# Monitoring COMPaaS

Timothy Bargo

2022-04-29

## Background

The Electronic Visualization Laboratory at the University of Illinois Chicago acquired a 24 compute node, 64 GPU composable infrastructure compute cluster, COMPaaS DLV - Composable Platform as a Service: Instrument for Deep Learning & Visualization in 2019. Since then 40+ users have designed and run computations on COMPaaS. To track cluster utilization, evaluate hardware performance, and simply analyze application resource utilization of this experimental composable infrastructure system we set out to collect metrics on the hardware and software infrastructure. Users deploy their applications on the cluster using Kubernetes, a container orchestration platform that enables users to deploy applications on the cluster without the need to have complete knowledge of the underlying hardware and with minimal assistance from system administrators. Following this, we utilized tooling already adapted for Kubernetes based clusters to collect hardware and software metrics. The metrics provided by this tooling is useful for system administrators in evaluating cluster performance and utilization, as well as useful for users in determining their application's performance and discovering bottlenecks.

## How data is collected

Our data sources are amongst host machine metrics, Kubernetes metrics, and NVIDIA GPU metrics. The data is aggregated into Prometheus, and displayed in Grafana.
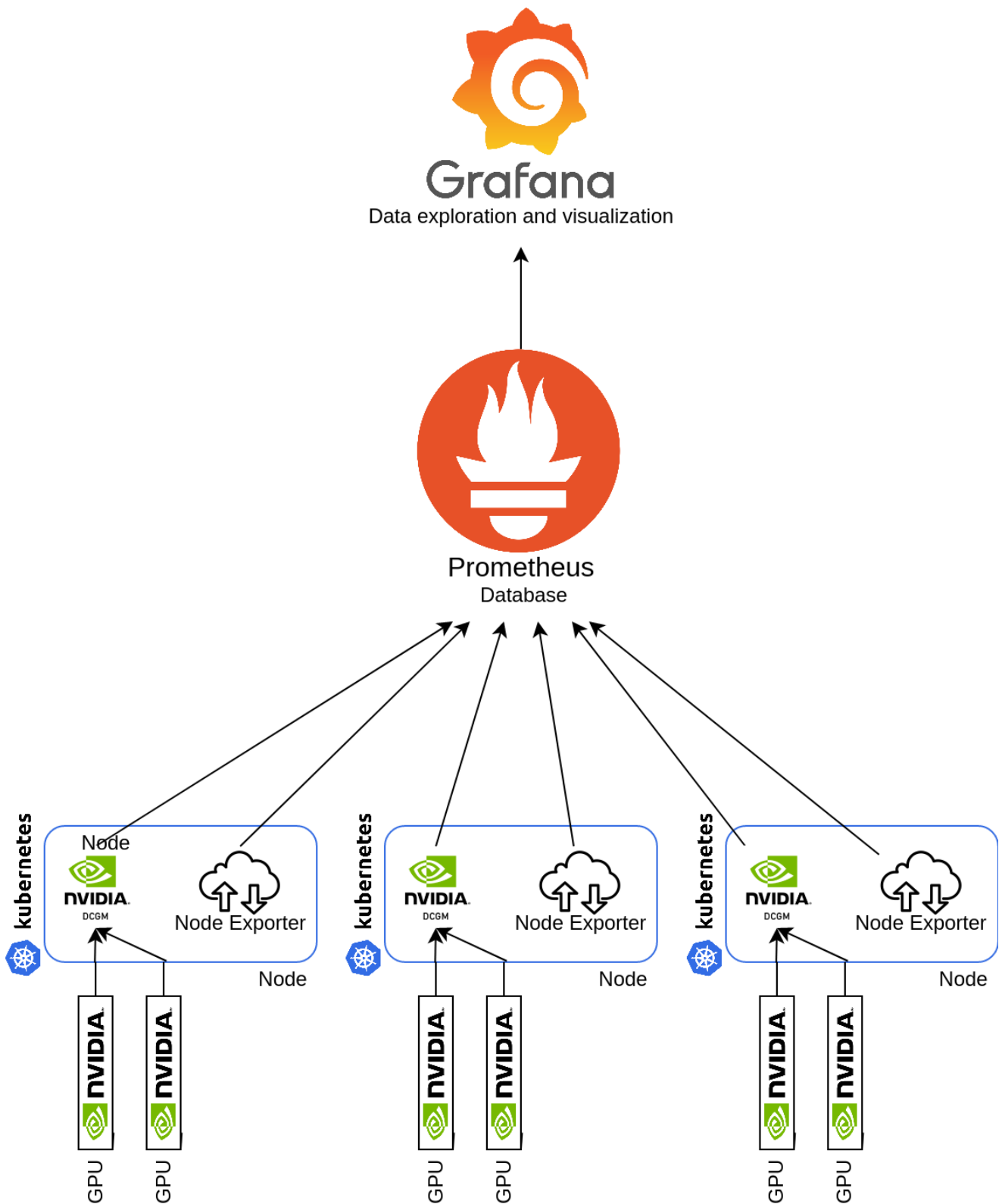
For host machine metrics we leverage the Node Exporter service. Node Exporter collects machine metrics such as overall CPU, memory, network, storage utilization, and more. Node Exporter is deployed on all machines in the cluster by a Kubernetes deployment and automatically periodically scrapes the machines for metrics data.

