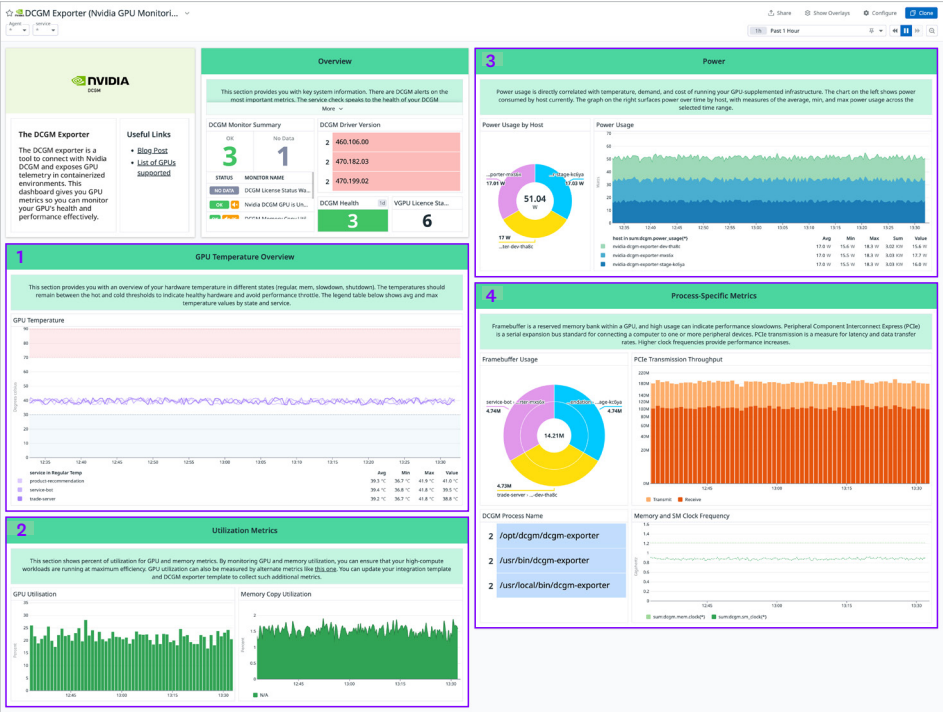


Cheatsheet: Nvidia DCGM



Why Datadog?

Datadog is a leading SaaS-based observability and security platform that brings together telemetry from across your tech environment—including infrastructure metrics, application traces, and logs—together in a single platform. Our monitoring capabilities include customizable alerting, monitoring reports, and visualization tools like out-of-the-box dashboards, making it easy and fast to investigate and resolve issues. **With 700+ vendor-backed integrations, including Nvidia DCGM, Triton, Jetson and NVML, you can gain visibility into any part of your tech stack.**

1. GPU Temperature Overview (DCGM)

— This section provides you with an overview of your hardware temperature in different states (regular, mem, slowdown, shutdown). The temperatures should remain between the hot and cold thresholds to indicate healthy hardware and avoid performance throttle. The legend table below shows avg and max temperature values by state and service.

METRIC DESCRIPTION	NVIDIA NAME	DATADOG NAME
GPU temperature (in C)	dcgm_fi_dev_gpu_temp	dcgm.temperature

2. Utilization Metrics

— This section shows percent of utilization for GPU and memory metrics. By monitoring GPU and memory utilization, you can ensure that your high-compute workloads are running at maximum efficiency. GPU utilization can also be measured by alternate metrics like this one. You can update your integration template and DCGM exporter template to collect such additional metrics.

METRIC DESCRIPTION	NVIDIA NAME	DATADOG NAME
GPU Utilization in Percentage	dcgm_fi_dev_gpu_util	dcgm.gpu_utilization
Memory Copy Utilization in Percentage	dcgm_fi_dev_mem_copy_util	dcgm.mem.copy_utilization

3. Power Metrics

— Power usage is directly correlated with temperature, demand, and cost of running your GPU-supplemented infrastructure. The chart on the left shows power consumed by host currently. The graph on the right surfaces power over time by host, with measures of the average, min, and max power usage across the selected time range.

METRIC DESCRIPTION	NVIDIA NAME	DATADOG NAME
Power draw (in W)	dcgm_fi_dev_power_usage	dcgm.power_usage

4. Process-Specific Metrics

— Framebuffer is a reserved memory bank within a GPU, and high usage can indicate performance slowdowns. Peripheral Component Interconnect Express (PCIe) is a serial expansion bus standard for connecting a computer to one or more peripheral devices. PCIe transmission is a measure for latency and data transfer rates. Higher clock frequencies provide performance increases.

METRIC DESCRIPTION	NVIDIA NAME	DATADOG NAME
Framebuffer Usage, memory used in MB	dcgm_fi_dev_fb_used	dcgm.fb_used
Framebuffer Usage, memory free in MB	dcgm_fi_dev_fb_free	dcgm.fb_free
PCIe Tx utilization information	dcgm_fi_dev_pcie_tx_throughput	dcgm.pcie_tx_throughput
PCIe Rx utilization information	dcgm_fi_dev_pcie_rx_throughput	dcgm.pcie_rx_throughput
SM Clock Frequency. SM clock frequency in MHz	dcgm_fi_dev_sm_clock	dcgm.sm_clock
Memory Clock Frequency. SM clock frequency in MHz	dcgm_fi_dev_mem_clock	dcgm.mem.clock
vGPU License status	dcgm_vgpu_license_status	dcgm.vgpu_license_status

Useful Links

- [Blog Post](#)
- [List of GPUs supported](#)

Cheatsheet: Nvidia Triton



1. Inference Metrics

— Counts: The following metrics show Triton server inference numbers, such as failure, success, and pending, along with the total count and batch execution.

