# Motivation



File IO Lustre Bandwidth (mean across nodes)

— scratch_read_bw    — scratch_write_bw

Infiniband Bandwidth (mean across nodes)

TECHNISCHE
UNIVERSITÄT
DRESDEN

ZIH
Center for Information Services &
High Performance Computing

# Throughput test with a sane I/O pattern



I/O throughput BeeGFS

# I/O performance impact factors

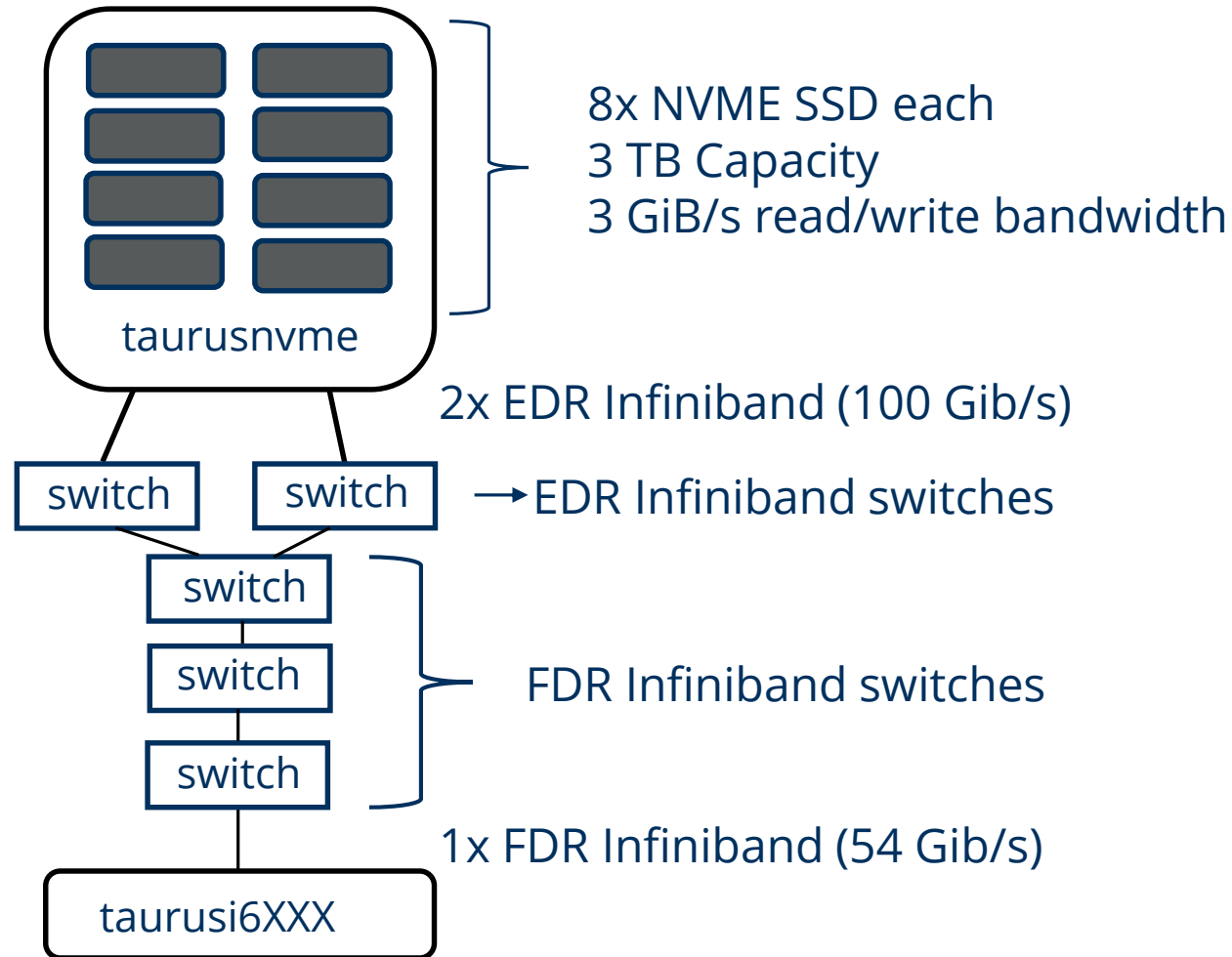## Application

- Request sizes
- Access pattern