For Kubernetes metrics, we use the built-in Kubernetes metrics API. This provides information such as what pods are running, the CPU utilization of each pod, Kubernetes container network utilization, and more. For NVIDIA GPU metrics, we use the NVIDIA Data Center GPU Manager (DCGM). The NVIDIA DCGM is deployed on all nodes by a Kubernetes deployment and collects information such as GPU processor, memory, power utilization, and more.

Prometheus is a monitoring tool that provides a time-series database. Prometheus pulls data from each of the individual data sources and stores the data in an aggregated form in a central database.

Grafana is an analytics and interactive visualization web application and primarily serves to present time-series data. Grafana can be configured to pull data from a number of sources including the Prometheus database for analysis and presentation.

We deployed the above components of the monitoring stack with an existing package built for monitoring Kubernetes clusters, see here https://github.com/prometheus-operator/kube-prometheus.

This diagram shows where each service is deployed. NVIDIA GPUs are attached to machines. The machines run Kubernetes. NVIDIA DCGM and Node Exporter run inside of Kubernetes. Prometheus pulls data from the Kubernetes node, NVIDIA DCGM, and Node Exporter to store in a database. Grafana pulls data from Prometheus for use in analysis.

# Application monitoring

To acquire an understanding of the types of metrics we can analyze, we will walk through graphed metrics of an individual user's application: This user is using a single machine with 72 CPU cores, 376 GB of memory, and 7x NVIDIA Tesla V100 GPUs. The user's application is designed to use 3x GPUs and has been running for over 5 days.
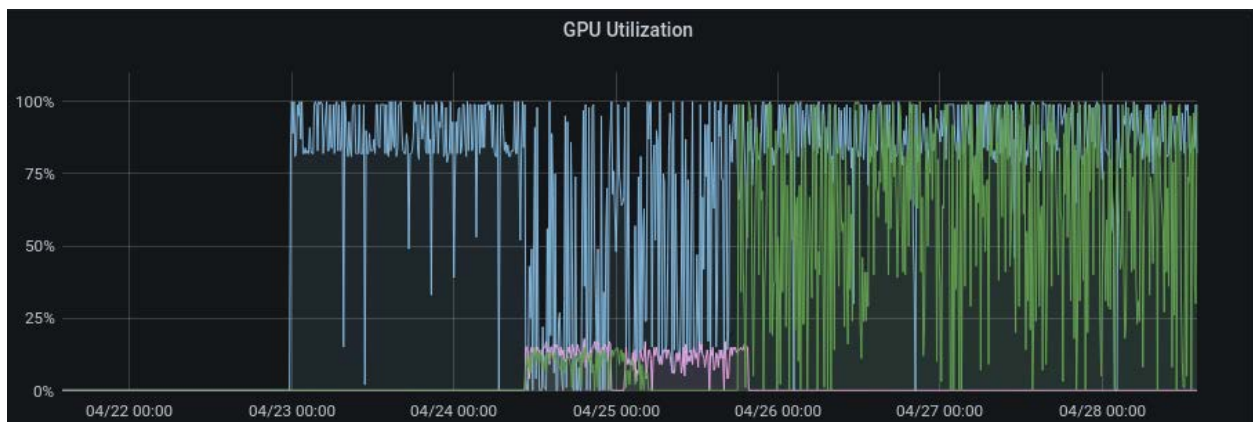
## Using Grafana to Analyze Metrics

As described before, we will use Grafana to analyze the metrics we have collected. To generate a graph in Grafana we first need to create a new dashboard and add a new panel. Once in the edit panel window we select the data source, in this case Prometheus, and generate a query. The query asks for the particular data we are interested in. If you click on the "Metrics" drop down within the query Grafana will show all data available for the query.
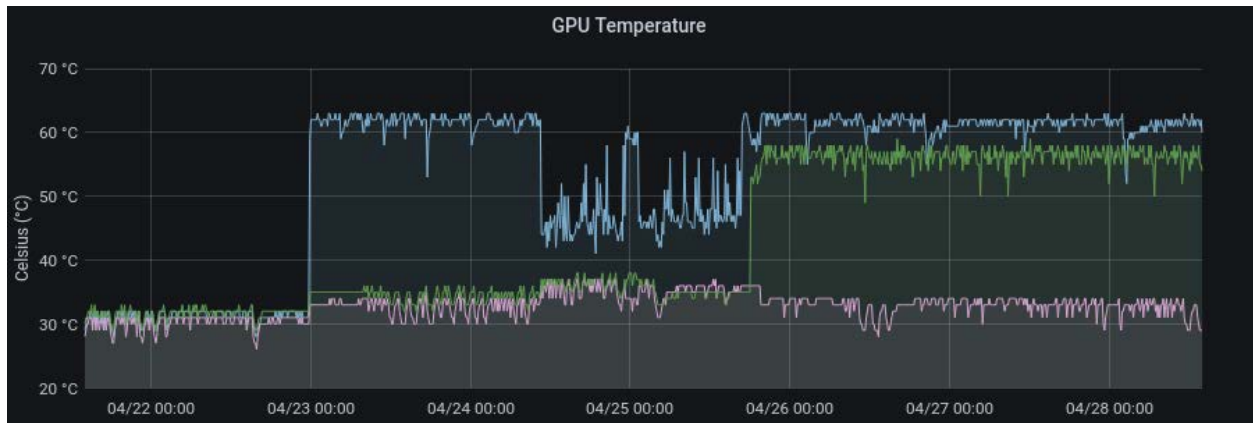
If we enter "DCGM_FI_DEV_GPU_UTIL" for the query value we will see a graph generated showing GPU utilization. Use the right hand pane to set the visualization type, labels, axes, legends, and set data overrides.

