correctly:



metric_relabel_debug

Let me test another unused metric to verify the pattern works:



metric_relabel_debug

Perfect! Now let me test that we don't accidentally drop metrics that should be kept:



metric_relabel_debug



Great idea! Let me first get a complete list of completely unused metrics (0 queries), then we'll create relabeling rules to drop them.



metric_stastistics

Perfect! I found several completely unused metrics (0 queries). Let me create relabeling rules to drop them. Here are the metrics with exactly **0 queries**:

Relabeling Rules To Drop Unused Metrics



Now let me test these relabeling rules using the debug tool to ensure they work correctly:



metric_relabel_debug

Let me test another unused metric to verify the pattern works:



metric_relabel_debug

Perfect! Now let me test that we don't accidentally drop metrics that should be kept:



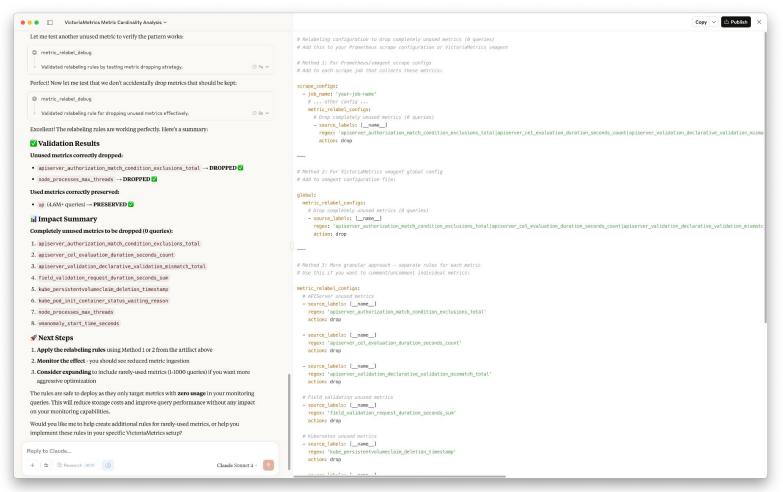
metric_relabel_debug

Excellent! The relabeling rules are working perfectly. Here's a summary:

Validation Results

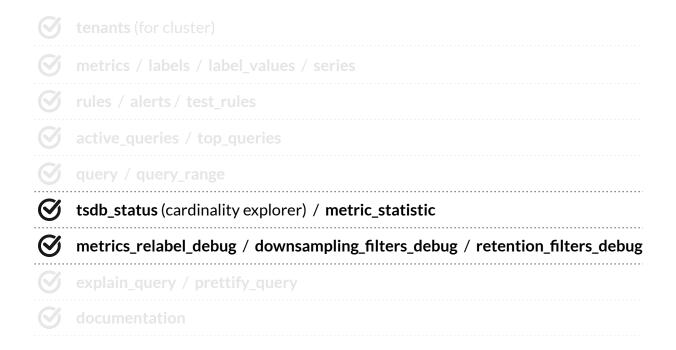








Available tools





Available tools

	tenants (for cluster)
	metrics / labels / label_values / series
	rules / alerts / test_rules
	active_queries / top_queries
	query / query_range
	tsdb_status (cardinality explorer) / metric_statistic
	metrics_relabel_debug / downsampling_filters_debug / retention_filters_debug
	explain_query / prettify_query
\odot	documentation











Great idea! Let me first get a complete list of completely unused metrics (0 queries), then we'll create relabeling rules to drop them.

Thinking about writing metrics to VictoriaMetrics Cloud platform

AAA



documentation



Great idea! Let me first get a complete list of completely unused metrics (0 queries), then we'll create relabeling rules to drop them.

Thinking about writing metrics to Victoria Metrics Cloud platform

AAA

- **documentation**
- Based on the documentation search results, i can see information about various ways to...
- **documentation**

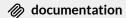


Great idea! Let me first get a complete list of completely unused metrics (0 queries), then we'll create relabeling rules to drop them.

• Thinking about writing metrics to Victoria Metrics Cloud platform



Based on the documentation search results, i can see information about various ways to...