- I/O operation

## Network

- Message sizes
- Network paths

## File system
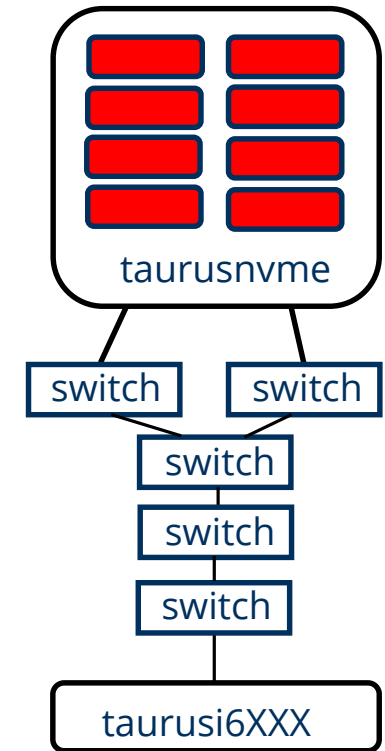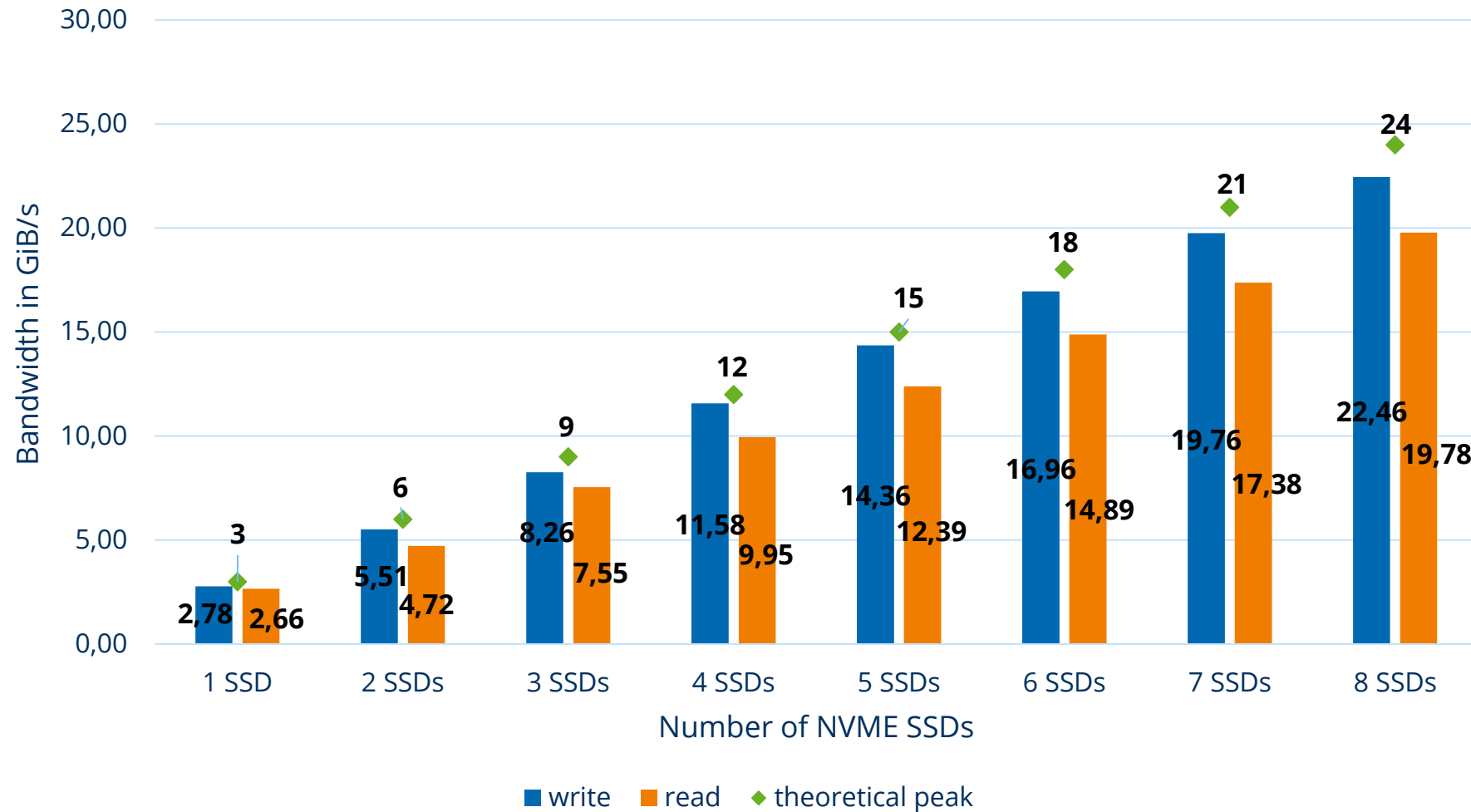
- Stripe sizes
- File hierarchy

# Setup

8x NVME SSD each
3 TB Capacity
3 GiB/s read/write bandwidth

taurusnvme

2x EDR Infiniband (100 Gib/s)

switch    switch    → EDR Infiniband switches

switch

switch    FDR Infiniband switches

switch

1x FDR Infiniband (54 Gib/s)