## Analyzing Metrics with Grafana

We will start with GPU metrics analyzing one week of data. Here we see the GPU processor utilization. The three GPUs are denoted by color, one blue, pink, and green. These colors will remain constant throughout our analysis of the GPU performance. First only one GPU is active, the blue GPU. Around midday on 4/24 the pink GPU and green GPUs become active with low utilization, and drop the utilization of the blue GPU. After midday of 4/25 the pink GPU stops its activity and the green and blue GPUs increase to 100% utilization.

We see that GPU temperature follows the same pattern as the GPU utilization where low utilization shows a temperature of 30C and high utilization at 60C.



GPU power follows the same pattern where low utilization shows power consumption of 50W and high utilization reaching as high as 250W.

The application started on 4/23 at 00:00 and used 80%-100% of the blue GPU processor. Mid-day of 4/24 the pink GPU becomes active, however at only 15% and the blue GPU's utilization decreases to 50%-70%. However, CPU load drastically increases at the same time and memory utilization steadily climbs. Performance profiling showing days-long performance impacts may be useful for a user in determining bottlenecks in their application and useful for a system administrator in understanding how their cluster resources are utilized.

## Cluster-wide metrics

Now that we are familiar with our basic measurements, we can look at cluster-wide metrics. To guide us through the metrics we will set a goal to attempt to determine overall utilization of all GPUs in the cluster.

This is a graph of the power consumption of all GPUs in the cluster that are known by Kubernetes. There are three types of GPUs in the cluster, NVIDIA Ampere A100, NVIDIA Tesla V100 and NVIDIA Tesla T4. With this graph we can see how much less power the Tesla T4 consumes. The peaks of Tesla T4 GPUs are circled in red, peaking at around 100W of power. The larger peaks ranging from 200W to 250W are Tesla V100s GPUs. The Ampere A100 GPU is not significant enough to see in this graph.



To further show the power difference, we can manipulate the data in Grafana to group each GPU type's power consumption.

Here we can see the V100 GPUs using twice as much power as the T4 GPUs and the single A100 GPU using a comparatively negligible amount of power at 35W. However, it is important to note that not all GPUs are used to their full capacity and this power difference can change.

When we look at the graph for GPU utilization grouped by GPU type we can see the correlation between utilization and power consumption.



If we return to the power graph and mark a time when GPU utilization was 0% cluster-wide, all GPUs were idle, we can sum the individual power values and see that the idle power consumption of all GPUs was 1.1 kW.

We can also use Grafana to count the number of GPUs attached to a machine and detected by Kubernetes. For the V100 GPUs, the graph below the value "21" drops from 22 to 21, noting that a GPU was removed from the system, a feature of this composable cluster.



Looking at power consumption cluster-wide, that is power consumption of all GPUs. We can see that over the last 7 days our peak power consumption, just for GPUs, was around 2 kiloWatts. If we subtract the 1.1 kW of idle power we find that peak cluster power consumption variance is 900W and is equivalent to 4x V100 GPUs at 100% utilization.



With this we can roughly estimate our overall cluster utilization. If we assume that at peak load a V100 GPU uses 225W and a T4 GPU uses 100W, with 21 V100 GPUs and 30 T4 GPUs attached to machines and available in Kubernetes we would expect a power consumption of around 7.7kW. If we remove the 1.1kW of idle power from 7.7kW, as idle power is a constant, our power variance is 6.6kW.

Then over a 7 day period our peak power variance of all GPUs was 900W.

If we divide 900W by 6600W we get a power consumption load of ~14% of peak. It is possible to extrapolate the 14% as a *very* rough estimate of overall GPU utilization cluster-wide.

## Conclusion

We walked through a small portion of the variety of metrics available to a system administrator through a Kubernetes cluster monitoring stack composed of Node Exporter, Prometheus with Kubernetes plugin, and NVIDIA DCGM. We demonstrated the effectiveness of cluster metrics through attempting to find an estimate of cluster utilization of GPU resources. Exploring the cluster metrics with a variety of graphs and measurements we concluded a very low overall utilization of GPU resources at 14%. A system administrator may consider this value in determining the effectiveness of their jobs scheduler or note that the cluster is far under capacity and can support more users. Metrics serve as a powerful tool to determine hardware and software performance and utilization. The information can be employed to determine application bottlenecks, understand the best machine composition for an application, and discover a variety of narratives in cluster performance.

