

HPC-I/O Analyse & Visualisierung mit Grafana

Simon Rosenberger

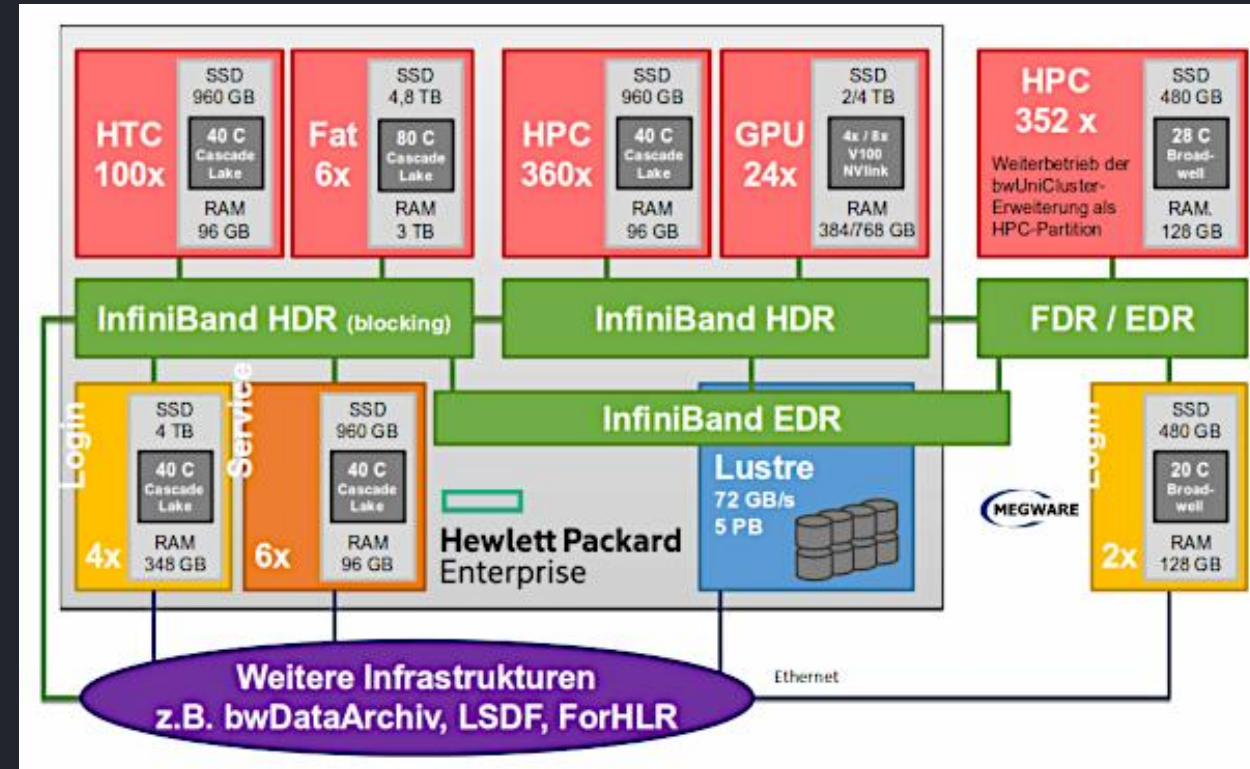
Gliederung

1. Motivation
2. Libiotrace
3. Was sind InfluxDB / Grafana?
4. Grafana Plugins
5. Grafana Dashboard

Motivation

Systeme wie BwUniCluster:

- Große Mengen an Speicher & Rechenkapazität
- Parallele Anbindung über InfiniBand
- Disc-IO Anbindung über IB-Netzwerk ist langsam