METRIC DESCRIPTION	NVIDIA NAME	DATADOG NAME
Number of successful inference requests, all batch sizes	nv_inference_request_success	nvidia_triton.inference.request_success.count
Number of failed inference requests, all batch sizes	nv_inference_request_failure	nvidia_triton.inference.request_failure.count
Number of inferences performed (does not include cached requests)	nv_inference_count	nvidia_triton.inference.count.count
Number of model executions performed (does not include cached requests)	nv_inference_exec_count	nvidia_triton.inference.exec.count.count
Instantaneous number of pending requests awaiting execution per-model	nv_inference_pending_request_count	nvidia_triton.inference.pending.request.count

2. Latency Metrics

— The following metrics show Triton server latency durations, such as request, queue, and computation.

METRIC DESCRIPTION	NVIDIA NAME	DATADOG NAME
Cumulative inference request duration in microseconds (includes cached requests)	nv_inference_request_duration_us	nvidia_triton.inference.request.duration_us.count
Cumulative inference queuing duration in microseconds (includes cached requests)	nv_inference_queue_duration_us	nvidia_triton.inference.queue.duration_us.count
Cumulative compute input duration in microseconds (does not include cached requests)	nv_inference_compute_input_duration_us	nvidia_triton.inference.compute.input.duration_us.count
Cumulative compute inference duration in microseconds (does not include cached requests)	nv_inference_compute_infer_duration_us	nvidia_triton.inference.compute.infer.duration_us.count
Cumulative inference compute output duration in microseconds (does not include cached requests)	nv_inference_compute_output_duration_us	nvidia_triton.inference.compute.output.duration_us.count

3. GPU Utilization Metrics

— GPU memory usage for your Triton server shows up here grouped by the GPU UID and host.

METRIC DESCRIPTION	NVIDIA NAME	DATADOG NAME
GPU instantaneous power, in watts	nv_gpu_power_usage	nvidia_triton.gpu.power.usage
GPU energy consumption since count in Joules	nv_energy_consumption	nvidia_triton.energy.consumption.count
GPU utilization rate (0.0 - 1.0)	nv_gpu_utilization	nvidia_triton.gpu.utilization
GPU power management limit in watts	nv_gpu_power_limit	nvidia_triton.gpu.power.limit

4. GPU Memory Metrics

— These metrics show the GPU power on your Triton server and the actual usage compared to the power limits.

METRIC DESCRIPTION	NVIDIA NAME	DATADOG NAME
Total GPU memory, in bytes	nv_gpu_memory_total_bytes	nvidia_triton.gpu.memory.total_bytes
Used GPU memory, in bytes	nv_gpu_memory_used_bytes	nvidia_triton.gpu.memory.used_bytes

Want to get started with Datadog? Click [here](#) for a free, 14-day trial, and [read more](#) about our product offerings.

Cheatsheet: Nvidia Triton



5. CPU Memory Metrics

— CPU Metrics are only supported on Linux. They collect information from the /proc filesystem.

METRIC DESCRIPTION	NVIDIA NAME	DATADOG NAME
CPU utilization rate [0.0 - 1.0]	nv_cpu_utilization	nvidia_triton.cpu.utilization
CPU used memory (RAM), in bytes	nv_cpu_memory_total_bytes	nvidia_triton.cpu.memory.used_bytes
GPU total memory, in bytes	nv_cpu_memory_used_bytes	nvidia_triton.gpu.memory.total_bytes

6. Cache Metrics

— These are the cache metrics reported directly by Triton, such as the cache hit/miss counts and durations.

METRIC DESCRIPTION	NVIDIA NAME	DATADOG NAME
Cache utilization [0.0 - 1.0]	nv_cache_util	nvidia_triton.cache.util
Number of responses stored in response cache	nv_cache_num_entries	nvidia_triton.cache.num.entries
Number of cache lookups in response cache	nv_cache_num_lookups	nvidia_triton.cache.num.lookups
Number of cache hits in response cache	nv_cache_num_hits	nvidia_triton.cache.num.hits
Number of cache misses in response cache	nv_cache_num_misses	nvidia_triton.cache.num.misses
Number of cache evictions in response cache	nv_cache_num_evictions	nvidia_triton.cache.num.evictions
Total cache lookup duration (hit and miss), in microseconds	nv_cache_lookup_duration	nvidia_triton.cache.lookup.duration
Total cache insertion duration, in microseconds	nv_cache_insertion_duration	nvidia_triton.cache.insertion.duration

Useful Links

- [Nvidia Triton Integration doc](#)
- [AI Stack Integrations blog](#)
- [Nvidia Triton metric collection](#)

Want to get started with Datadog? Click [here](#) for a free, 14-day trial, and [read more](#) about our product offerings.

Real-Time Nvidia GPU Monitoring

Track the performance of all your GPU workloads, regardless of whether they are containerized, hosted locally, or deployed in the cloud. Correlate GPU performance and usage with other technologies that support AI, including large language models use cases.

[TRY DATADOG FOR FREE](#)