## References

https://compaas.evl.uic.edu/

https://github.com/prometheus-operator/kube-prometheus

https://github.com/prometheus/node_exporter

https://grafana.com/

https://kubernetes.io/

https://catalog.ngc.nvidia.com/orgs/nvidia/teams/k8s/containers/dcgm-exporter

https://developer.nvidia.com/blog/monitoring-gpus-in-kubernetes-with-dcgm/

https://prometheus.io/

## Appendix

This is a JSON export of the Grafana dashboard used to generate the graphs shown in this paper.

```
{
  "__inputs": [
      {
      "name": "DS_PROMETHEUS",
      "label": "Prometheus",
      "description": "",
      "type": "datasource",
      "pluginId": "prometheus",
      "pluginName": "Prometheus"
      }
  ],
  "__requires": [
      {
      "type": "grafana",
      "id": "grafana",
      "name": "Grafana",
      "version": "7.4.5"
      },
      {
      "type": "panel",
      "id": "graph",
      "name": "Graph",
      "version": ""
      },
      {
      "type": "datasource",
      "id": "prometheus",
      "name": "Prometheus",
      "version": "1.0.0"
      },
      {
      "type": "panel",
      "id": "stat",
      "name": "Stat",
      "version": ""
      }
  ],
  "annotations": {
      "list": [
      {
      "builtIn": 1,
      "datasource": "-- Grafana --",
```

```
    "enable": true,
    "hide": true,
    "iconColor": "rgba(0, 211, 255, 1)",
    "name": "Annotations & Alerts",
    "type": "dashboard"
    }
    ]
},
"editable": true,
"gnetId": null,
"graphTooltip": 0,
"id": null,
"links": [],
"panels": [
    {
    "datasource": "${DS_PROMETHEUS}",
    "fieldConfig": {
    "defaults": {
        "color": {
        "mode": "thresholds"
        },
        "custom": {},
        "mappings": [],
        "thresholds": {
        "mode": "absolute",
        "steps": []
        }
    },
    "overrides": []
    },
    "gridPos": {
    "h": 8,
    "w": 12,
    "x": 0,
    "y": 0
    },
    "id": 30,
    "options": {
    "colorMode": "value",
    "graphMode": "area",
    "justifyMode": "auto",
    "orientation": "auto",
    "reduceOptions": {
        "calcs": [
        "lastNotNull"
        ],
```

```
        "fields": "",
        "values": false
    },
    "text": {},
    "textMode": "auto"
    },
    "pluginVersion": "7.4.5",
    "targets": [
    {
        "expr": "count(DCGM_FI_DEV_GPU_UTIL{modelName=\"NVIDIA
A100-PCIE-40GB\"})",
        "interval": "",
        "legendFormat": "A100",
        "refId": "A"
    },
    {
        "expr": "count(DCGM_FI_DEV_GPU_UTIL{modelName=\"Tesla T4\"})",
        "hide": false,
        "interval": "",
        "legendFormat": "T4",
        "refId": "B"
    },
    {
        "expr": "count(DCGM_FI_DEV_GPU_UTIL{modelName=\"Tesla
V100-PCIE-32GB\"})",
        "hide": false,
        "interval": "",
        "legendFormat": "V100",
        "refId": "C"
    }
    ],
    "timeFrom": null,
    "timeShift": null,
    "title": "Number of Detected GPUs by Type",
    "type": "stat"
    },
    {
    "aliasColors": {},
    "bars": false,
    "dashLength": 10,
    "dashes": false,
    "datasource": "${DS_PROMETHEUS}",
    "fieldConfig": {
    "defaults": {
        "custom": {}
    },
```

```
      "overrides": []
    },
    "fill": 1,
    "fillGradient": 0,
    "gridPos": {
    "h": 9,
    "w": 12,
    "x": 12,
    "y": 0
    },
    "hiddenSeries": false,
    "id": 24,
    "legend": {
    "avg": false,
    "current": false,
    "max": false,
    "min": false,
    "show": true,
    "total": false,
    "values": false
    },
    "lines": true,
    "linewidth": 1,
    "nullPointMode": "null",
    "options": {
    "alertThreshold": true
    },
    "percentage": false,
    "pluginVersion": "7.4.5",
    "pointradius": 2,
    "points": false,
    "renderer": "flot",
    "seriesOverrides": [],
    "spaceLength": 10,
    "stack": false,
    "steppedLine": false,
    "targets": [
    {
          "expr": "sum(DCGM_FI_DEV_POWER_USAGE{modelName=\"Tesla
V100-PCIE-32GB\"})",
          "interval": "",
          "legendFormat": "",
          "refId": "A"
    }
    ],
    "thresholds": [],
```

```
"timeFrom": null,
"timeRegions": [],
"timeShift": null,
"title": "V100 Power Consumption Watt (W)",
"tooltip": {
"shared": true,
"sort": 0,
"value_type": "individual"
},
"type": "graph",
"xaxis": {
"buckets": null,
"mode": "time",
"name": null,
"show": true,
"values": []
},
"yaxes": [
{
        "format": "watt",
        "label": null,
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
},
{
        "format": "short",
        "label": null,
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
}
],
"yaxis": {
"align": false,
"alignLevel": null
}
},
{
"aliasColors": {},
"bars": false,
"dashLength": 10,
"dashes": false,
"datasource": "${DS_PROMETHEUS}",
```

```json
"fieldConfig": {
"defaults": {
    "custom": {}
},
"overrides": []
},
"fill": 1,
"fillGradient": 0,
"gridPos": {
"h": 8,
"w": 12,
"x": 0,
"y": 8
},
"hiddenSeries": false,
"id": 28,