Motivation

Problem Livetracing

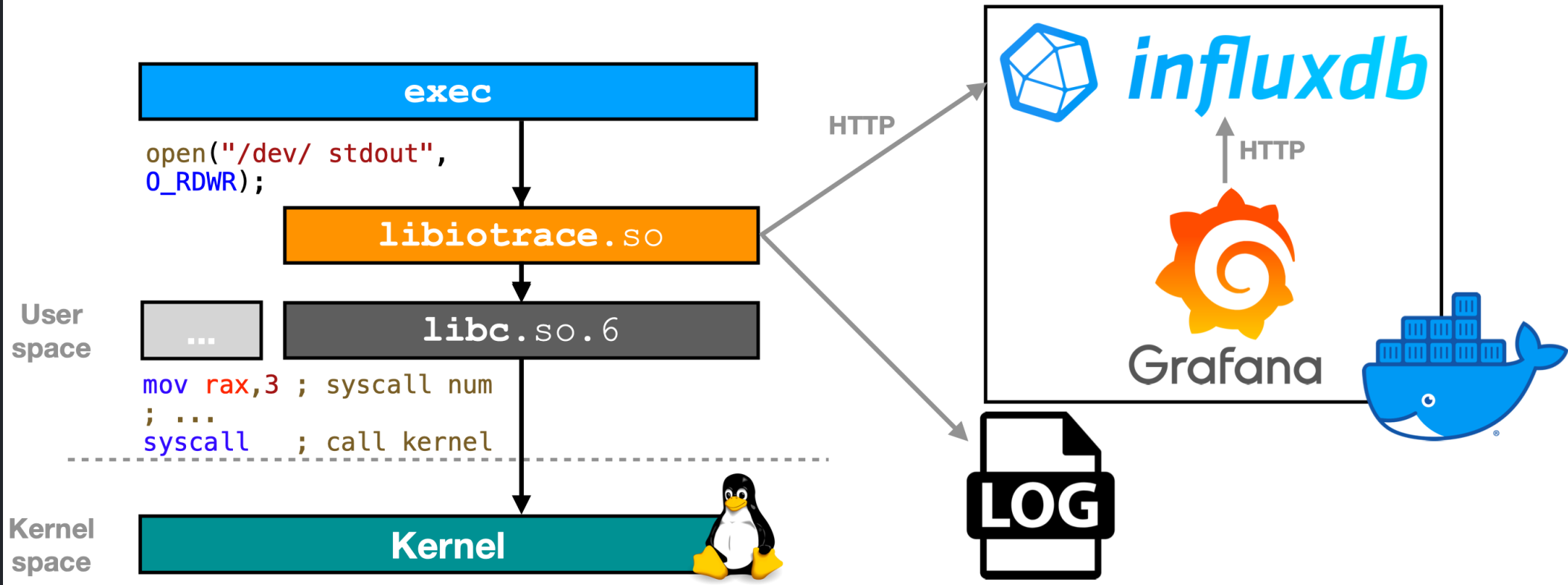
- Nicht erkennbar welche Prozesse und Threads im System viel schreiben & lesen und die Performance bremsen
→ unerkannte Bottlenecks

Lösungsansatz

- „HeatMap“ mit Gesamtübersicht die die Auslastung von Prozessen und unterlagerter Threads anzeigt
- Detailliertere Daten mithilfe eines Force-Graphen

Libiotrace

→ ~ **LD_PRELOAD**=path/to/libiotrace.so ./exec





- Opensource, nicht relationale Datenbank
- Datenabfrage über eigene Sprache: „Flux“
- Einfache Integration in Visualisierung über Queries
- Optimierte für High Performance Computing mit vielen Daten



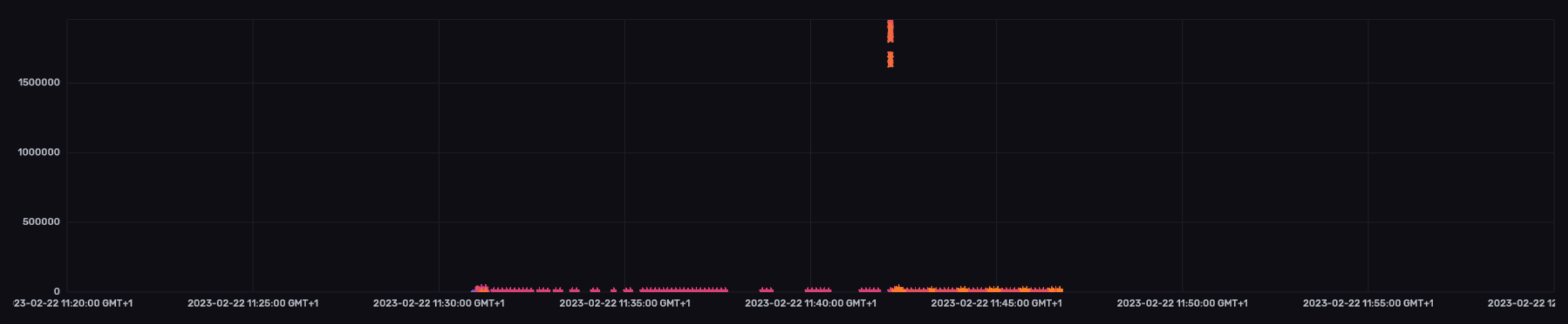
Data Explorer

Scatter

Customize

Local

Save As



Query 1 (0.39s)

+

View Raw Data

CSV

II

2023-02-22 11:20 - 2023-02-22 12:00

Script Editor

Submit

FROM

Search for a bucket

_monitoring

_tasks

hsebucket

hsebucket_fileres

hsebucket_randomIO

hsebucket_randomIO2

+ Create Bucket

Filter

_measurement

Search _measurement tag values

libiotrace

libiotrace_filesystem

Filter

_field

Search _field tag values

file_type_descriptor

file_type_stream

function_data_read_bytes

function_data_written_bytes

function_name

hostname

process_id

return_state

return_state_detail_errno...

