

# Healthy GPU use together

Roel Aaij, Emily Kooistra

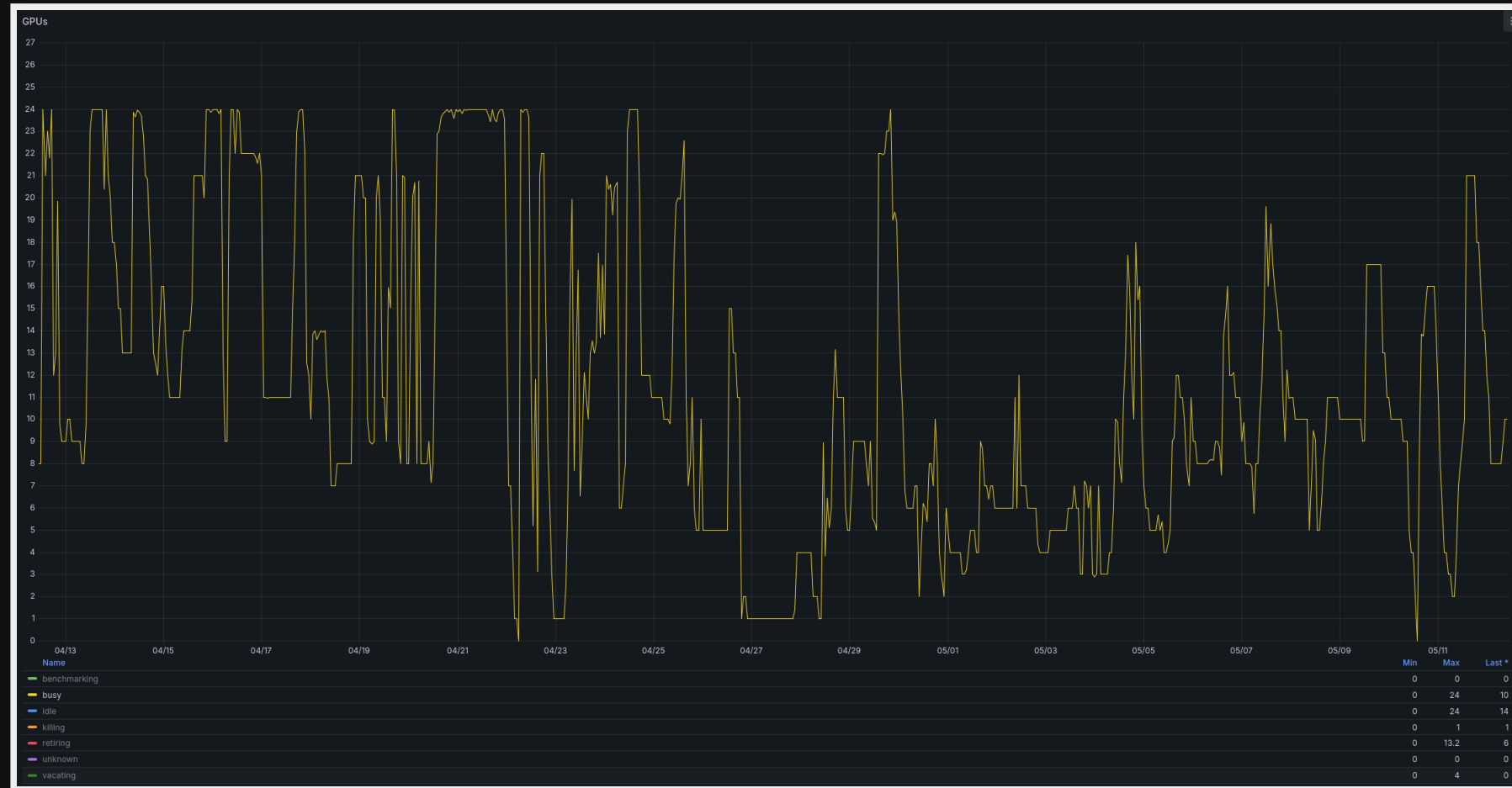
# Inventory

- 4x V100
- 20x L40S

# Goal

- High utilisation of the GPUs
- High job throughput
- More science per user
- More science across the board

