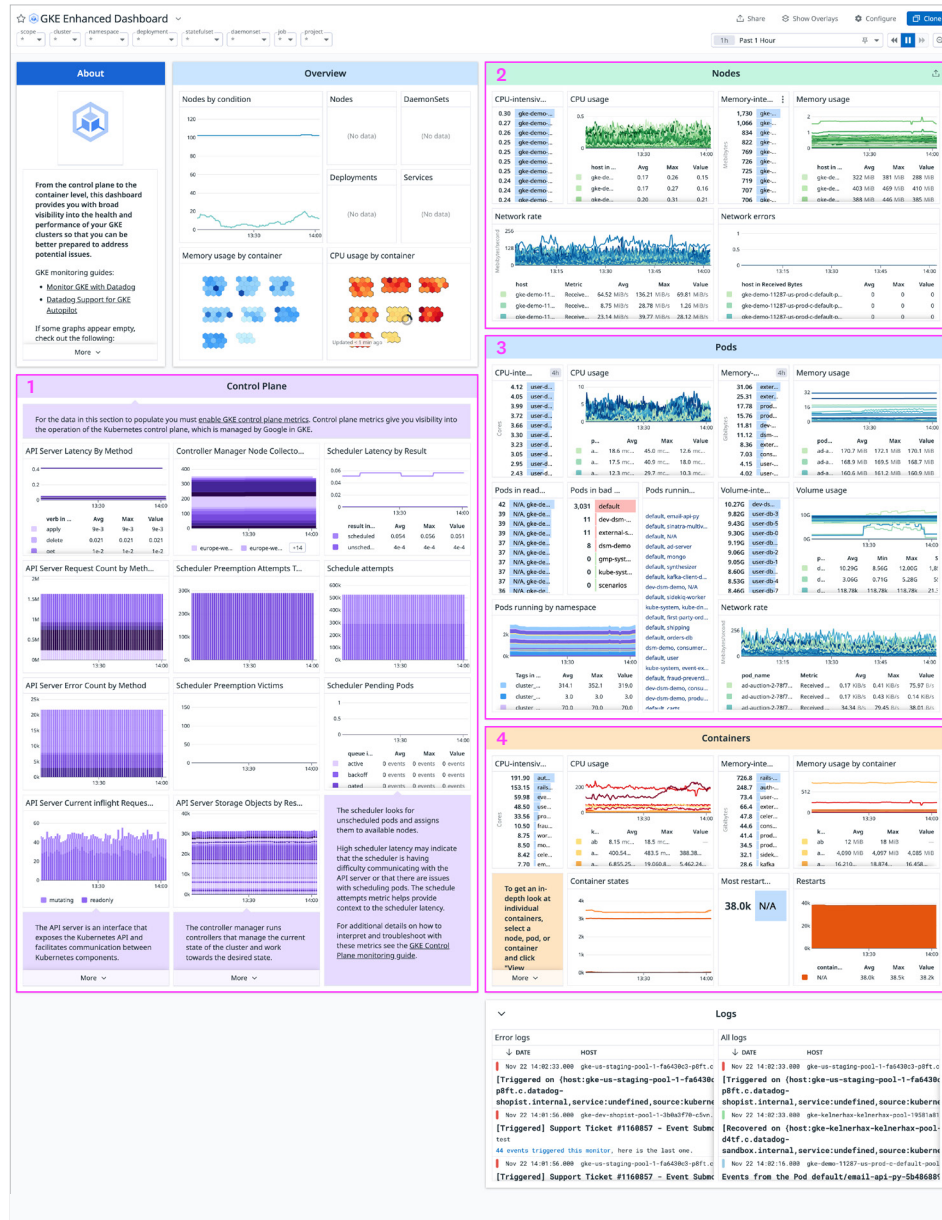# Cheatsheet: GKE Monitoring

Google Kubernetes Engine (GKE) is a managed Kubernetes service on Google Cloud that allows cluster operators to focus on running containerized applications without the overhead of managing the Kubernetes control plane.

With **Datadog's preconfigured dashboard**, you can get deep visibility into application performance data from across your cluster—such as CPU and memory usage, container and pod events, network throughput, and individual request traces—so you can be ready to tackle any issue you might encounter.

— Visualize the performance of your GKE containers and GKE control plane.
— Correlate the performance of your GKE containers with your applications.

Learn more about Datadog's GKE monitoring capabilities here.