Now i have comprehensive information about writing data to Victoria Metrics Cloud.



documentation



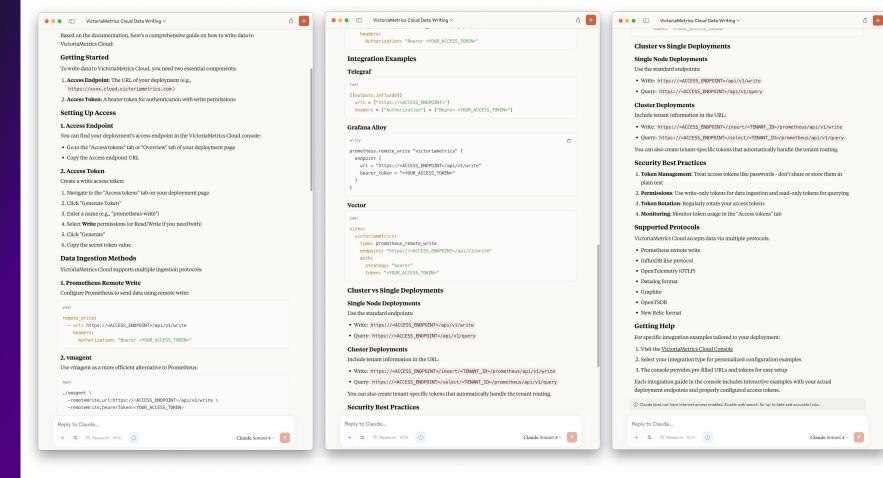
Great idea! Let me first get a complete list of completely unused metrics (0 queries), then we'll create relabeling rules to drop them.

- Thinking about writing metrics to Victoria Metrics Cloud platform
- **documentation**
- Based on the documentation search results, i can see information about various ways to...
- **documentation**
- Now i have comprehensive information about writing data to Victoria Metrics Cloud.
- **documentation**
- Synthesized comprehensive strategy for Victoria Metrics Cloud data writing.

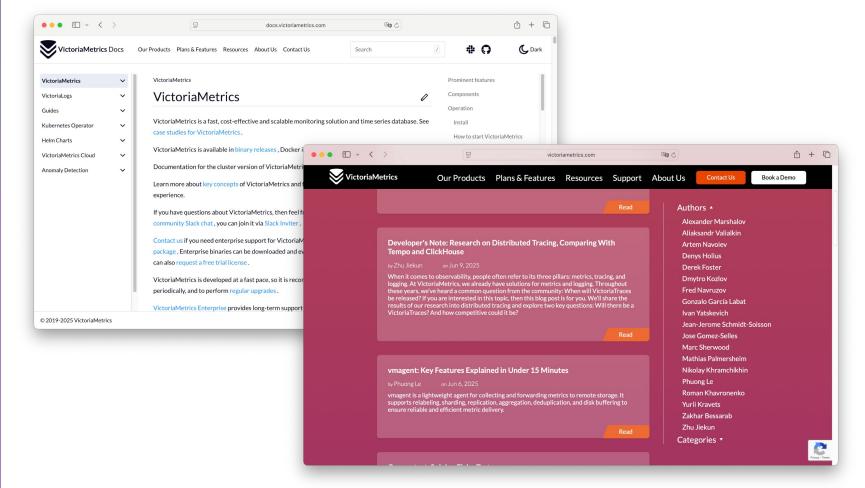
Based on the documentation, here's a comprehensive guide on how to write data to Victoria Metrics Cloud:

• • •



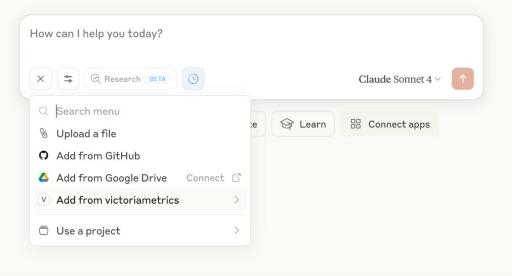






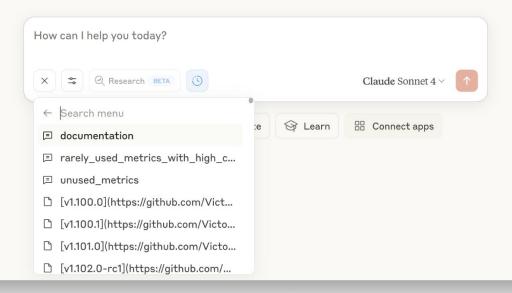


* What's new, Alexander?