"legend": {
"avg": false,
"current": false,
"max": false,
"min": false,
"show": true,
"total": false,
"values": false
},
"lines": true,
"linewidth": 1,
"nullPointMode": "null",
"options": {
"alertThreshold": true
},
"percentage": false,
"pluginVersion": "7.4.5",
"pointradius": 2,
"points": false,
"renderer": "flot",
"seriesOverrides": [],
"spaceLength": 10,
"stack": false,
"steppedLine": false,
"targets": [
{
    "expr": "sum(DCGM_FI_DEV_GPU_UTIL{modelName=\"NVIDIA
A100-PCIE-40GB\"})",
    "interval": "",
    "legendFormat": "",
```

```json
                "refId": "A"
        },
        {
                "expr": "sum(DCGM_FI_DEV_GPU_UTIL{modelName=\"Tesla T4\"})",
                "hide": false,
                "interval": "",
                "legendFormat": "",
                "refId": "B"
        },
        {
                "expr": "sum(DCGM_FI_DEV_GPU_UTIL{modelName=\"Tesla
V100-PCIE-32GB\"})",
                "hide": false,
                "interval": "",
                "legendFormat": "",
                "refId": "C"
        }
        ],
        "thresholds": [],
        "timeRegions": [],
        "title": "GPU Utilization by Type",
        "tooltip": {
        "shared": true,
        "sort": 0,
        "value_type": "individual"
        },
        "type": "graph",
        "xaxis": {
        "buckets": null,
        "mode": "time",
        "name": null,
        "show": true,
        "values": []
        },
        "yaxes": [
        {
                "format": "percent",
                "label": null,
                "logBase": 1,
                "max": null,
                "min": null,
                "show": true
        },
        {
                "format": "short",
                "label": null,
```

```json
            "logBase": 1,
            "max": null,
            "min": null,
            "show": true
    }
    ],
    "yaxis": {
    "align": false,
    "alignLevel": null
    }
    },
    {
    "aliasColors": {},
    "bars": false,
    "dashLength": 10,
    "dashes": false,
    "datasource": "${DS_PROMETHEUS}",
    "fieldConfig": {
    "defaults": {
        "custom": {}
    },
    "overrides": []
    },
    "fill": 1,
    "fillGradient": 0,
    "gridPos": {
    "h": 12,
    "w": 12,
    "x": 12,
    "y": 9
    },
    "hiddenSeries": false,
    "id": 4,
    "legend": {
    "avg": false,
    "current": false,
    "max": false,
    "min": false,
    "show": true,
    "total": false,
    "values": false
    },
    "lines": true,
    "linewidth": 1,
    "nullPointMode": "null",
    "options": {
```

```
"alertThreshold": true
},
"percentage": false,
"pluginVersion": "7.4.5",
"pointradius": 2,
"points": false,
"renderer": "flot",
"seriesOverrides": [],
"spaceLength": 10,
"stack": false,
"steppedLine": false,
"targets": [
{
        "expr": "DCGM_FI_DEV_GPU_TEMP",
        "interval": "",
        "legendFormat": "",
        "refId": "A"
}
],
"thresholds": [],
"timeFrom": null,
"timeRegions": [],
"timeShift": null,
"title": "GPU Temperature",
"tooltip": {
"shared": true,
"sort": 0,
"value_type": "individual"
},
"type": "graph",
"xaxis": {
"buckets": null,
"mode": "time",
"name": null,
"show": true,
"values": []
},
"yaxes": [
{
        "format": "celsius",
        "label": "Celsius (°C)",
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
},
```

```
{
      "format": "short",
      "label": null,
      "logBase": 1,
      "max": null,
      "min": null,
      "show": true
}
],
"yaxis": {
"align": false,
"alignLevel": null
}
},
{
"aliasColors": {},
"bars": false,
"dashLength": 10,
"dashes": false,
"datasource": "${DS_PROMETHEUS}",
"fieldConfig": {
"defaults": {
      "custom": {}
},
"overrides": []
},
"fill": 1,
"fillGradient": 0,
"gridPos": {
"h": 11,
"w": 12,
"x": 0,
"y": 16
},
"hiddenSeries": false,
"id": 22,
"legend": {
"avg": false,
"current": false,
"max": false,
"min": false,
"show": true,
"total": false,
"values": false
},
"lines": true,
```

```
"linewidth": 1,
"nullPointMode": "null",
"options": {
"alertThreshold": true
},
"percentage": false,
"pluginVersion": "7.4.5",
"pointradius": 2,
"points": false,
"renderer": "flot",
"seriesOverrides": [],
"spaceLength": 10,
"stack": false,
"steppedLine": false,
"targets": [
{
        "expr": "sum(DCGM_FI_DEV_POWER_USAGE{modelName=\"NVIDIA
A100-PCIE-40GB\"})",
        "interval": "",
        "legendFormat": "",
        "refId": "A"
},
{
        "expr": "sum(DCGM_FI_DEV_POWER_USAGE{modelName=\"Tesla T4\"})",
        "hide": false,
        "interval": "",
        "legendFormat": "",
        "refId": "B"
},
{
        "expr": "sum(DCGM_FI_DEV_POWER_USAGE{modelName=\"Tesla
V100-PCIE-32GB\"})",
        "hide": false,
        "interval": "",
        "legendFormat": "",
        "refId": "C"
}
],
"thresholds": [],
"timeRegions": [],
"title": "GPU Power Consumption by Type - Watt (W)",
"tooltip": {
"shared": true,
"sort": 0,
"value_type": "individual"
},
```

```