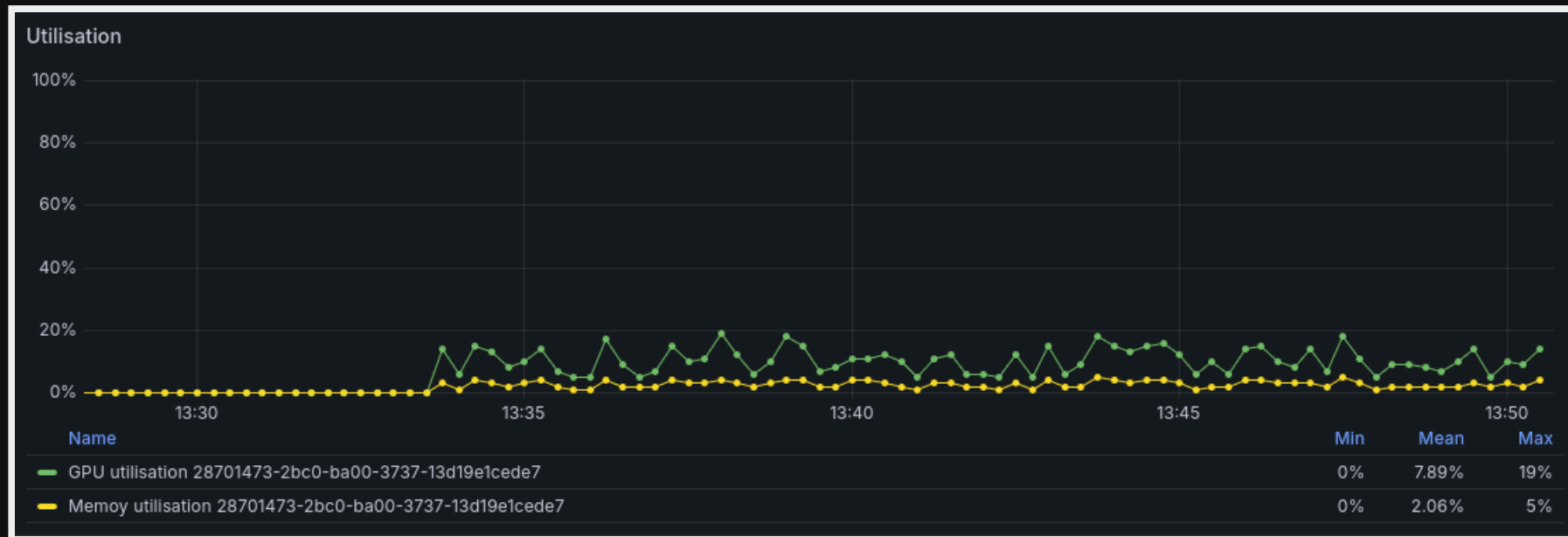
# GPUs are popular!



# Recurring issues

- Low GPU compute utilisation
- Low GPU memory footprint
- Suboptimal datamovement

# Compute utilisation



# Compute utilisation, what to do?

- Preprocess data in separate jobs
- Use a DAG if this is needed for every submission
- Run multiple GPU workloads in a single job
- Have a look at optimisation tips for e.g. [PyTorch](#), [JAX](#) or [Tensorflow](#)