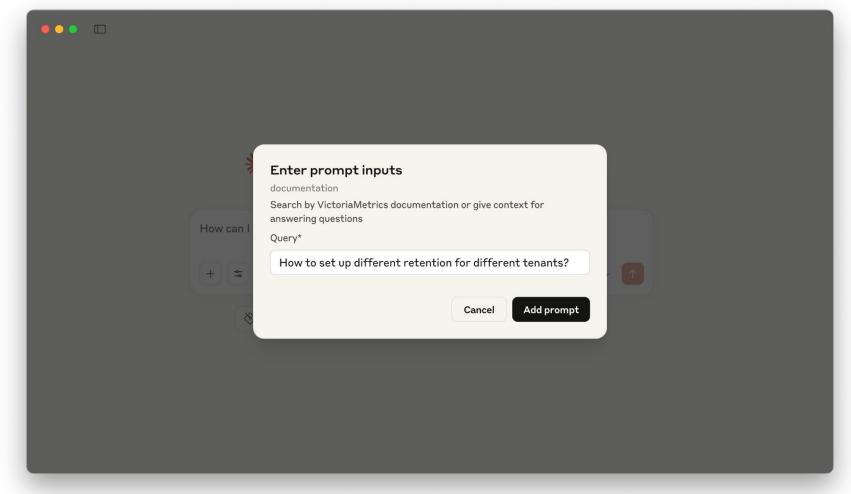




* What's new, Alexander?

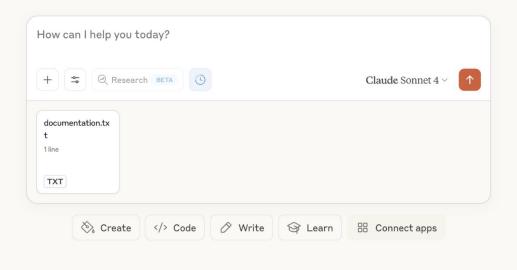








* What's new, Alexander?







Based on the VictoriaMetrics documentation, here's how to set up different retention periods for different tenants:

Enterprise Solution (Recommended)

VictoriaMetrics Enterprise supports tenant-specific retention through **retention filters**.

This is the most straightforward approach:

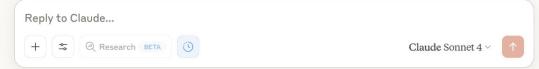
Configuration

Use the **-retentionFilter** command-line flag with tenant pseudo-labels:

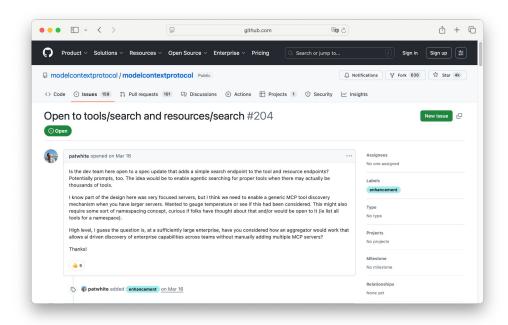
```
bash
-retentionFilter='{vm account id=~"42.*"}:1d' \
-retentionFilter='{vm account id="5",env="dev"}:5d' \
-retentionPeriod=4w
```

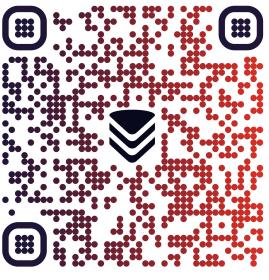
This configuration:

- Sets 1-day retention for tenants with accountID starting with "42"
- Sets 5-day retention for development data from tenant accountID=5
- Sets 4-week retention for all other tenants









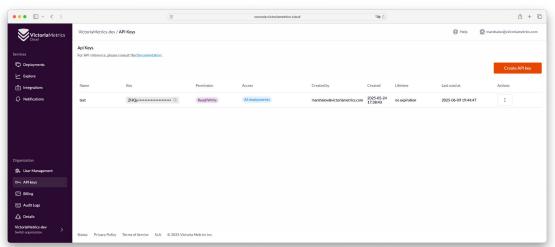
https://aithub.com/modelcontextprotocol/modelcontextprotocol/issues/204



Integration with VictoriaMetrics Cloud



mcp-victoriametrics + VictoriaMetrics Cloud









How to **set up** in your **client**?

```
Claude → claude_desktop_config.json
         "mcpServers": {
            "VictoriaMetrics": {
               "command": "/path/to/mcp-victoriametrics",
               "env": {
                  "VMC_API_KEY": "<YOUR_CLOUD_API_KEY>"
```