taurusi6XXX

## Application: IOR

- Blocked sequential I/O pattern
- File per process
- 2 MiB request sizes

TECHNISCHE
UNIVERSITÄT
DRESDEN

ZIH
Center for Information Services &
High Performance Computing

# Single NVME SSD performance



I/O throughput local NVME SSDs

# Check message sizes of network links



Message size test from taurusi6001 to taurusnvme22

Characterization of Infiniband routes to support data intensive I/O
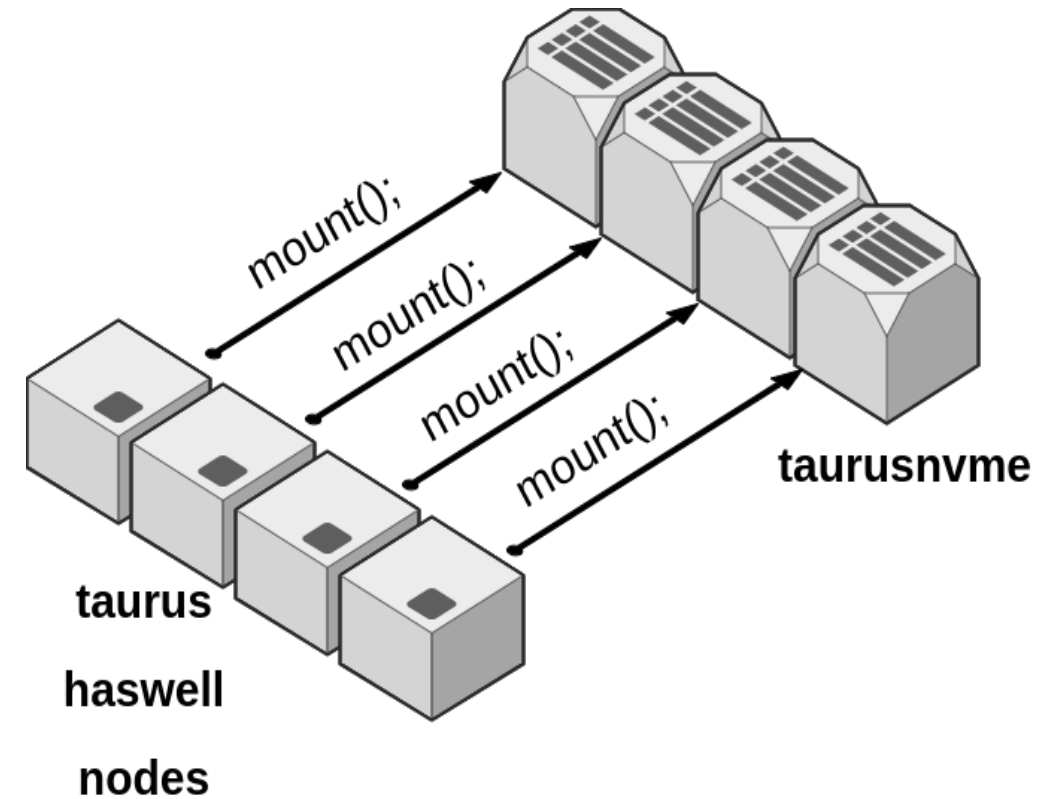Sebastian Oeste

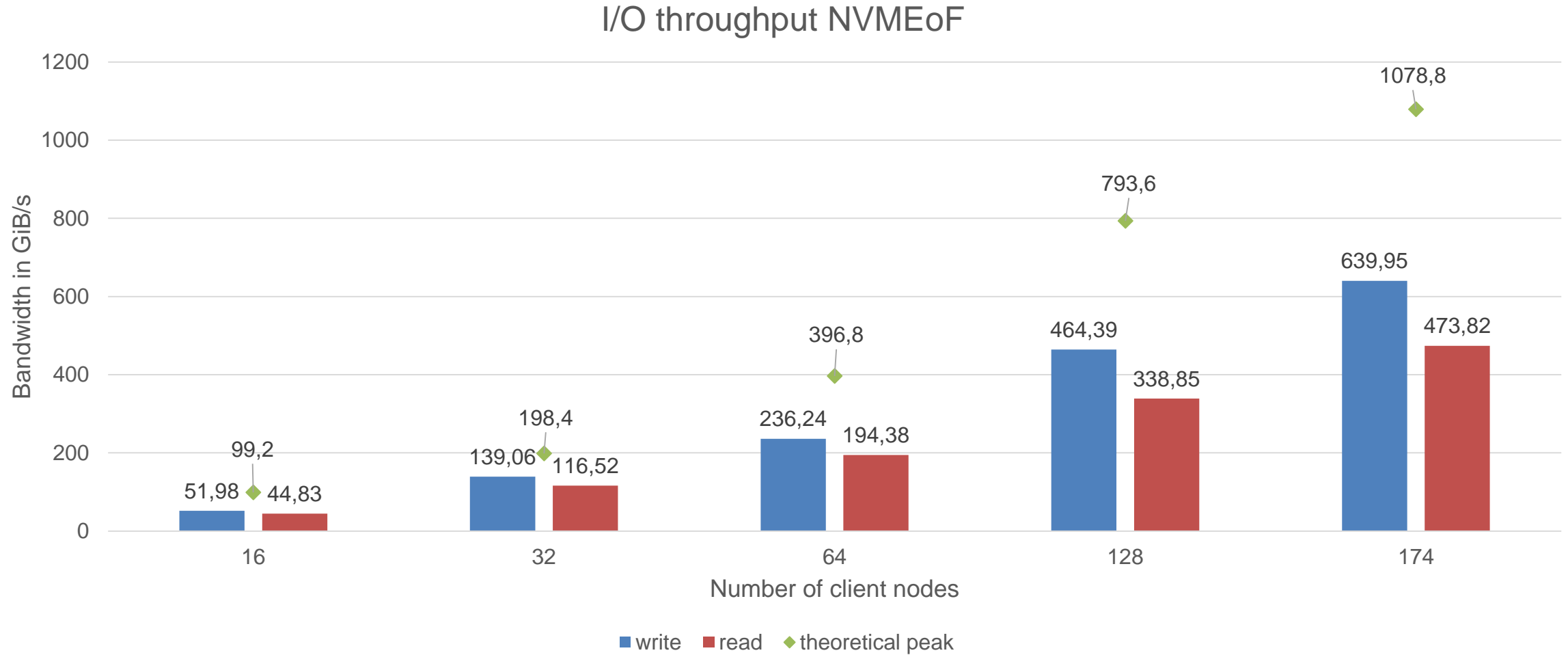# Eliminate the parallel file system as shared medium using NVME over Fabrics

- Using `nvme-connect` to connect NVME SSDs directly to compute nodes
- Server-side SSD appears as block device on compute node
- Use a local file system on that block device (e.g. ext4, XFS, ...)
- No shared view across compute nodes