# Memory utilisation



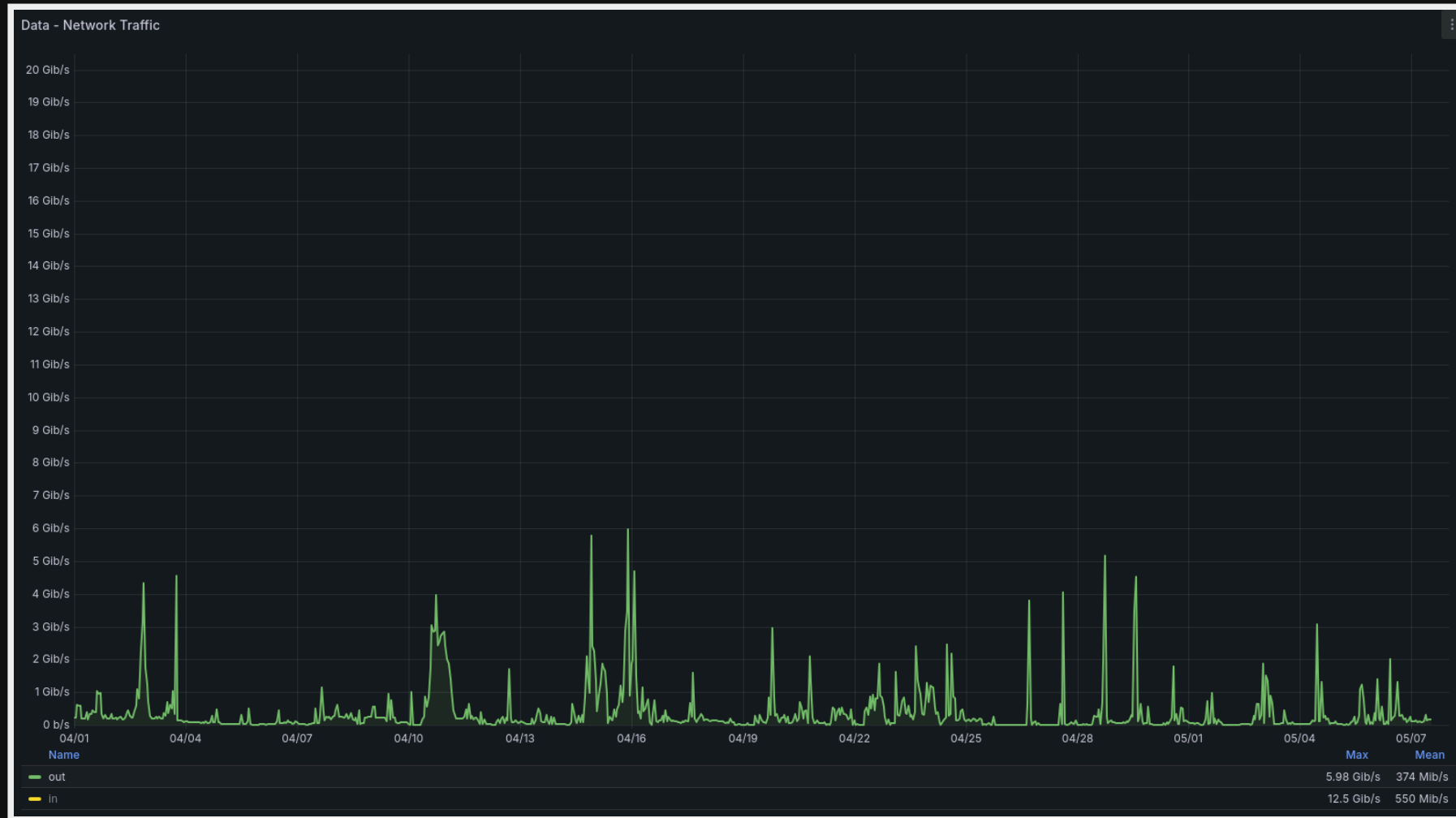
# Memory utilisation, what to do?

- *Low* memory utilisation is not really a problem in itself
- Allows multiple workloads per GPU per job, which gives higher throughput
- High memory usage can lead to lower compute utilisation
- Optimise memory allocations in e.g. **PyTorch**, **JAX** or **Tensorflow**

