# Analysis of Vlasiator with BSC Tools

**Valentin Seitz, BSC**

*48th VI-HPS Tuning Workshop 09.02.2026*

# Content

- The ecosystem of BSC performance tools
    - TALP
    - Extrae
    - Paraver

- Show case of a performance analysis using the BSC performance tools

# TALP

Catalan for mole
*A spy monitoring your applications resource usage.*

# TALP

- Profiler that collects POP efficiency-metrics.

- Supported programming models:
  - MPI
  - OpenMP
  - CUDA + HIP

- Provides a API to:
  - query metrics at runtime
  - annotate code regions

```
############### Monitoring Region POP Metrics ###############
### Name:                              Global
### Elapsed Time:                      10.16 s
### Host
### ----
### Parallel efficiency:               0.99
###   - MPI Parallel efficiency:       1.00
###       - Communication efficiency:  1.00
###       - Load Balance:              1.00
###            - In:                   1.00
###            - Out:                  1.00
###   - Device Offload efficiency:     0.99
###
### NVIDIA Device
### ------------
### Parallel efficiency:               0.50
###   - Load Balance:                  1.00
###   - Communication efficiency:      1.00
###   - Orchestration efficiency:      0.50
```

*Example output of TALP of an application on 2 GPUs with MPI*

EuroHPC
Joint Undertaking

# Extrae

Spanish "Extracting"
*The tracer collecting the information from hardware and the programming models*

# Extrae

**No need to recompile or relink**

- Trace producer that is transparent to the application
- Parallel programming models:
  - MPI, OpenMP, pthreads, OmpSs, CUDA, HIP
- Hardware counters (PAPI)
- Link to source:
  - Callstack at MPI routines
  - OpenMP outlined routines
  - Selected user functions (Dyninst)
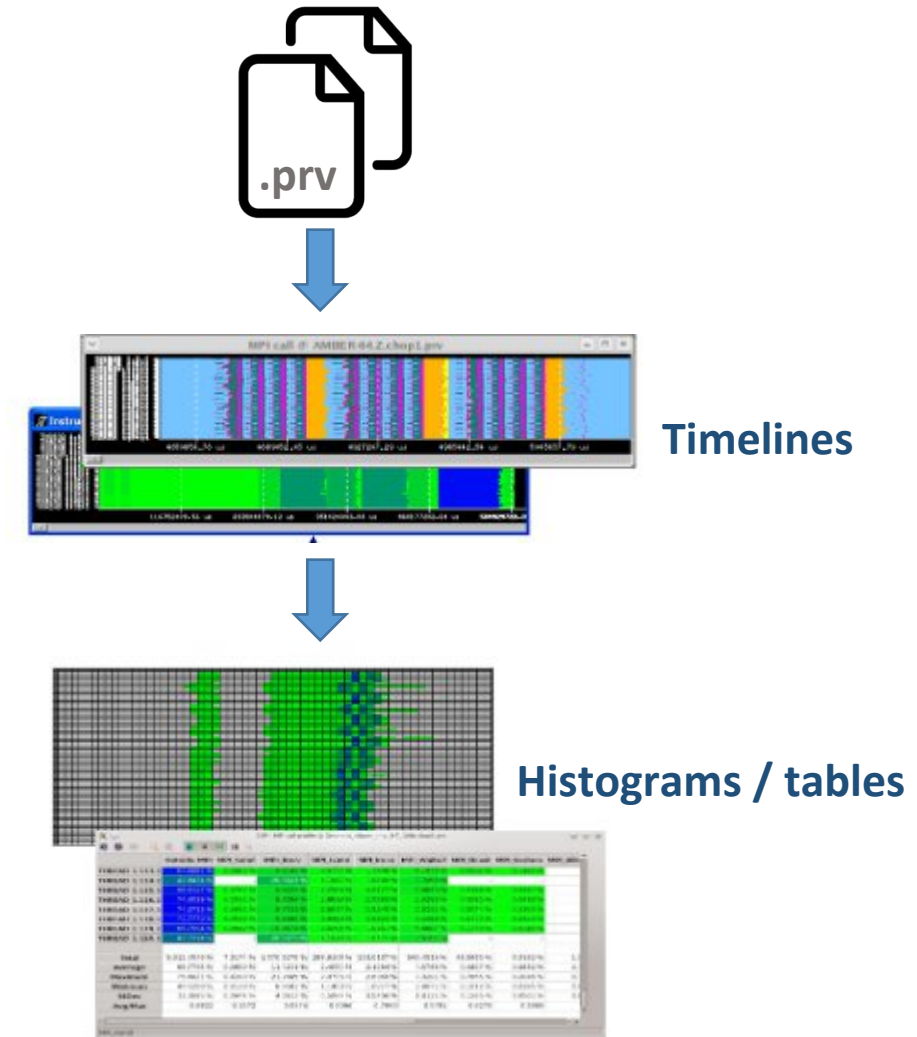- Periodic sampling
- API to annotate your code with custom events