TECHNISCHE UNIVERSITÄT DRESDEN

ZIH
Center for Information Services & High Performance Computing

# Throughput test with NVME over Fabrics



I/O throughput NVMEoF

# Need a measurement for the quality of network routes

- Infiniband routes are managed by the subnet manager (openSM)

- Routes may change over time (e.g. if a host crashes or a switch port becomes unavailable)

- In Infiniband networks each device get a GUID (Global Unique Identifier)

- Combination of GUID and port number refers to a unique physical link between two devices

- Multiple paths sharing the same physical link results in congestion

- No tool available

TECHNISCHE
UNIVERSITÄT
DRESDEN

ZIH
Center for Information Services &
High Performance Computing

# Select nodes with low overlap in paths

Build the routingtable
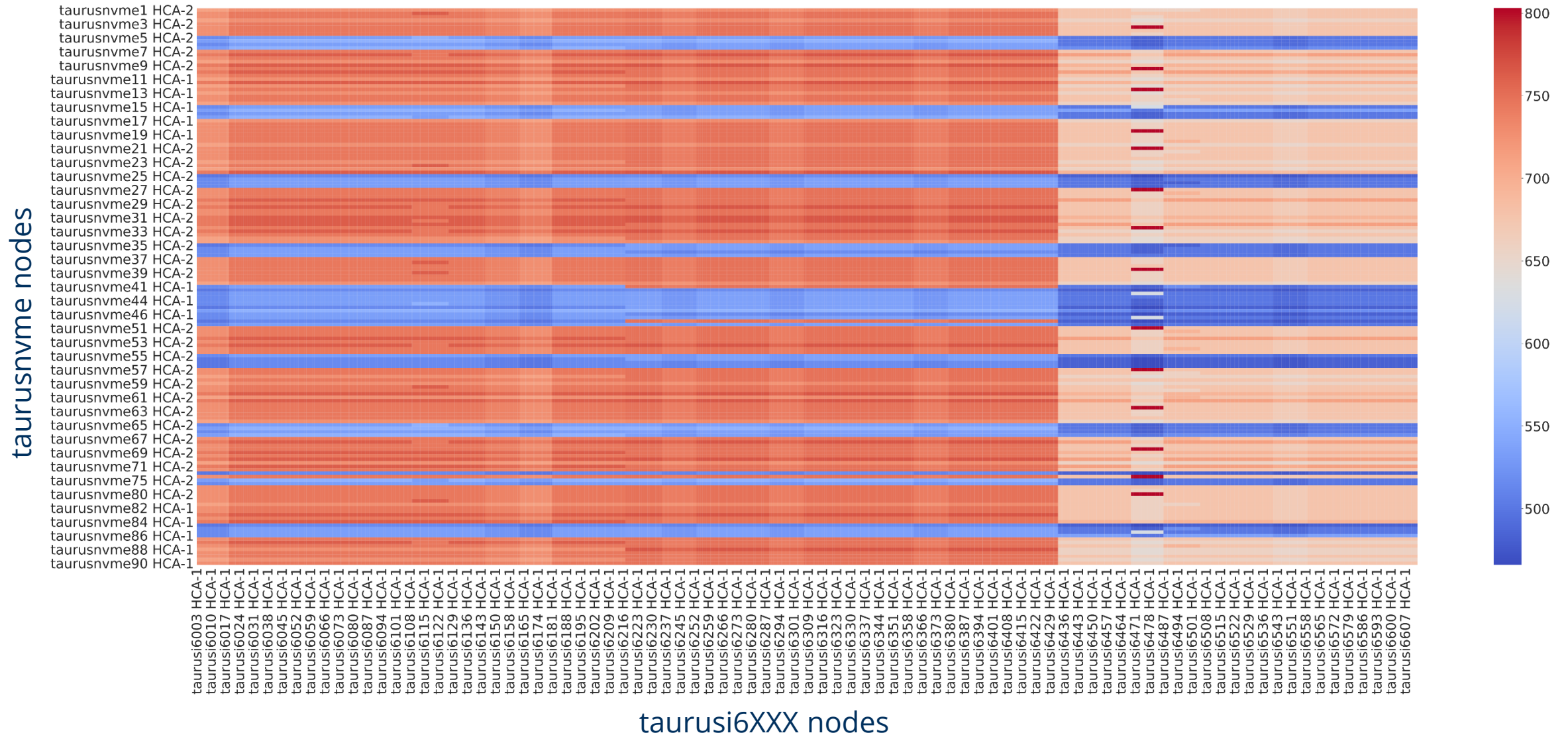
Count each GUID port combination