Additional cloud-specific tools

\otimes	deployments	available VM instances
\otimes	cloud_providers	supported cloud providers
\otimes	regions	supported regions
\odot	tiers	supported tier for VM instances
8	access_tokens	access tokens of VM instance
8	rule_filenames	alerting and recording rules filenames of VM instance
\otimes	rule_file	content of specific alerting or recording rules file of VM instance



Future plans



Roadmap

Support tool for query execution stats analysis

Add extra knowledge to the server

Implement multi-instance version

Support tools for VMAgent API

Support Alertmanager API

Support metrics metadata

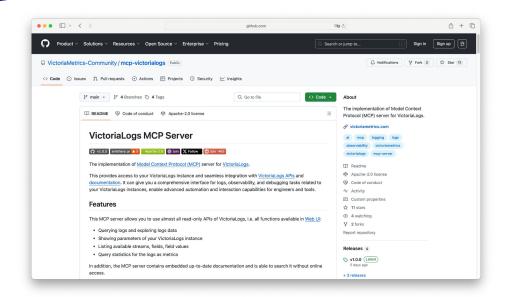
Improve tools for VMAlert API

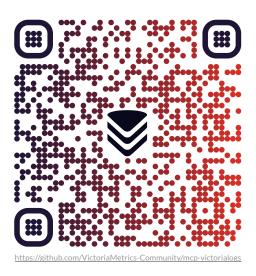


mcp-victorialogs



mcp-victorialogs







Available tools

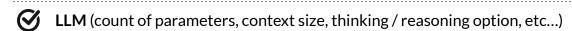
\otimes	facets	most frequent values per each log field
\otimes	field_names / field_values	list field names or values for specified query
\otimes	flags	non-default flags of the VictoriaLogs instance
\otimes	hits	number of matching log entries grouped by buckets
\otimes	query	query logs with LogsQL
\otimes	stats_query / stats_query_range	instant or range stats query (for alerting)
\otimes	stream_field_names / stream_field_values	list field names or values for specified log stream
\otimes	stream_ids / streams	list stream IDs or streams for specified query
\otimes	documentation	search in embedded VictoriaLogs documentation



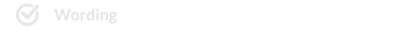
Before the end



Quality of responses strongly depends on:









Quality of responses strongly depends on:









Quality of responses strongly depends on:



Client (how well it works with the MCP tools)

Wording



- Power of LLMs in the keeping context, move step-by-step and combine different tools
- If you know the available capabilities (tools), terms and keywords and form queries based on them rather than using abstract queries the results will be more accurate
- Be tolerant of LLM "hallucinations", try different approaches to improve you results



- Power of LLMs in the keeping context, move step-by-step and combine different tools
- If you know the available capabilities (tools), terms and keywords and form queries based on them rather than using abstract queries the results will be more accurate
- Be tolerant of LLM "hallucinations", try different approaches to improve your results



- Power of LLMs in the keeping context, move step-by-step and combine different tools
- If you know the available capabilities (tools), terms and keywords and form queries based on them rather than using abstract queries the results will be more accurate
- Be tolerant of LLM "hallucinations", try different approaches to improve your results



- you should help LLMs with plan/context in abstract queries
- It works much accurate with known metrics from OpenSource software, because LLMs already has context about them in internal knowledge
- Combining the capabilities of different MCP servers from the same domain (for example mcp-victoriametrics + mcp-victorialogs + mcp for k8s) should enhance the capabilities and open up new complex scenarios



- you should help LLMs with plan/context in abstract queries
- It works much accurate with known metrics from OpenSource software, because LLMs already has context about them in internal knowledge
- Combining the capabilities of different MCP servers from the same domain (for example mcp-victoriametrics + mcp-victorialogs + mcp for k8s) should enhance the capabilities and open up new complex scenarios

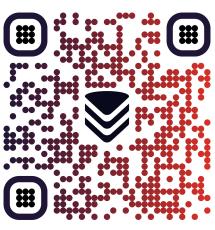


- you should help LLMs with plan/context in abstract queries
- It works much accurate with known metrics from OpenSource software because LLMs already has context about them in internal knowledge
- Combining the capabilities of different MCP servers from the same domain (for example mcp-victoriametrics + mcp-victorialogs + mcp for k8s) should enhance the capabilities and open up new complex scenarios



Can AI/ML Guide Observability?





https://www.voutube.com/watch?v=7t0ZMhX7OIY



VictoriaMetrics Cloud Thank you for listening

marshalov@victoriametrics.com