# Paraver

Spanish "to see"
*The flexible trace viewer visualizing the data.*

# Paraver

- (Performance) Data Visualizer
  - Any kind of timestamped data
  - Trace Visualization and analysis
  - Flexible
    - No pre-assumed semantics
    - Programmable
  - 2 Kind of views:
    - Timeline
    - Histograms, 2D and 3D tables
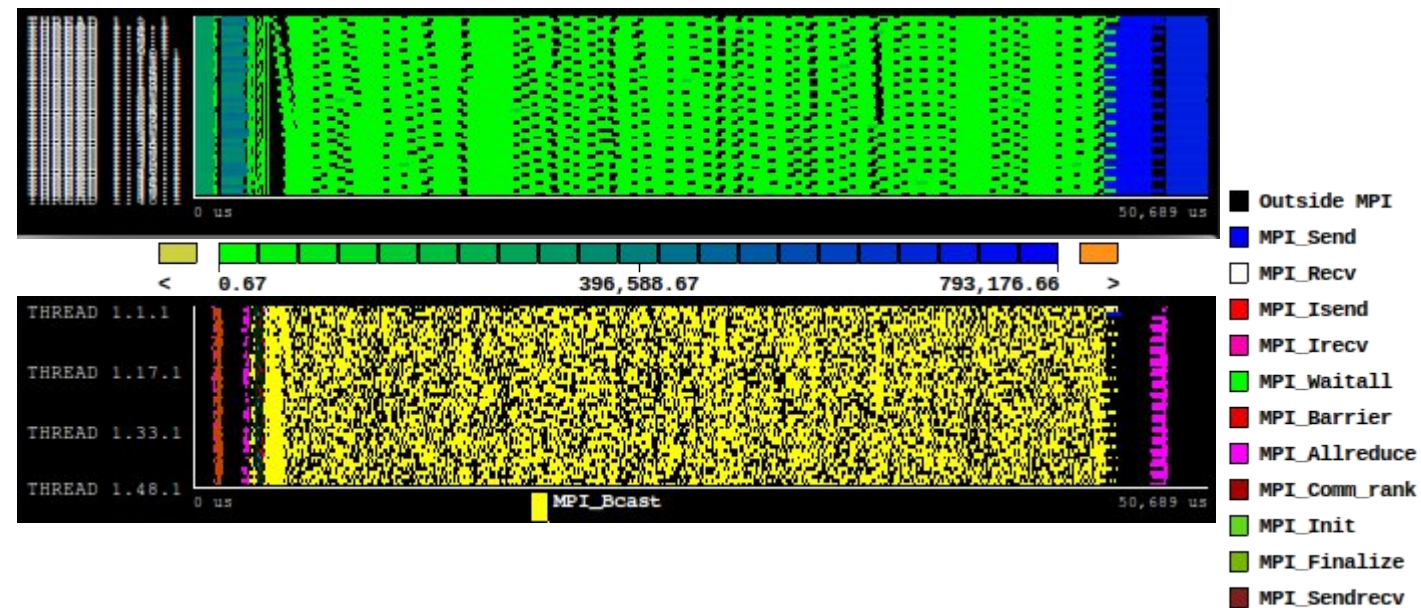


**Timelines**

**Histograms / tables**

# Paraver: Timelines

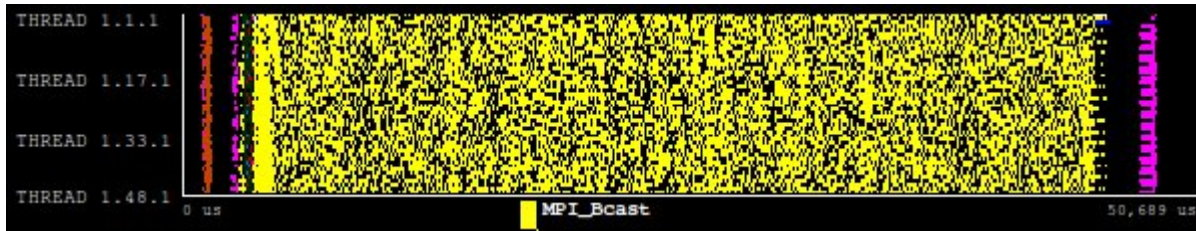Every trace is a function of time of certain data (MPI calls, CPU Frequency,..)