"type": "graph",
"xaxis": {
"buckets": null,
"mode": "time",
"name": null,
"show": true,
"values": []
},
"yaxes": [
{
        "format": "watt",
        "label": null,
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
},
{
        "format": "short",
        "label": null,
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
}
],
"yaxis": {
"align": false,
"alignLevel": null
}
},
{
"aliasColors": {},
"bars": false,
"dashLength": 10,
"dashes": false,
"datasource": "${DS_PROMETHEUS}",
"fieldConfig": {
"defaults": {
        "custom": {}
},
"overrides": []
},
"fill": 1,
"fillGradient": 0,
"gridPos": {
```

```
        "h": 12,
        "w": 12,
        "x": 12,
        "y": 21
      },
      "hiddenSeries": false,
      "id": 8,
      "legend": {
        "avg": false,
        "current": false,
        "max": false,
        "min": false,
        "show": true,
        "total": false,
        "values": false
      },
      "lines": true,
      "linewidth": 1,
      "nullPointMode": "null",
      "options": {
        "alertThreshold": true
      },
      "percentage": false,
      "pluginVersion": "7.4.5",
      "pointradius": 2,
      "points": false,
      "renderer": "flot",
      "seriesOverrides": [],
      "spaceLength": 10,
      "stack": false,
      "steppedLine": false,
      "targets": [
        {
            "expr": "DCGM_FI_DEV_POWER_USAGE",
            "interval": "",
            "legendFormat": "",
            "refId": "A"
        }
      ],
      "thresholds": [],
      "timeFrom": null,
      "timeRegions": [],
      "timeShift": null,
      "title": "GPU Power Consumption",
      "tooltip": {
        "shared": true,
```

```
"sort": 0,
"value_type": "individual"
},
"type": "graph",
"xaxis": {
"buckets": null,
"mode": "time",
"name": null,
"show": true,
"values": []
},
"yaxes": [
{
        "format": "watt",
        "label": "Watt (W)",
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
},
{
        "format": "short",
        "label": null,
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
}
],
"yaxis": {
"align": false,
"alignLevel": null
}
},
{
"aliasColors": {},
"bars": false,
"dashLength": 10,
"dashes": false,
"datasource": "${DS_PROMETHEUS}",
"fieldConfig": {
"defaults": {
        "custom": {}
},
"overrides": []
},
```

```
"fill": 1,
"fillGradient": 0,
"gridPos": {
"h": 9,
"w": 12,
"x": 0,
"y": 27
},
"hiddenSeries": false,
"id": 26,
"legend": {
"avg": false,
"current": false,
"max": false,
"min": false,
"show": true,
"total": false,
"values": false
},
"lines": true,
"linewidth": 1,
"nullPointMode": "null",
"options": {
"alertThreshold": true
},
"percentage": false,
"pluginVersion": "7.4.5",
"pointradius": 2,
"points": false,
"renderer": "flot",
"seriesOverrides": [],
"spaceLength": 10,
"stack": false,
"steppedLine": false,
"targets": [
{
      "expr": "sum(DCGM_FI_DEV_POWER_USAGE{modelName=\"Tesla T4\"})",
      "interval": "",
      "legendFormat": "",
      "refId": "A"
}
],
"thresholds": [],
"timeRegions": [],
"title": "T4 Power Consumption Watt (W)",
"tooltip": {
```

```
"shared": true,
"sort": 0,
"value_type": "individual"
},
"type": "graph",
"xaxis": {
"buckets": null,
"mode": "time",
"name": null,
"show": true,
"values": []
},
"yaxes": [
{
        "format": "watt",
        "label": null,
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
},
{
        "format": "short",
        "label": null,
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
}
],
"yaxis": {
"align": false,
"alignLevel": null
}
},
{
"aliasColors": {},
"bars": false,
"dashLength": 10,
"dashes": false,
"datasource": "${DS_PROMETHEUS}",
"fieldConfig": {
"defaults": {
        "custom": {}
},
"overrides": []
```

```json
},
"fill": 1,
"fillGradient": 0,
"gridPos": {
"h": 12,
"w": 12,
"x": 12,
"y": 33
},
"hiddenSeries": false,
"id": 6,
"legend": {
"avg": false,
"current": false,
"max": false,
"min": false,
"show": true,
"total": false,
"values": false
},
"lines": true,
"linewidth": 1,
"nullPointMode": "null",
"options": {
"alertThreshold": true
},
"percentage": false,
"pluginVersion": "7.4.5",
"pointradius": 2,
"points": false,
"renderer": "flot",
"seriesOverrides": [],
"spaceLength": 10,
"stack": false,
"steppedLine": false,
"targets": [
{
      "expr": "DCGM_FI_DEV_GPU_UTIL",
      "interval": "",
      "legendFormat": "",
      "refId": "A"
}
],
"thresholds": [],
"timeFrom": null,
"timeRegions": [],
```

```json
"timeShift": null,
"title": "GPU Utilization",
"tooltip": {
"shared": true,
"sort": 0,
"value_type": "individual"
},
"type": "graph",
"xaxis": {
"buckets": null,
"mode": "time",
"name": null,
"show": true,
"values": []
},
"yaxes": [
{
        "format": "percent",
        "label": null,
        "logBase": 1,
        "max": "110",
        "min": "0",
        "show": true
},
{
        "format": "short",
        "label": null,
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
}
],
"yaxis": {
"align": false,