Filter

functionname

Search functionname tag values

fread

fwrite

getdelim

read

write

writew

WINDOW PERIOD

Custom

Auto

auto (6s667ms)

Fill missing values

AGGREGATE FUNCTION

Custom

Auto

mean

median

last



- Open Source Software für Visualisierung von Zeitreihen
- Daten aus Influx, SQL, etc. verarbeitbar
- Dashboards können aus vorgefertigten und selbst erstellten Plugins zusammengestellt werden

Plugins – Flux Queries

Written Bytes (InfluxDB-LiveTracing)

```
1 from(bucket: "hsebucket_fileres")
2   |> range(start: v.timeRangeStart, stop: v.timeRangeStop)
3   |> filter(fn: (r) => r["_measurement"] == "libiotrace")
4   |> filter(fn: (r) => r["_field"] == "function_data_written_bytes")
5   |> filter(fn: (r) => r["functionname"] == "fwrite" or r["functionname"] == "write" or r["functionname"] == "writev")
6   |> group(columns: ["thread", "processid"])
7   |> aggregateWindow(every: inf, fn: sum, createEmpty: false)
8   |> yield(name: "sum")
```

[Flux language syntax](#)[Sample Query](#)[Help](#)

> **Read Bytes** (InfluxDB-LiveTracing) `from(bucket: "hsebucket_fileres") |> range(start: v.timeRangeStart, stop: v.timeRangeStop) |> filter(fn: (r) => r["_measurement"] == "libiotrace") |> filter(fn: (r) => r["_fi...`

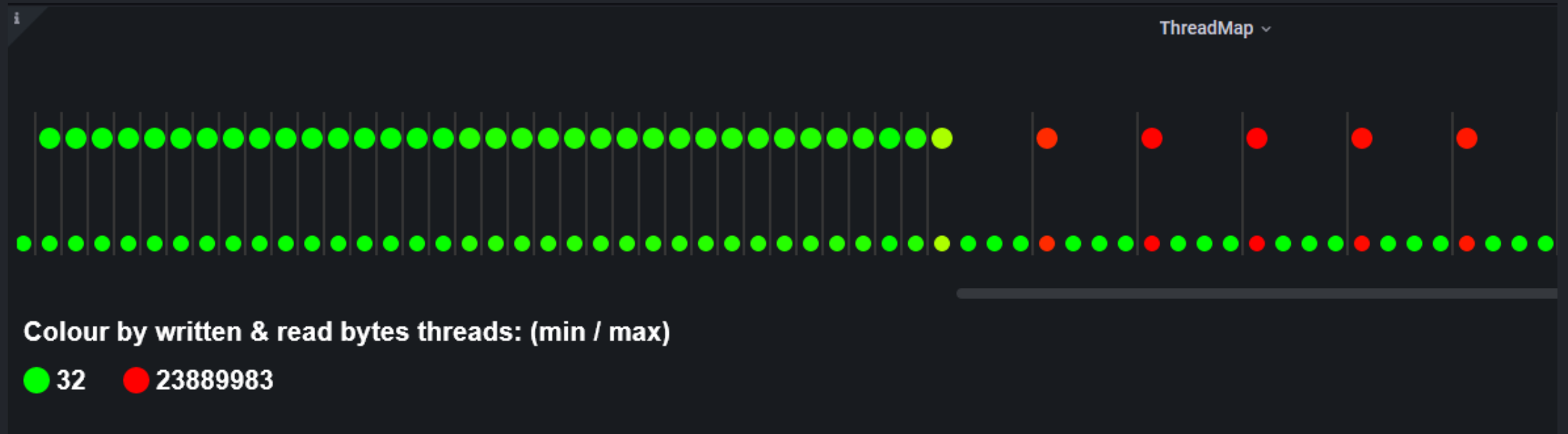
Filename Write (InfluxDB-LiveTracing)

```
1 from(bucket: "hsebucket_fileres")
2   |> range(start: v.timeRangeStart, stop: v.timeRangeStop)
3   |> filter(fn: (r) => r["_measurement"] == "libiotrace")
4   |> filter(fn: (r) => r["_field"] == "traced_filename")
5   |> filter(fn: (r) => r["functionname"] == "fwrite" or r["functionname"] == "write" or r["functionname"] == "writev")
6   |> group(columns: ["thread", "processid"])
7   |> aggregateWindow(every: inf, column: "_value", fn: distinct, createEmpty: false)
8   |> yield(name: "distinct")
```