Different semantics require different visualizations.

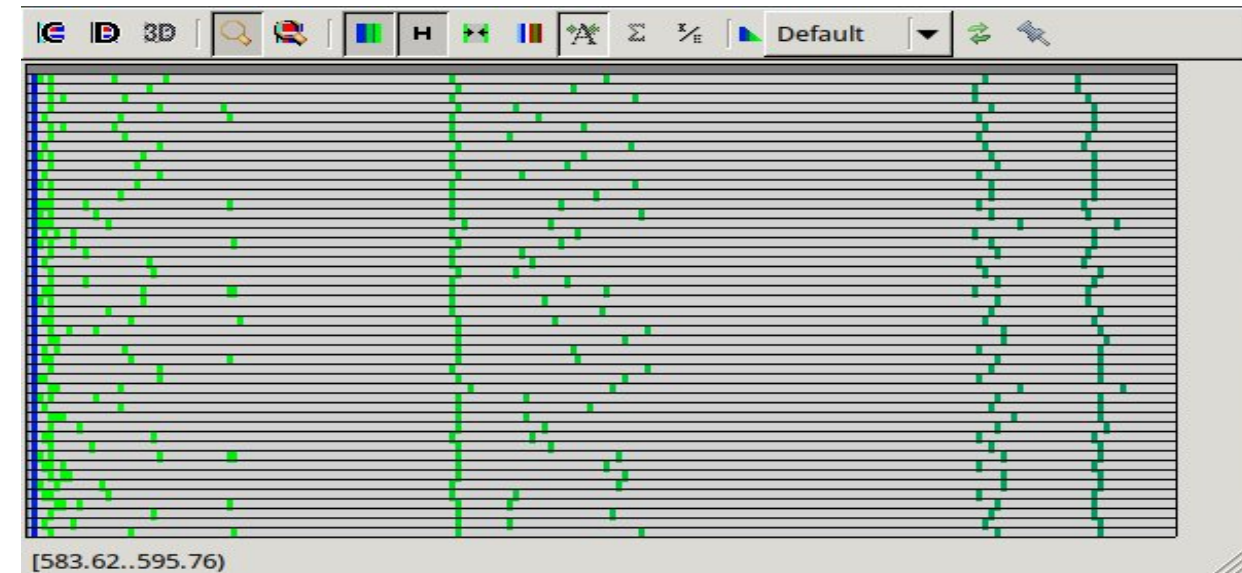Granularity between runtime calls;
IPC; Network Bandwidth ...

MPI Calls; Functions; CUDA kernels ...

# Paraver: Histograms

# Performance analysis of Vlasiator

*using BSC performance tools*

# Vlasiator

- Scientific field: near-earth plasma simulations

- Programming: C++, MPI+OpenMP

- Input: Restart from timestep 82848 of recent unpublished production run. 4.6TB

- Production run:
  - 6400 MPI x 16 OpenMP (SMT enabled) @Mahti (500 Nodes)
- Our runs:
  - 3500 MPI x 8 OpenMP (SMT disabled) (250 Nodes)
  - 7168 MPI x 4 OpenMP (SMT disabled) (256 Nodes)
  - 7000 MPI x 8 OpenMP (SMT disabled) (500 Nodes)

# Vlasiator: Timings

- Runtime on 250 nodes: 36.8s

- Runtime on 500 nodes: 17.3s
    - —> Superlinear Speed-up of 2.1


- Runtime on 1000 nodes: ~18s

Parallel efficiency 250 nodes: 0.15 (measured by TALP)
Parallel efficiency 500 nodes: 0.14 (measured by TALP)
Parallel efficiency 1000 nodes: 0.11 (measured by TALP)