"alignLevel": null
}
},
{
"aliasColors": {},
"bars": false,
"dashLength": 10,
"dashes": false,
"datasource": "${DS_PROMETHEUS}",
"fieldConfig": {
"defaults": {
```

```
        "custom": {}
},
"overrides": []
},
"fill": 1,
"fillGradient": 0,
"gridPos": {
"h": 12,
"w": 12,
"x": 0,
"y": 36
},
"hiddenSeries": false,
"id": 20,
"legend": {
"avg": false,
"current": false,
"max": false,
"min": false,
"show": true,
"total": false,
"values": false
},
"lines": true,
"linewidth": 1,
"nullPointMode": "null",
"options": {
"alertThreshold": true
},
"percentage": false,
"pluginVersion": "7.4.5",
"pointradius": 2,
"points": false,
"renderer": "flot",
"seriesOverrides": [],
"spaceLength": 10,
"stack": false,
"steppedLine": false,
"targets": [
{
        "expr": "node_memory_Active_bytes",
        "interval": "",
        "legendFormat": "",
        "refId": "A"
}
],
```

```
"thresholds": [],
"timeRegions": [],
"title": "Memory Utilization",
"tooltip": {
"shared": true,
"sort": 0,
"value_type": "individual"
},
"type": "graph",
"xaxis": {
"buckets": null,
"mode": "time",
"name": null,
"show": true,
"values": []
},
"yaxes": [
{
        "format": "decbytes",
        "label": null,
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
},
{
        "format": "short",
        "label": null,
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
}
],
"yaxis": {
"align": false,
"alignLevel": null
}
},
{
"aliasColors": {},
"bars": false,
"dashLength": 10,
"dashes": false,
"datasource": "${DS_PROMETHEUS}",
"fieldConfig": {
```

```
"defaults": {
      "custom": {}
},
"overrides": []
},
"fill": 1,
"fillGradient": 0,
"gridPos": {
"h": 12,
"w": 12,
"x": 0,
"y": 48
},
"hiddenSeries": false,
"id": 18,
"legend": {
"avg": false,
"current": false,
"max": false,
"min": false,
"show": true,
"total": false,
"values": false
},
"lines": true,
"linewidth": 1,
"nullPointMode": "null",
"options": {
"alertThreshold": true
},
"percentage": false,
"pluginVersion": "7.4.5",
"pointradius": 2,
"points": false,
"renderer": "flot",
"seriesOverrides": [],
"spaceLength": 10,
"stack": false,
"steppedLine": false,
"targets": [
{
      "expr": "node_load15",
      "interval": "",
      "legendFormat": "",
      "refId": "A"
}
```

```
],
"thresholds": [],
"timeFrom": null,
"timeRegions": [],
"timeShift": null,
"title": "CPU Load 15 min",
"tooltip": {
"shared": true,
"sort": 0,
"value_type": "individual"
},
"type": "graph",
"xaxis": {
"buckets": null,
"mode": "time",
"name": null,
"show": true,
"values": []
},
"yaxes": [
{
        "format": "short",
        "label": null,
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
},
{
        "format": "short",
        "label": null,
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
}
],
"yaxis": {
"align": false,
"alignLevel": null
}
},
{
"aliasColors": {},
"bars": false,
"dashLength": 10,
```

```json
"dashes": false,
"datasource": "${DS_PROMETHEUS}",
"fieldConfig": {
"defaults": {
      "custom": {}
},
"overrides": []
},
"fill": 1,
"fillGradient": 0,
"gridPos": {
"h": 8,
"w": 12,
"x": 0,
"y": 60
},
"hiddenSeries": false,
"id": 16,
"legend": {
"avg": false,
"current": false,
"max": false,
"min": false,
"show": true,
"total": false,
"values": false
},
"lines": true,
"linewidth": 1,
"nullPointMode": "null",
"options": {
"alertThreshold": true
},
"percentage": false,
"pluginVersion": "7.4.5",
"pointradius": 2,
"points": false,
"renderer": "flot",
"seriesOverrides": [],
"spaceLength": 10,
"stack": false,
"steppedLine": false,
"targets": [
{
      "expr": "node_load5",
      "interval": "",
```

```
        "legendFormat": "",
        "refId": "A"
}
],
"thresholds": [],
"timeFrom": null,
"timeRegions": [],
"timeShift": null,
"title": "CPU Load 5 min",
"tooltip": {
"shared": true,
"sort": 0,
"value_type": "individual"
},
"type": "graph",
"xaxis": {
"buckets": null,
"mode": "time",
"name": null,
"show": true,
"values": []
},
"yaxes": [
{
      "format": "short",
      "label": null,
      "logBase": 1,
      "max": null,
      "min": null,
      "show": true
},
{
      "format": "short",
      "label": null,
      "logBase": 1,
      "max": null,
      "min": null,
      "show": true
}
],
"yaxis": {
"align": false,
"alignLevel": null
}
},
{
```

```
"aliasColors": {},
"bars": false,
"dashLength": 10,
"dashes": false,
"datasource": "${DS_PROMETHEUS}",
"fieldConfig": {
"defaults": {
      "custom": {}
},
"overrides": []
},
"fill": 1,
"fillGradient": 0,
"gridPos": {
"h": 8,
"w": 12,
"x": 0,
"y": 68
},
"hiddenSeries": false,
"id": 14,
"legend": {
"avg": false,
"current": false,
"max": false,
"min": false,
"show": true,
"total": false,
"values": false
},
"lines": true,
"linewidth": 1,
"nullPointMode": "null",
"options": {
"alertThreshold": true
},
"percentage": false,
"pluginVersion": "7.4.5",
"pointradius": 2,
"points": false,
"renderer": "flot",
"seriesOverrides": [],
"spaceLength": 10,
"stack": false,
"steppedLine": false,
"targets": [