[Flux language syntax](#)[Sample Query](#)[Help](#)

> **Filename Read** (InfluxDB-LiveTracing) `from(bucket: "hsebucket_fileres") |> range(start: v.timeRangeStart, stop: v.timeRangeStop) |> filter(fn: (r) => r["_measurement"] == "libiotrace") |> filter(fn: (r) => r...`

Plugins - Thread Map



Plugins - Thread Map



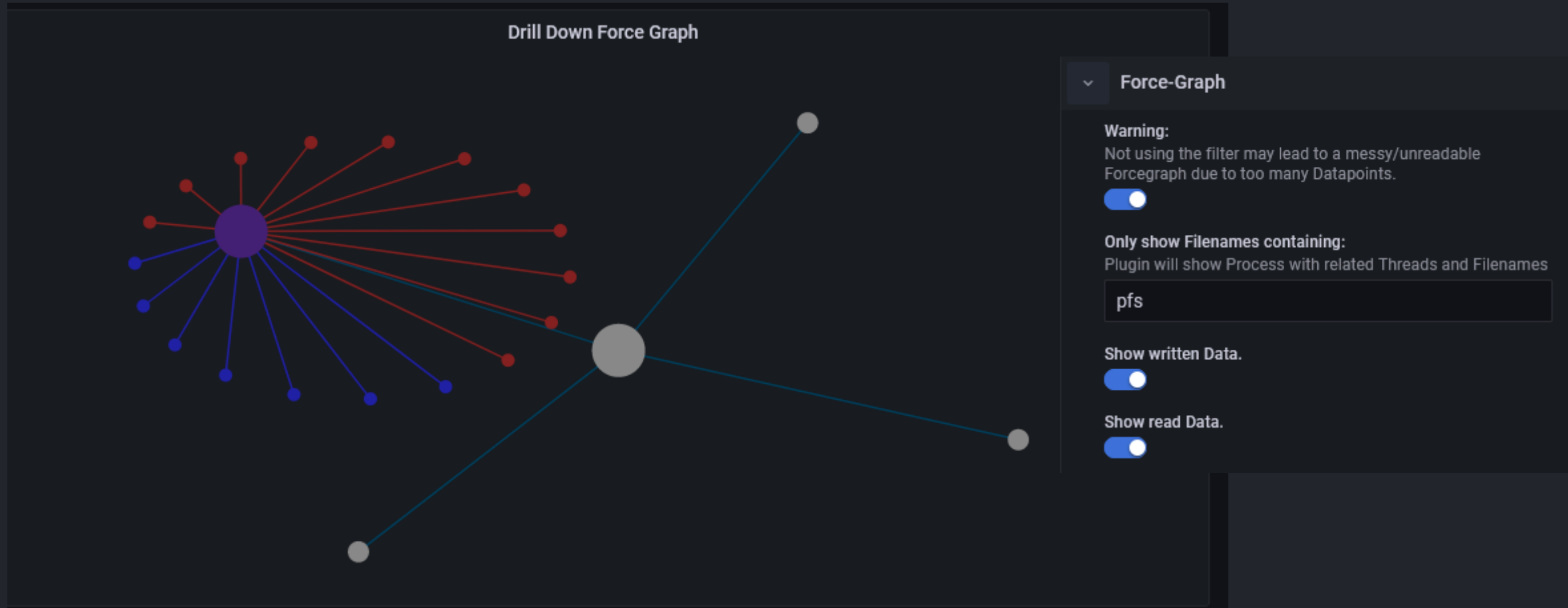
Plugins - Thread Map



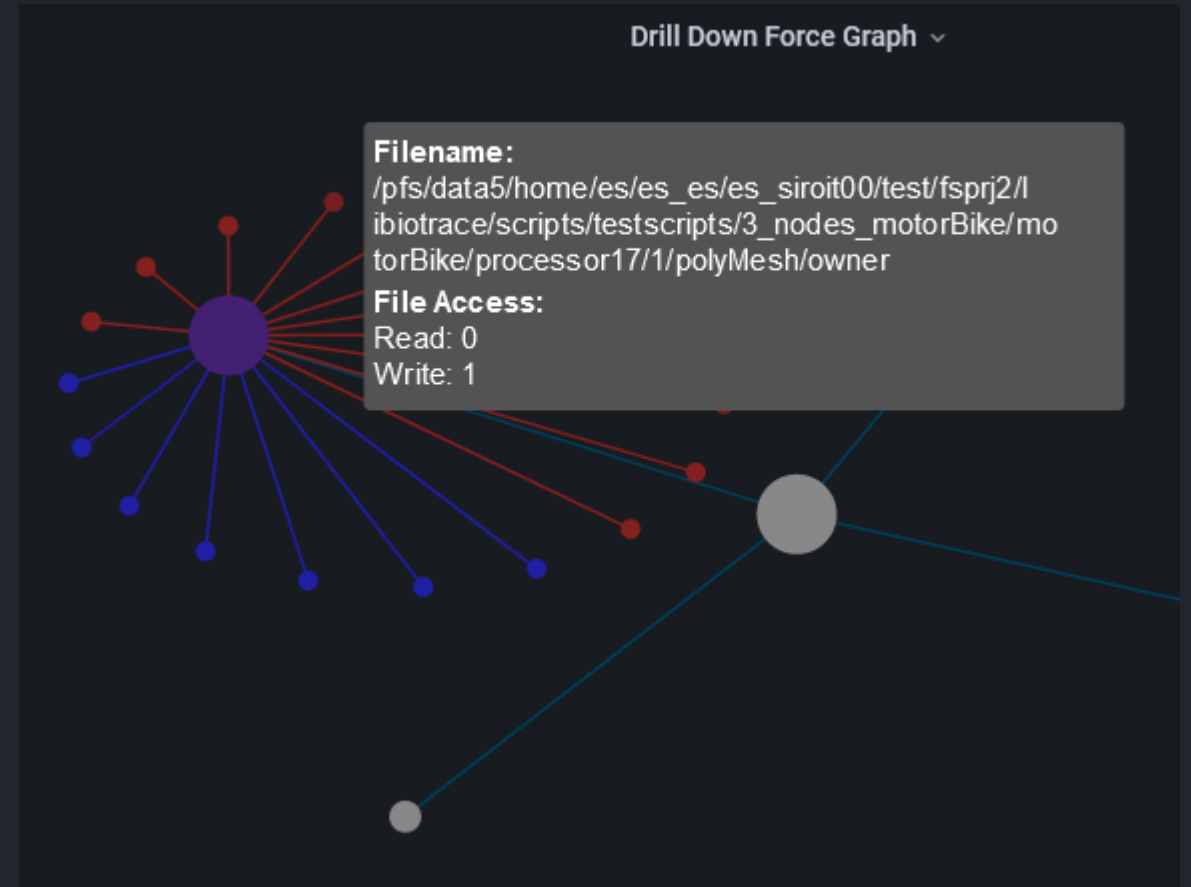
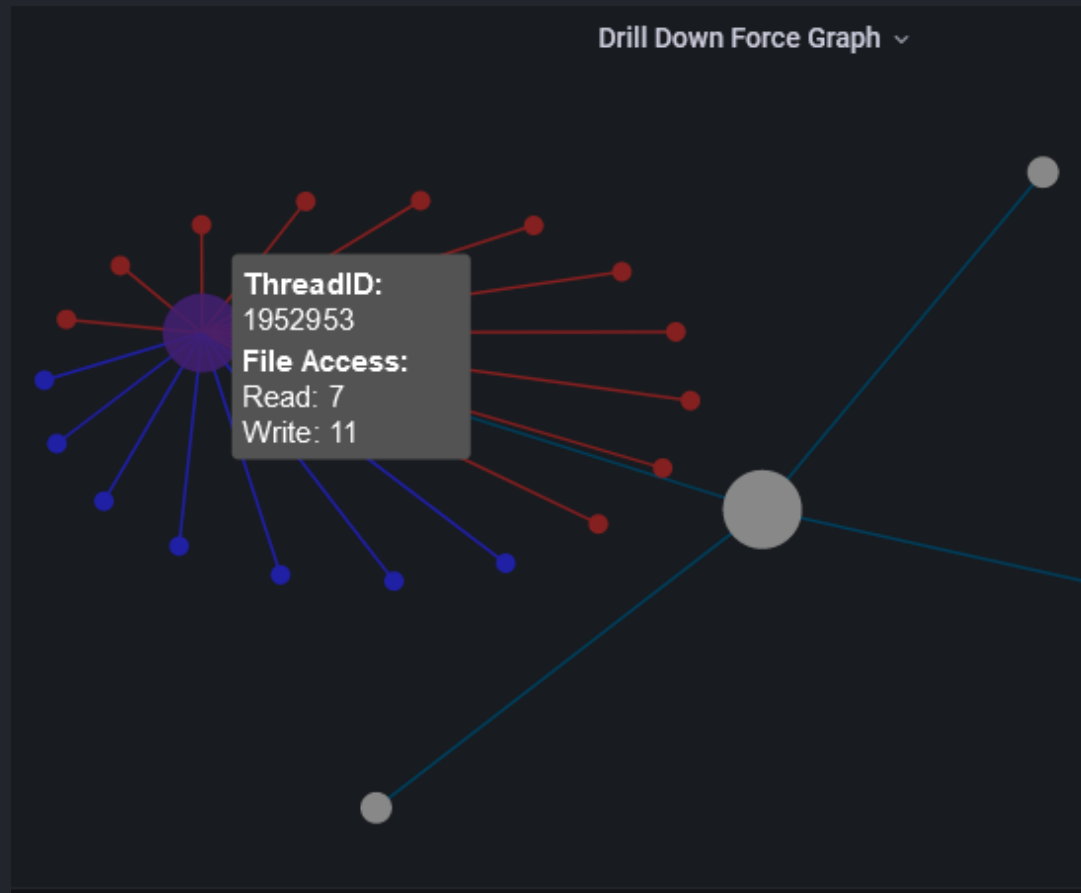
Plugins - Tables