## 1. Control plane metrics

| METRIC DESCRIPTION | METRIC NAME |
| --- | --- |
| Admission controller latency histogram in seconds, identified by name and broken out for each operation and API resource and type (validate or admit). *Shown as second* | gcp.gke.control_plane.apiserver.admission_controller_admission_duration_seconds (gauge) |
| Admission sub-step latency histogram in seconds, broken out for each operation and API resource and step type (validate or admit). *Shown as second* | gcp.gke.control_plane.apiserver.admission_step_admission_duration_seconds (gauge) |
| Admission webhook latency histogram in seconds, identified by name and broken out for each operation and API resource and type (validate or admit). *Shown as second* | gcp.gke.control_plane.apiserver.admission_webhook_admission_duration_seconds (gauge) |
| Maximal number of currently used inflight request limit of this apiserver per request kind. *Shown as request* | gcp.gke.control_plane.apiserver.current_inflight_requests (gauge) |
| Response latency distribution in seconds for each verb, dry run value, group, version, resource, subresource, scope and component. *Shown as second* | gcp.gke.control_plane.apiserver.request_duration_seconds (gauge) |
| Counter of apiserver requests broken out for each verb, dry run value, group, version, resource, scope, component, and HTTP response code. *Shown as request* | gcp.gke.control_plane.apiserver.request_total (gauge) |
| Response size distribution in bytes for each group, version, verb, resource, subresource, scope and component. *Shown as byte* | gcp.gke.control_plane.apiserver.response_sizes (gauge) |
| Number of stored objects at the time of last check split by kind. *Shown as object* | gcp.gke.control_plane.apiserver.storage_objects (gauge) |
| Number of Node evictions that happened since current instance of NodeController started. *Shown as event* | gcp.gke.control_plane.controller_manager.node_collector_evictions_number (count) |
| Number of pending pods, by the queue type. *Shown as event* | gcp.gke.control_plane.scheduler.pending_pods (gauge) |
| E2e latency for a pod being scheduled. *Shown as second* | gcp.gke.control_plane.scheduler.pod_scheduling_duration_seconds (gauge) |
| Total preemption attempts in the cluster till now. *Shown as attempt* | gcp.gke.control_plane.scheduler.preemption_attempts_total (count) |
| Number of selected preemption victims. *Shown as event* | gcp.gke.control_plane.scheduler.preemption_victims (gauge) |
| Scheduling attempt latency in seconds. *Shown as second* | gcp.gke.control_plane.scheduler.scheduling_attempt_duration_seconds (gauge) |
| Number of attempts to schedule pods. *Shown as attempt* | gcp.gke.control_plane.scheduler.schedule_attempts_total (gauge) |

## 2. Node metrics

| METRIC DESCRIPTION | METRIC NAME |
| --- | --- |
| Number of allocatable CPU cores on the node. *Shown as core* | gcp.gke.node.cpu.allocatable_cores (gauge) |
| Fraction of the allocatable CPU that is currently in use on the instance. *Shown as fraction* | gcp.gke.node.cpu.allocatable_utilization (gauge) |
| Cumulative CPU usage on all cores used on the node. *Shown as second* | gcp.gke.node.cpu.core_usage_time (count) |
| Total number of CPU cores on the node. *Shown as core* | gcp.gke.node.cpu.total_cores (gauge) |
| Local ephemeral storage bytes allocatable on the node. *Shown as byte* | gcp.gke.node.ephemeral_storage.allocatable_bytes (gauge) |
| Free number of inodes on local ephemeral storage. | gcp.gke.node.ephemeral_storage.inodes_free (gauge) |
| Total number of inodes on local ephemeral storage. | gcp.gke.node.ephemeral_storage.inodes_total (gauge) |
| Total ephemeral storage bytes on the node. Shown as byte | gcp.gke.node.ephemeral_storage.total_bytes (gauge) |