Sum the GUID port weight for each route

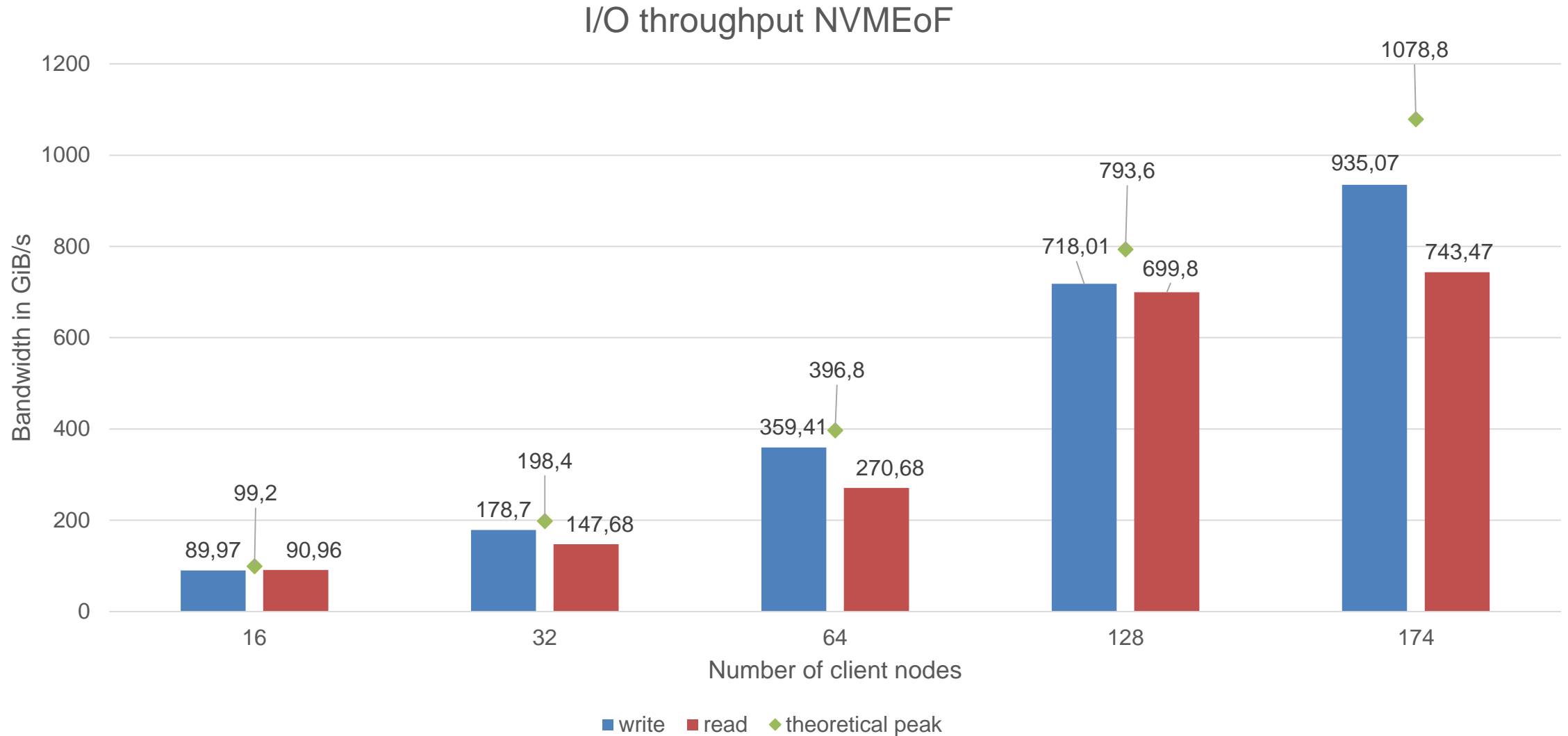Select nodes with lowest weight on routes

# Issues with the creation of the routingtable

- OFED Tool *ibtracert* to query routes between hosts
  - Query the whole fabric takes several hours

- OpenSM provides an option to dump the switch forwarding tables
  - Dump occurs every time routing changed

- Read switch fowarding tables and calculate routes manually
  - Takes ~10min for the whole fabric