```

```
            {
                  "expr": "node_load1",
                  "interval": "",
                  "legendFormat": "",
                  "refId": "A"
            }
            ],
            "thresholds": [],
            "timeFrom": null,
            "timeRegions": [],
            "timeShift": null,
            "title": "CPU Load 1 min",
            "tooltip": {
            "shared": true,
            "sort": 0,
            "value_type": "individual"
            },
            "type": "graph",
            "xaxis": {
            "buckets": null,
            "mode": "time",
            "name": null,
            "show": true,
            "values": []
            },
            "yaxes": [
            {
                  "format": "short",
                  "label": null,
                  "logBase": 1,
                  "max": null,
                  "min": null,
                  "show": true
            },
            {
                  "format": "short",
                  "label": null,
                  "logBase": 1,
                  "max": null,
                  "min": null,
                  "show": true
            }
            ],
            "yaxis": {
            "align": false,
            "alignLevel": null
```

```
        }
    },
    {
    "aliasColors": {},
    "bars": false,
    "dashLength": 10,
    "dashes": false,
    "datasource": "${DS_PROMETHEUS}",
    "fieldConfig": {
    "defaults": {
            "custom": {}
    },
    "overrides": []
    },
    "fill": 1,
    "fillGradient": 0,
    "gridPos": {
    "h": 8,
    "w": 12,
    "x": 0,
    "y": 76
    },
    "hiddenSeries": false,
    "id": 10,
    "legend": {
    "avg": false,
    "current": false,
    "max": false,
    "min": false,
    "show": true,
    "total": false,
    "values": false
    },
    "lines": true,
    "linewidth": 1,
    "nullPointMode": "null",
    "options": {
    "alertThreshold": true
    },
    "percentage": false,
    "pluginVersion": "7.4.5",
    "pointradius": 2,
    "points": false,
    "renderer": "flot",
    "seriesOverrides": [],
    "spaceLength": 10,
```

```
"stack": false,
"steppedLine": false,
"targets": [
{
        "expr": "node_infiniband_rate_bytes_per_second",
        "interval": "",
        "legendFormat": "",
        "refId": "A"
}
],
"thresholds": [],
"timeFrom": null,
"timeRegions": [],
"timeShift": null,
"title": "Panel Title",
"tooltip": {
"shared": true,
"sort": 0,
"value_type": "individual"
},
"type": "graph",
"xaxis": {
"buckets": null,
"mode": "time",
"name": null,
"show": true,
"values": []
},
"yaxes": [
{
        "format": "short",
        "label": null,
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
},
{
        "format": "short",
        "label": null,
        "logBase": 1,
        "max": null,
        "min": null,
        "show": true
}
],
```

```
"yaxis": {
"align": false,
"alignLevel": null
}
},
{
"aliasColors": {},
"bars": false,
"dashLength": 10,
"dashes": false,
"datasource": "${DS_PROMETHEUS}",
"fieldConfig": {
"defaults": {
      "custom": {}
},
"overrides": []
},
"fill": 1,
"fillGradient": 0,
"gridPos": {
"h": 9,
"w": 12,
"x": 0,
"y": 84
},
"hiddenSeries": false,
"id": 2,
"legend": {
"avg": false,
"current": false,
"max": false,
"min": false,
"show": true,
"total": false,
"values": false
},
"lines": true,
"linewidth": 1,
"nullPointMode": "null",
"options": {
"alertThreshold": true
},
"percentage": false,
"pluginVersion": "7.4.5",
"pointradius": 2,
"points": false,
```

```
"renderer": "flot",
"seriesOverrides": [],
"spaceLength": 10,
"stack": false,
"steppedLine": false,
"targets": [
{
      "expr": "sum(DCGM_FI_DEV_POWER_USAGE)",
      "instant": false,
      "interval": "",
      "intervalFactor": 1,
      "legendFormat": "",
      "refId": "Total Power"
}
],
"thresholds": [],
"timeFrom": null,
"timeRegions": [],
"timeShift": null,
"title": "Cluster-wide GPU power consumption",
"tooltip": {
"shared": true,
"sort": 0,
"value_type": "individual"
},
"type": "graph",
"xaxis": {
"buckets": null,
"mode": "time",
"name": null,
"show": true,
"values": []
},
"yaxes": [
{
      "format": "watt",
      "label": "Watts",
      "logBase": 1,
      "max": null,
      "min": null,
      "show": true
},
{
      "format": "short",
      "label": null,
      "logBase": 1,
```

```
                    "max": null,
                    "min": null,
                    "show": true
                }
            ],
            "yaxis": {
            "align": false,
            "alignLevel": null
            }
            }
    ],
    "refresh": false,
    "schemaVersion": 27,
    "style": "dark",
    "tags": [],
    "templating": {
        "list": []
    },
    "time": {
        "from": "2022-04-20T13:00:00.000Z",
        "to": "2022-04-28T13:00:00.000Z"
    },
    "timepicker": {},
    "timezone": "",
    "title": "GPU Dashboard",
    "uid": "BEu-j-wnk",
    "version": 7
}
```