## 2. Node metrics cont'd

| METRIC DESCRIPTION | METRIC NAME |
| --- | --- |
| Local ephemeral storage bytes used by the node. *Shown as byte* | gcp.gke.node.ephemeral_storage.used_bytes (gauge) |
| Cumulative memory bytes used by the node. *Shown as byte* | gcp.gke.node.memory.allocatable_bytes (gauge) |
| Fraction of the allocatable memory that is currently in use on the instance. *Shown as fraction* | gcp.gke.node.memory.allocatable_utilization (gauge) |
| Number of bytes of memory allocatable on the node. *Shown as byte* | gcp.gke.node.memory.total_bytes (gauge) |
| Cumulative memory bytes used by the node. *Shown as byte* | gcp.gke.node.memory.used_bytes (gauge) |
| Cumulative number of bytes received by the node over the network. *Shown as byte* | gcp.gke.node.network.received_bytes_count (count) |
| Cumulative number of bytes transmitted by the node over the network. *Shown as byte* | gcp.gke.node.network.sent_bytes_count (count) |
| Max PID of OS on the node. | gcp.gke.node.pid_limit (gauge) |
| Number of running process in the OS on the node. | gcp.gke.node.pid_used (gauge) |
| Cumulative CPU usage on all cores used by the node level system daemon. *Shown as second* | gcp.gke.node_daemon.cpu.core_usage_time (count) |
| Memory usage by the system daemon. *Shown as byte* | gcp.gke.node_daemon.memory.used_bytes (gauge) |

## 3. Pod metrics

| METRIC DESCRIPTION | METRIC NAME |
| --- | --- |
| Cumulative number of bytes received by the pod over the network. *Shown as byte* | gcp.gke.pod.network.received_bytes_count (count) |
| Cumulative number of bytes transmitted by the pod over the network. *Shown as byte* | gcp.gke.pod.network.sent_bytes_count (count) |
| Total number of disk bytes available to the pod. *Shown as byte* | gcp.gke.pod.volume.total_bytes (gauge) |
| Number of disk bytes used by the pod. *Shown as byte* | gcp.gke.pod.volume.used_bytes (gauge) |
| Fraction of the volume that is currently being used by the instance. *Shown as fraction* | gcp.gke.pod.volume.utilization (gauge) |

## 4. Container metrics

| METRIC DESCRIPTION | METRIC NAME |
| --- | --- |
| Percent of time over the past sample period during which the accelerator was actively processing. *Shown as percent* | gcp.gke.container.accelerator.duty_cycle (gauge) |
| Total accelerator memory. *Shown as byte* | gcp.gke.container.accelerator.memory_total (gauge) |
| Total accelerator memory allocated. *Shown as byte* | gcp.gke.container.accelerator.memory_used (gauge) |
| Number of accelerator devices requested by the container. *Shown as device* | gcp.gke.container.accelerator.request (gauge) |
| Cumulative CPU usage on all cores used by the container. *Shown as second* | gcp.gke.container.cpu.core_usage_time (count) |
| CPU cores limit of the container. *Shown as core* | gcp.gke.container.cpu.limit_cores (gauge) |
| Fraction of the CPU limit that is currently in use on the instance. *Shown as fraction* | gcp.gke.container.cpu.limit_utilization(gauge) |
| Number of CPU cores requested by the container. *Shown as core* | gcp.gke.container.cpu.request_cores (gauge) |
| Fraction of the requested CPU that is currently in use on the instance. *Shown as fraction* | gcp.gke.container.cpu.request_utilization (gauge) |
| Local ephemeral storage limit. *Shown as byte* | gcp.gke.container.ephemeral_storage.limit_bytes (gauge) |
| Local ephemeral storage request. *Shown as byte* | gcp.gke.container.ephemeral_storage.request_bytes (gauge) |
| Local ephemeral storage usage. *Shown as byte* | gcp.gke.container.ephemeral_storage.used_bytes (gauge) |
| Memory limit of the container. *Shown as byte* | gcp.gke.container.memory.limit_bytes (gauge) |
| Fraction of the memory limit that is currently in use on the instance. *Shown as fraction* | gcp.gke.container.memory.limit_utilization (gauge) |
| Number of page faults, broken down by type. *Shown as fault* | gcp.gke.container.memory.page_fault_count (count) |
| Memory request of the container. *Shown as byte* | gcp.gke.container.memory.request_bytes (gauge) |
| Fraction of the requested memory that is currently in use on the instance. *Shown as fraction* | gcp.gke.container.memory.request_utilization (gauge) |
| Memory usage of the container. *Shown as byte* | gcp.gke.container.memory.used_bytes (gauge) |
| Number of times the container has restarted. *Shown as occurrence* | gcp.gke.container.restart_count (count) |
| Time in seconds that the container has been running. *Shown as second* | gcp.gke.container.uptime (gauge) |