# Steps

Add custom events through internal profiler

↓

Gather metrics with TALP
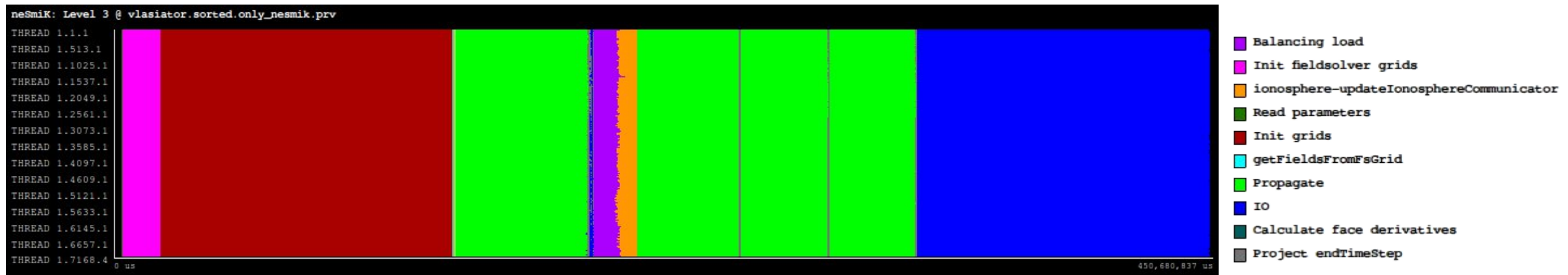
↓

Selective tracing of interesting regions

↓

Detailed performance analysis

# Custom events from internal profiler

## Main regions identified in trace:

# Metrics for **Propagate** slice

| Metrics | 250 Nodes |
|---|---|
| Global Effiency | 0.15 |
| - Parallel efficiency | 0.15 |
| -- MPI Parallel efficiency | 0.25 |
| --- MPI Communication efficiency | 0.79 |
| --- MPI Load balance | 0.32 |
| ---- MPI In-node load balance | 0.38 |
| ---- MPI Inter-node load balance | 0.83 |
| -- OpenMP Parallel efficiency | 0.36 |
| --- OpenMP Scheduling efficiency | 1 |
| --- OpenMP Load balance | 1 |
| --- OpenMP Serialization efficiency | 0.36 |
| - Computation Scalability | 1 |
| -- Instructions scaling | 1 |
| -- IPC scaling | 1 |
| -- Frequency scaling | 1 |
| Useful IPC | 1.43 |
| Frequency [GHz] | 2.91 |
| Elapsed time [s] | 226.79 |

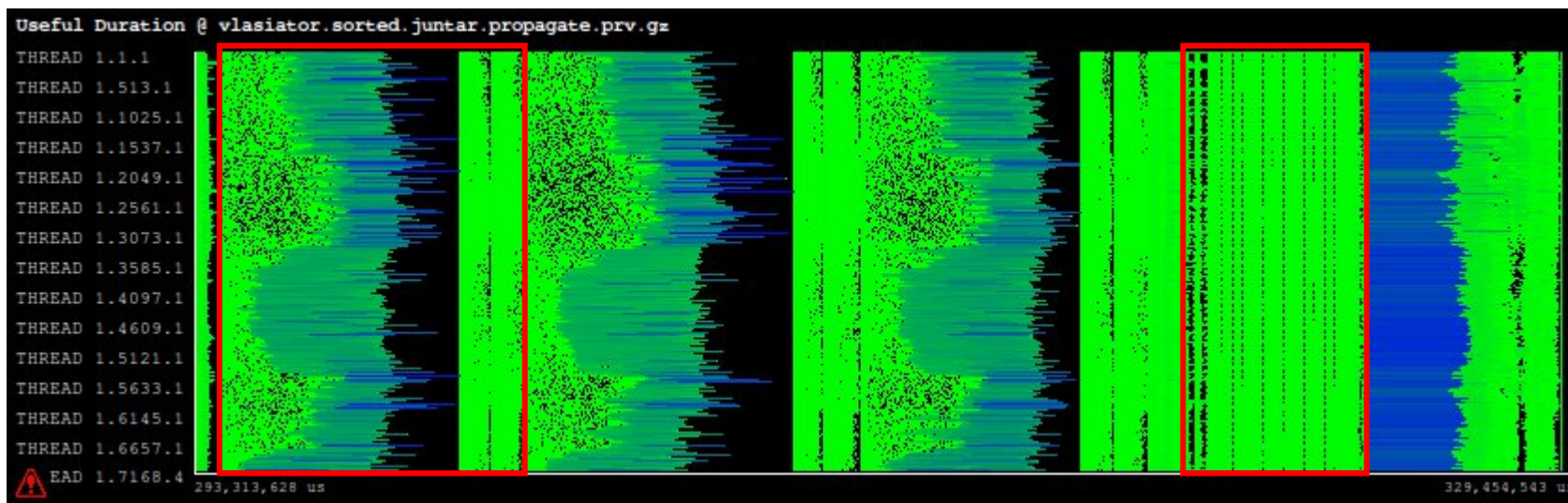- We check the metrics obtained with TALP for the relevant region