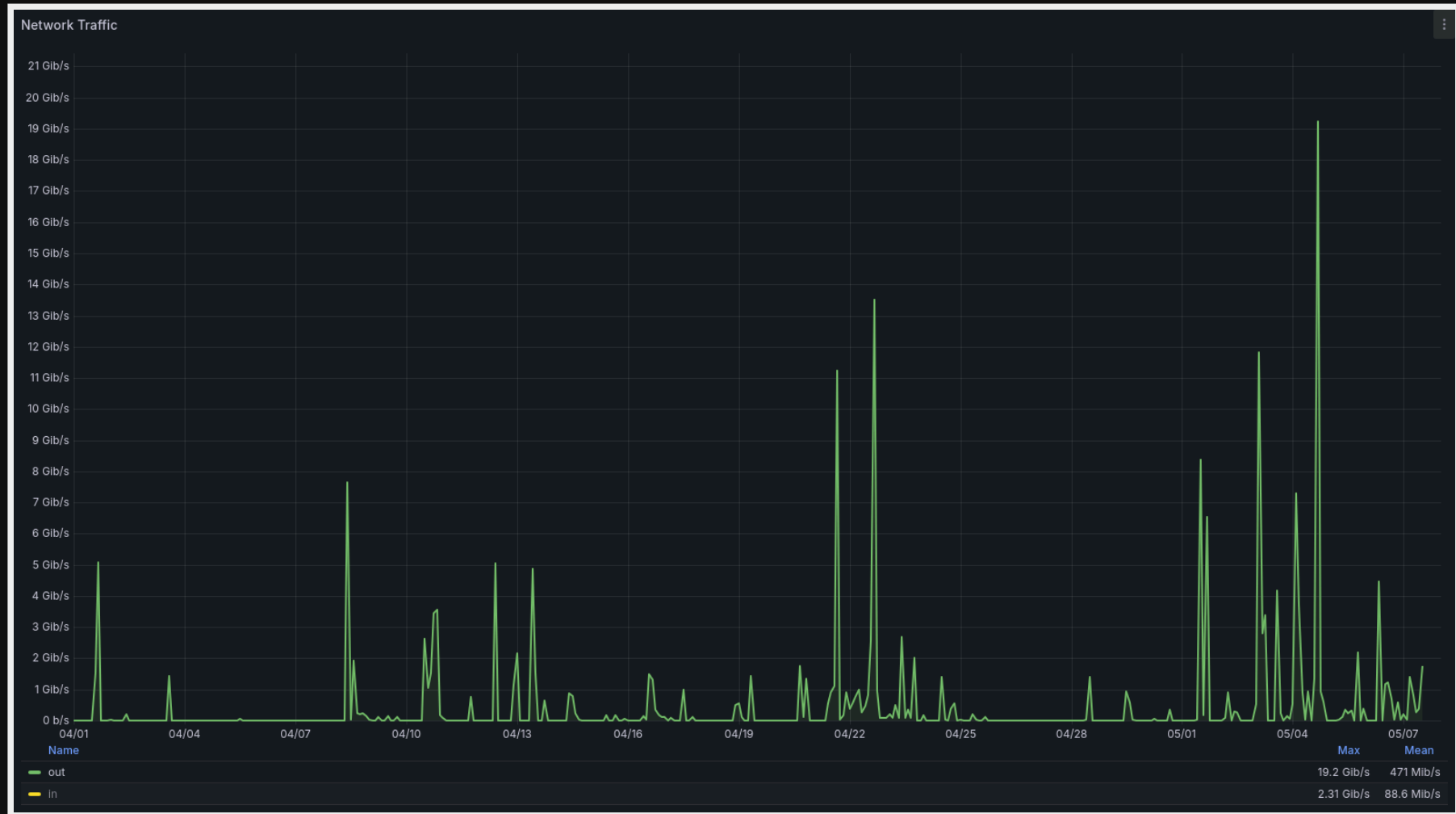
# Data Movement

- We see a lot of load on /data from GPU jobs
- Unnecessary load affects everybody
- dcache is underutilised
- Files bigger than 1GB -> dcache
- Avoid small files: bundle output as archives
- Use \$TMPDIR whenever possible
- Avoid writing logs (directly) to storage

# IO /data



# IO /dCache



# Storage usage

- Huge increase of data volume on /data
- On average 70% of all data is never used again after 7 days
- Many groups are even in the 80-90% range
- Help your group members: clean up or archive

# Talk to us

- Workflows with data bigger than 1TB?
- GPU workflows
- GPU optimisation
- Data processing pipelines
- Accessing data from CERN/other experiments
- <https://go.nikhef.nl/stbc-slot-usage>