TECHNISCHE
UNIVERSITÄT
DRESDEN

ZIH
Center for Information Services &
High Performance Computing

# Select nodes with lowest overlap in routes



taurusnvme nodes

taurusi6XXX nodes

TECHNISCHE
UNIVERSITÄT
DRESDEN

ZIH
Center for Information Services &
High Performance Computing

# Throughput test NVME over Fabrics with selected routes



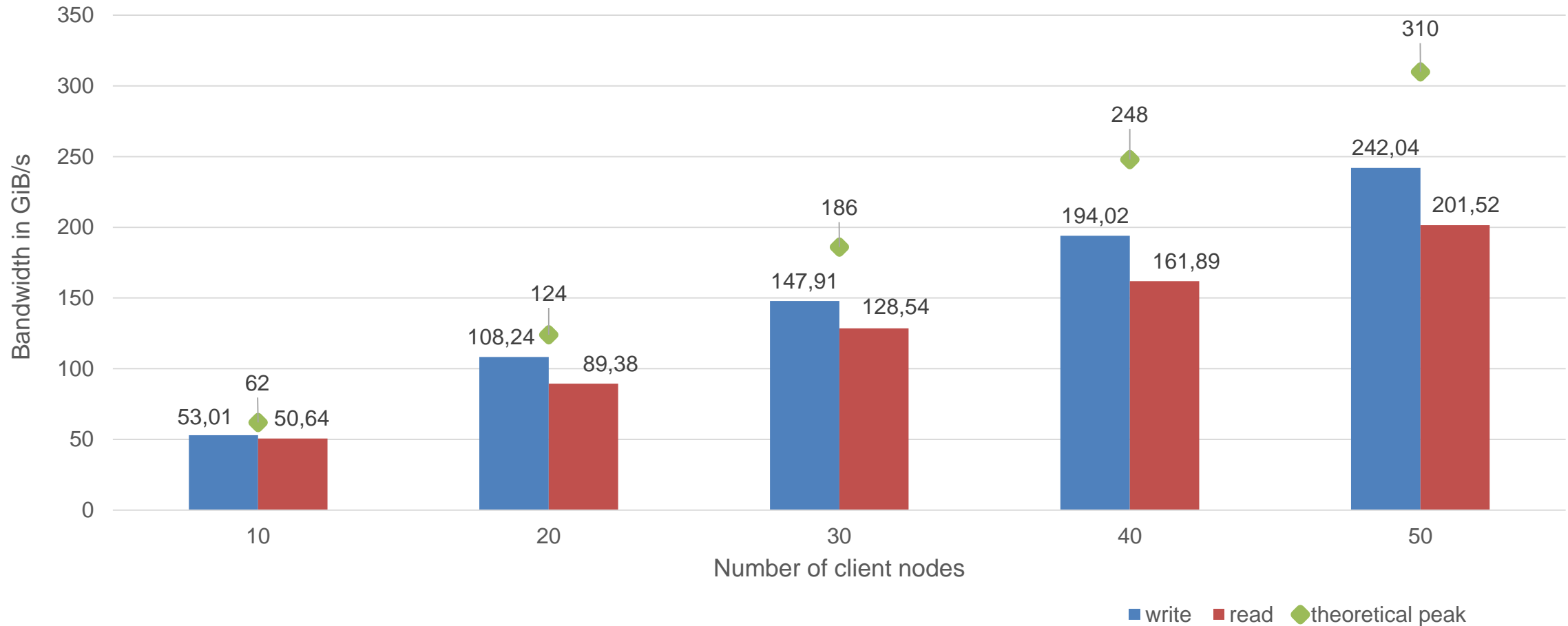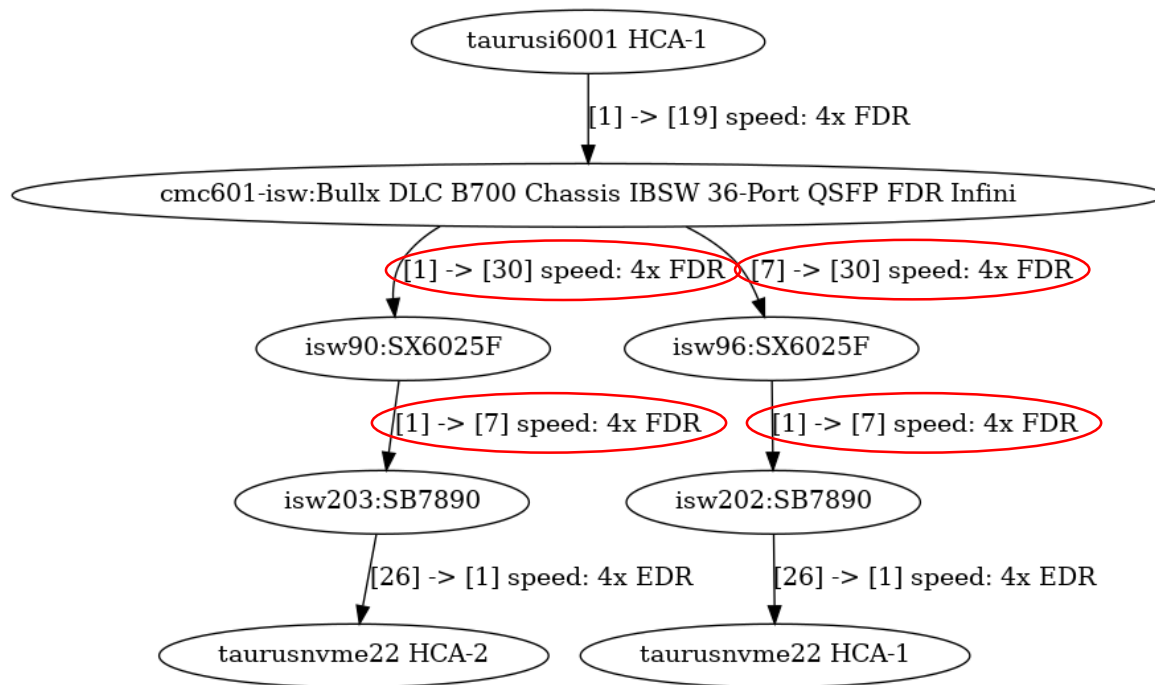I/O throughput NVMEoF