Filename Write
TracedFilenameWrite {processid="1952952", thread="1952952"}
_ NOT FOUND _
_ PSEUDO-FILE _
/dev/infiniband/rdma_cm
{processid="1952952", thread="1952952"} ▾
Filename Read
TracedFilenameRead {processid="1952952", thread="1952952"}
/opt/bwhpc/common/mpi/openmpi/4.1.4-gnu-12.1/etc/openmpi-mca-params.conf
/opt/bwhpc/common/mpi/openmpi/4.1.4-gnu-12.1/etc/pmix-mca-params.conf
/opt/bwhpc/common/mpi/openmpi/4.1.4-gnu-12.1/share/openmpi/mca-btl-openib-device-params.ini
{processid="1952952", thread="1952952"} ▾

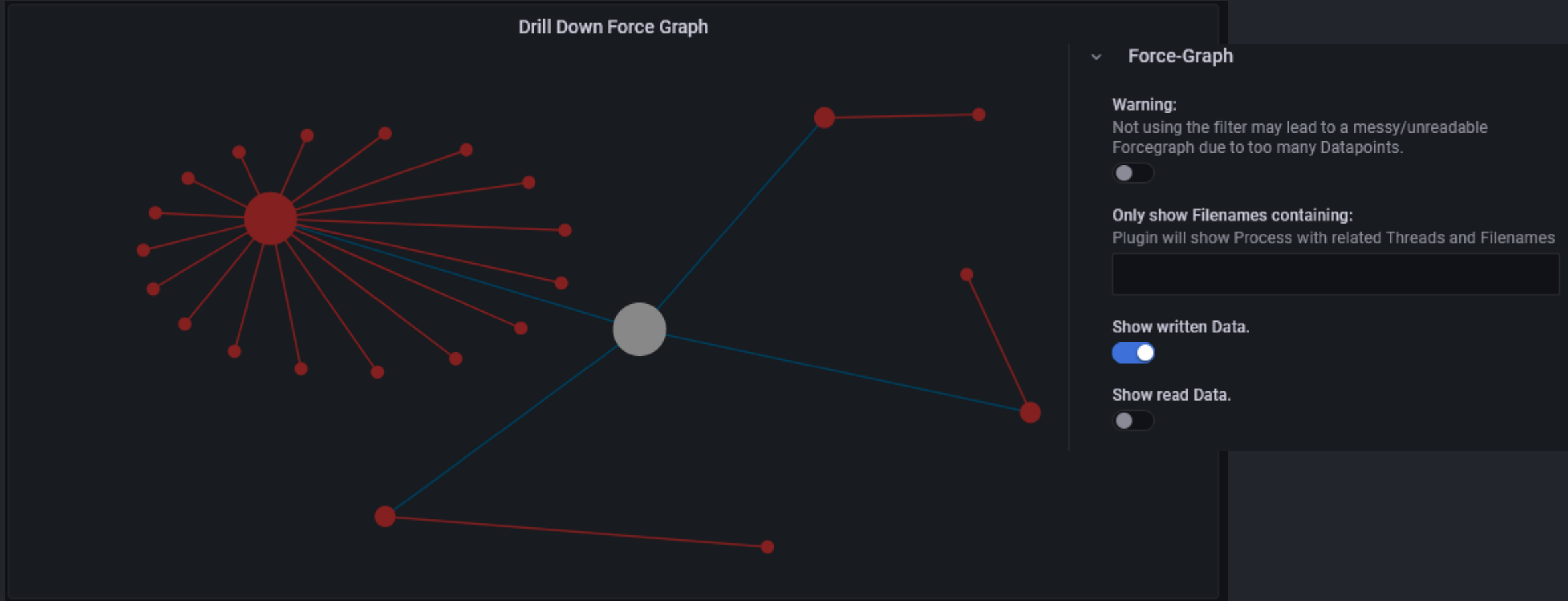
Plugins - Force Graph



Plugins - Force Graph

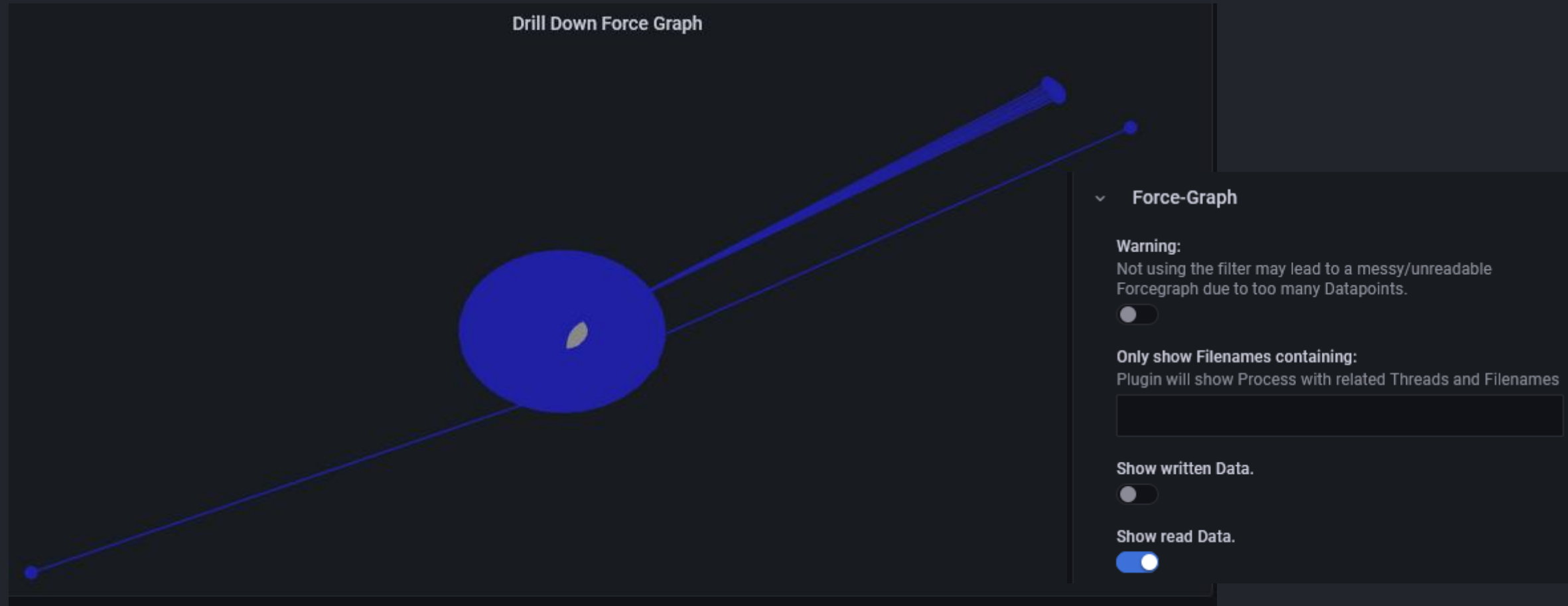


Plugins - Force Graph



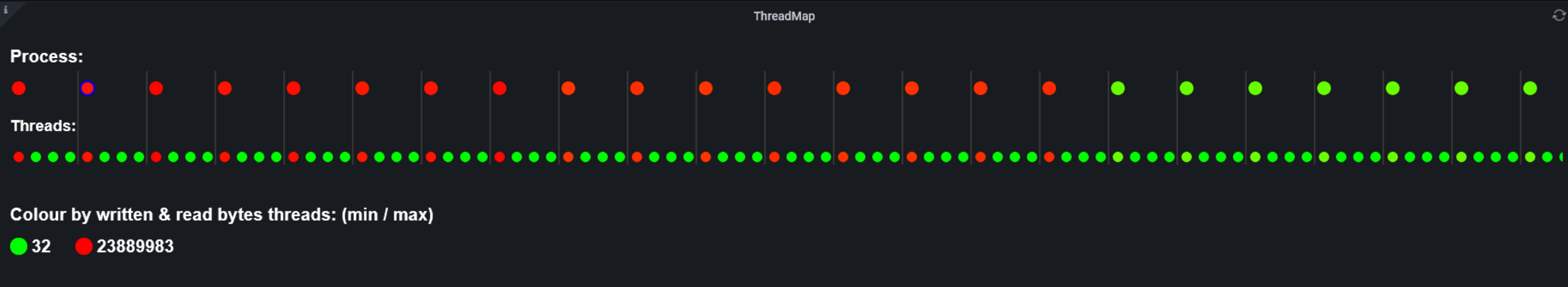
Plugins - Force Graph

Ca. 2000 Datenpunkte



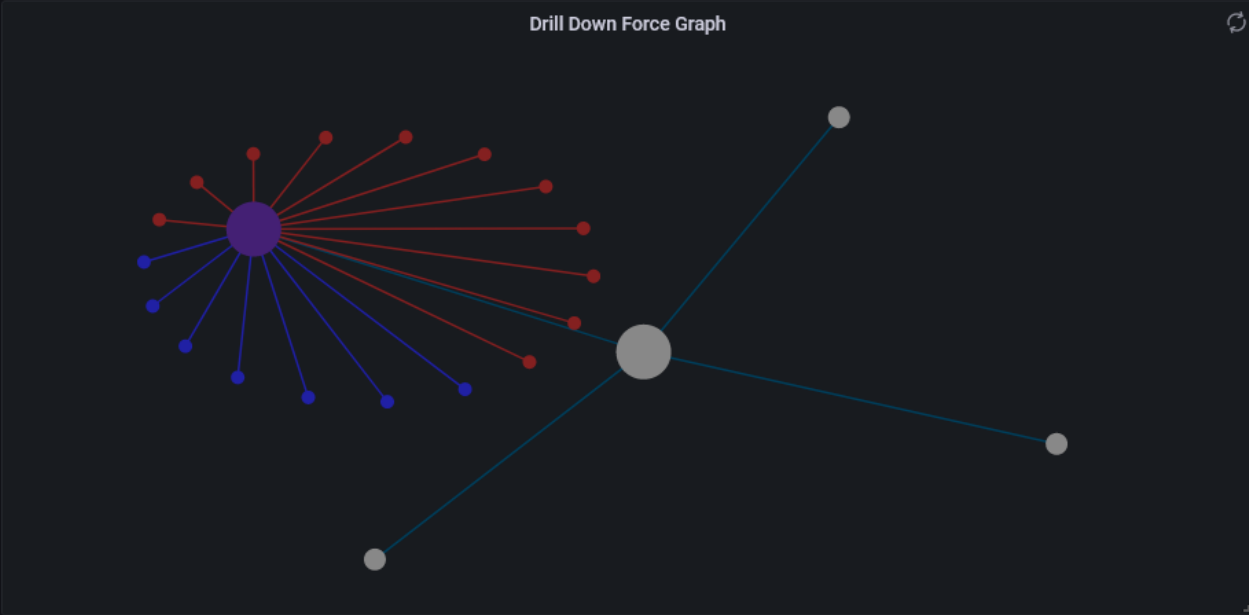
Dashboard

ThreadMap & Drilldown



Filename Write
TracedFilenameWrite {processid="1952952", thread="1952952"}
_ NOT FOUND _
_ PSEUDO-FILE _
/dev/infiniband/rdma_cm
{processid="1952952", thread="1952952"}

Filename Read
TracedFilenameRead {processid="1952952", thread="1952952"}
/opt/bwhpc/common/mpi/openmpi/4.1.4-gnu-12.1/etc/openmpi-mca-params.conf
/opt/bwhpc/common/mpi/openmpi/4.1.4-gnu-12.1/etc/pmix-mca-params.conf
/opt/bwhpc/common/mpi/openmpi/4.1.4-gnu-12.1/share/openmpi/mca-btl-openib-device-params.ini
{processid="1952952", thread="1952952"}



Vielen Dank für ihre
Aufmerksamkeit!

Fragen?