# Select routes to the parallel file system



I/O throughput BeeGFS

# Different routes for read and write direction



Write direction

Read direction

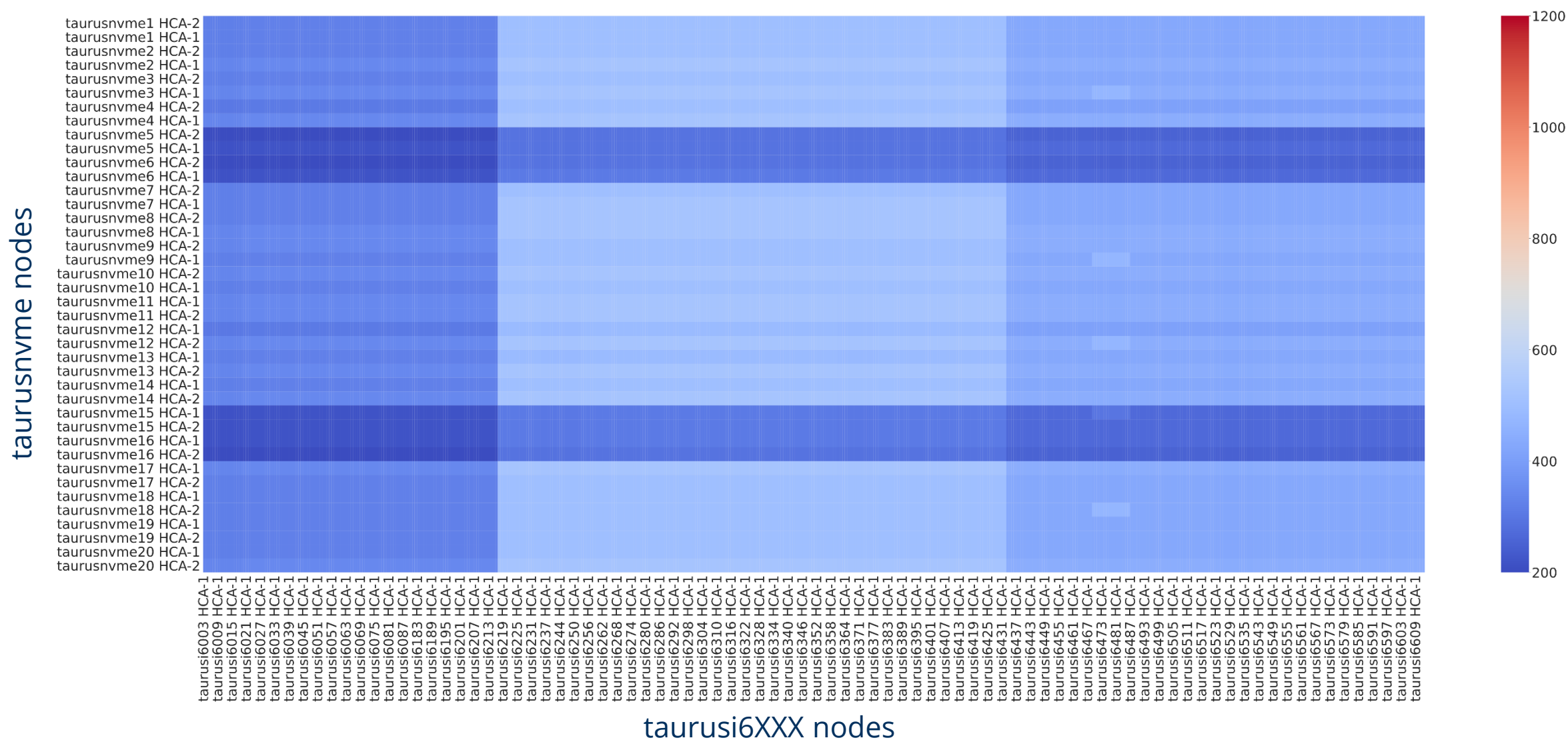# Routes for read direction

# Routes for write direction

# Summary and further work

Summary:
- Congestion on the network has a significant impact on I/O performance
- Select nodes with lower overlap in their routes can improve I/O performance
- Tool that weights paths
  - Without producing load on the fabric
  - Able to monitor path changes over time
  - https://github.com/blastmaster/IBspy

Further work:
- Enable route evaluation for live jobs together with PIKA
- Slurm plugin for integration in job scheduler

TECHNISCHE
UNIVERSITÄT
DRESDEN

ZIH
Center for Information Services &
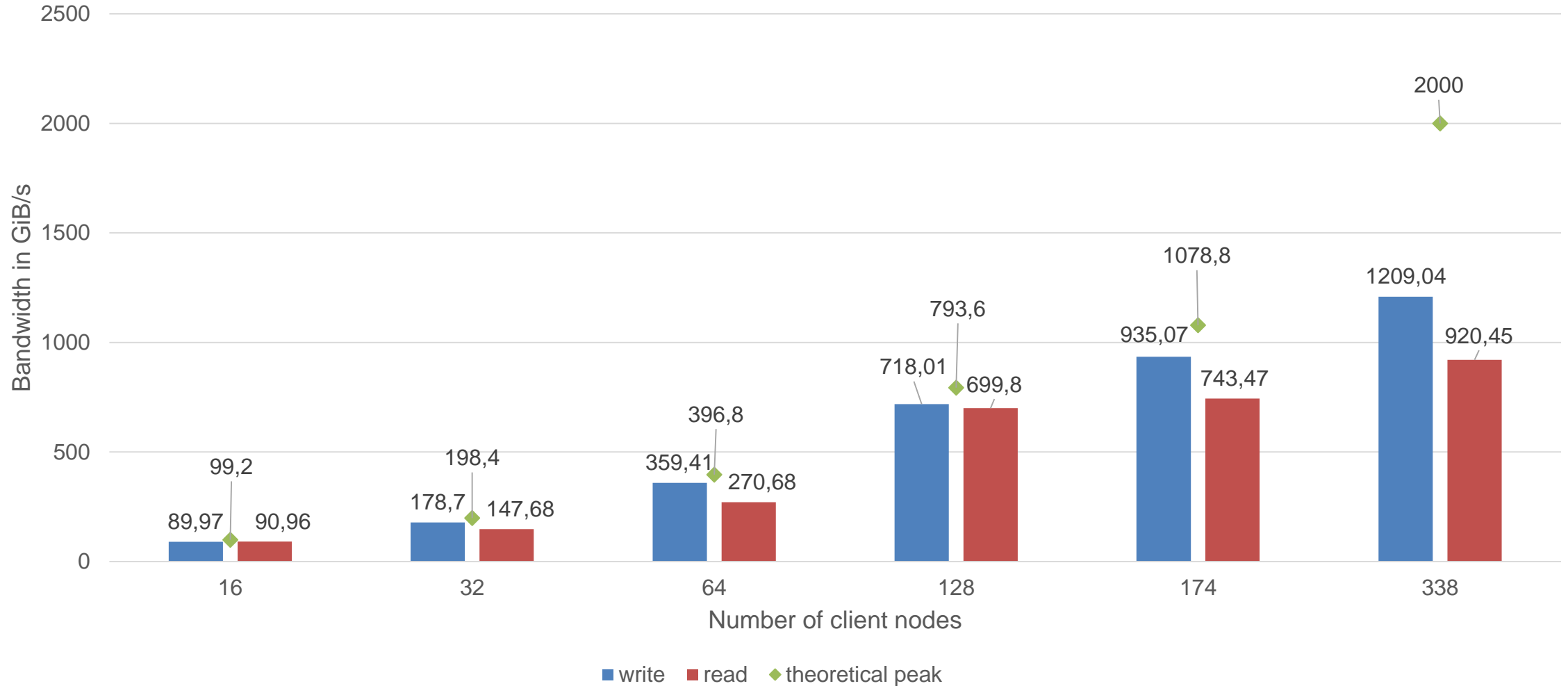High Performance Computing

# Discussion

- Are there other sides that discover similiar challenges with the network between storage and compute nodes?

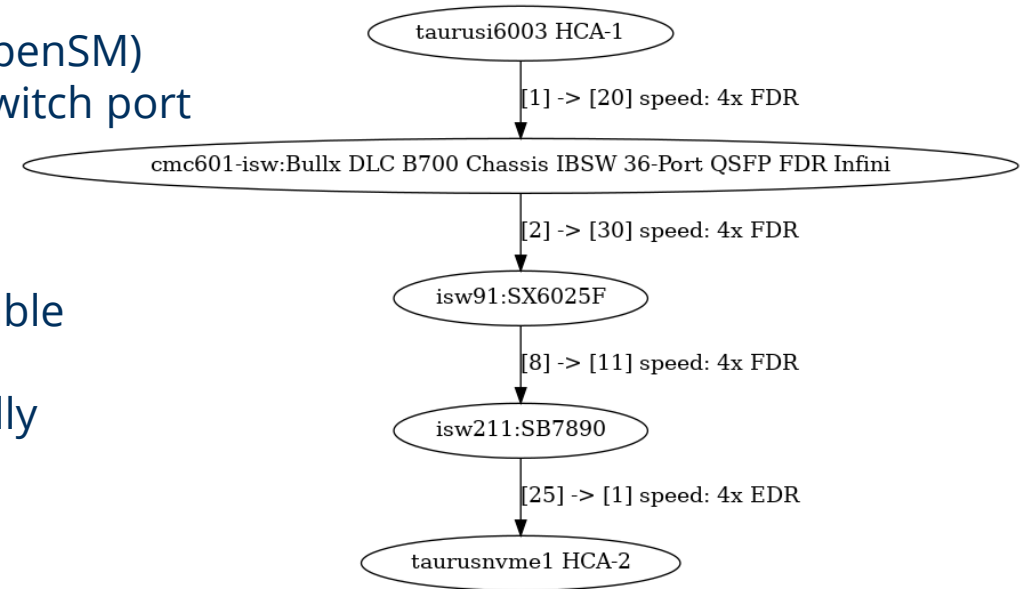- How do you monitor that?

- Possible solutions?

# Backup

# Throughput test NVME over Fabrics with selected routes



I/O throughput NVMEoF

Characterization of Infiniband routes to support data intensive I/O
Sebastian Oeste

# Collect the routing information

- Infiniband routes are managed by the subnet manager (openSM)
- Routes may change over time (e.g. if a host crashes or a switch port becomes unavailable)
- OFED Tool *ibtracert* to query routes between hosts
  - Query the whole fabric takes several hours
- OpenSM provides an option to dump switch forwarding table
  - Dump occurs every time routing changed
- Read switch fowarding tables and calculate routes manually
  - Takes ~10min for the whole fabric

```
taurusi6003 HCA-1
        |
        | [1] -> [20] speed: 4x FDR
        v
cmc601-isw:Bullx DLC B700 Chassis IBSW 36-Port QSFP FDR Infini
        |
        | [2] -> [30] speed: 4x FDR
        v
isw91:SX6025F
        |
        | [8] -> [11] speed: 4x FDR
        v
isw211:SB7890
        |
        | [25] -> [1] speed: 4x EDR
        v
taurusnvme1 HCA-2
```

```
> ibtracert -G 0x08003800013c368b 0x506b4b0300fbef06
From ca {0x08003800013c368b} portnum 1 lid 367-367 "taurusi6003 HCA-1"
[1] -> switch port {0x08003800023cceda}[20] lid 302-302 "cmc601-isw:Bullx DLC B700 Chassis IBSW 36-Port QSFP FDR Infini"
[2] -> switch port {0xf45214030369540}[30] lid 565-565 "isw91:SX6025F"
[8] -> switch port {0x248a070300bedb10}[11] lid 170-170 "isw211:SB7890"
[25] -> ca port {0x506b4b0300fbef06}[1] lid 73-73 "taurusnvme1 HCA-2"
```

TECHNISCHE
UNIVERSITÄT
DRESDEN

ZIH
Center for Information Services &
High Performance Computing

# The setup

Storage: taurusnvme nodes
- 2x EDR Infiniband (100 Gbit/s)
- 8x NVME SSD with 3TB capacity and ~3GiB/s read/write bandwidth

Client: taurus haswell nodes
- 1x FDR Infiniband (54 Gbit/s)
- 24 cores haswell CPU

Benchmark: IOR
- Using blocked sequential I/O pattern
- File per process
- 2MiB request sizes