Main issues:
- MPI Load imbalance
- OpenMP serialization

# Propagate

# POP metrics Spatial-space

| Metrics | Spatial-space |
|---|---|
| Global Effiency | 0.11 |
| - Parallel efficiency | 0.11 |
| -- MPI Parallel efficiency | 0.15 |
| --- MPI Communication efficiency | 0.86 |
| --- MPI Load balance | 0.18 |
| ---- MPI In-node load balance | 0.24 |
| ---- MPI Inter-node load balance | 0.74 |
| -- OpenMP Parallel efficiency | 0.33 |
| --- OpenMP Scheduling efficiency | 1 |
| --- OpenMP Load balance | 1 |
| --- OpenMP Serialization efficiency | 0.33 |
| - Computation Scalability | 1 |
| -- Instructions scaling | 1 |
| -- IPC scaling | 1 |
| -- Frequency scaling | 1 |
| Useful IPC | 0.58 |
| Frequency [GHz] | 2.96 |
| Elapsed time [s] | 148 |

- Metrics for the Spatial-space region
- Main issues:
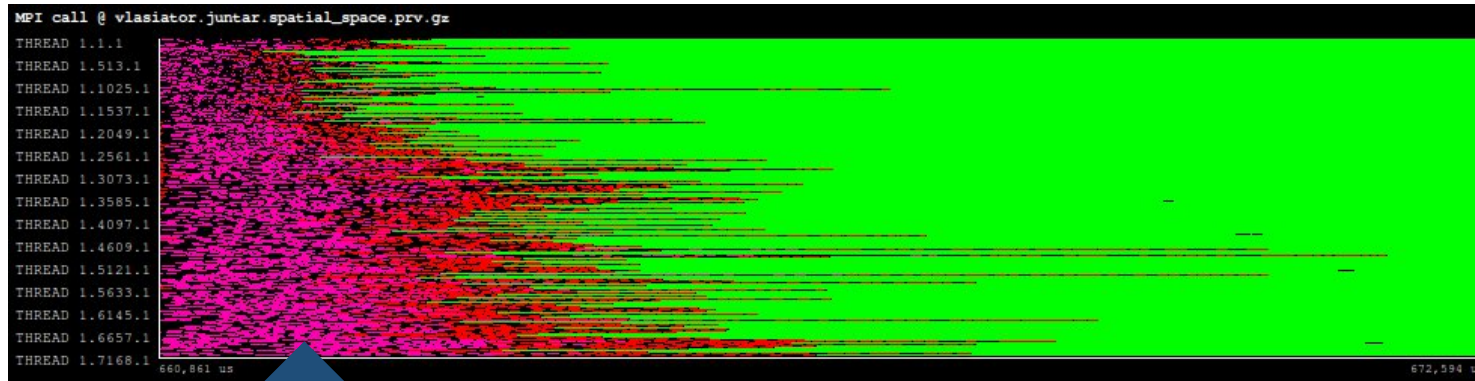  - MPI Load imbalance
  - OpenMP serialization

# Spatial-space: z-direction



transfer-stencil-data-z     compute-mapping-z     update_remote-z

# Spatial-space::transfer-stencil-data-z

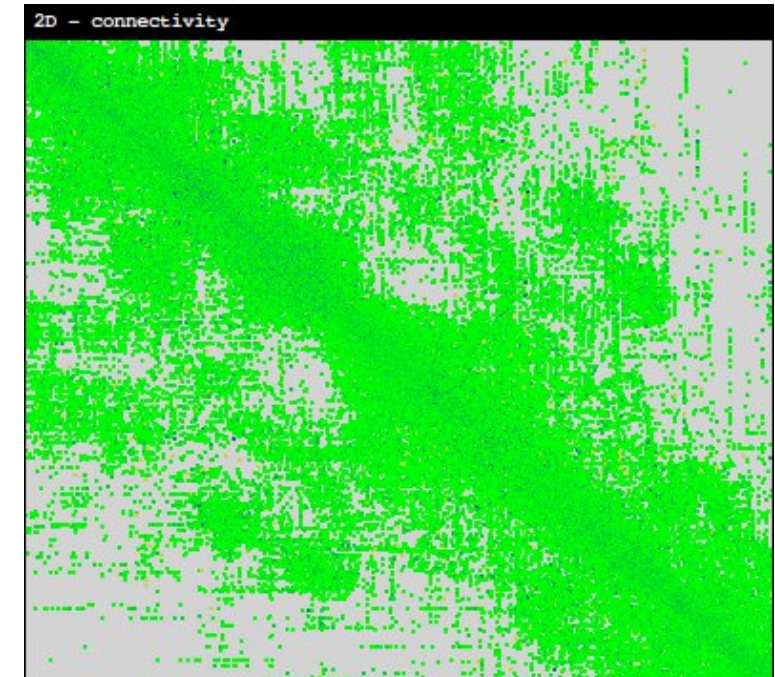# Spatial-space::transfer-stencil-data-z



Instantiation of non-blocking sends and receives

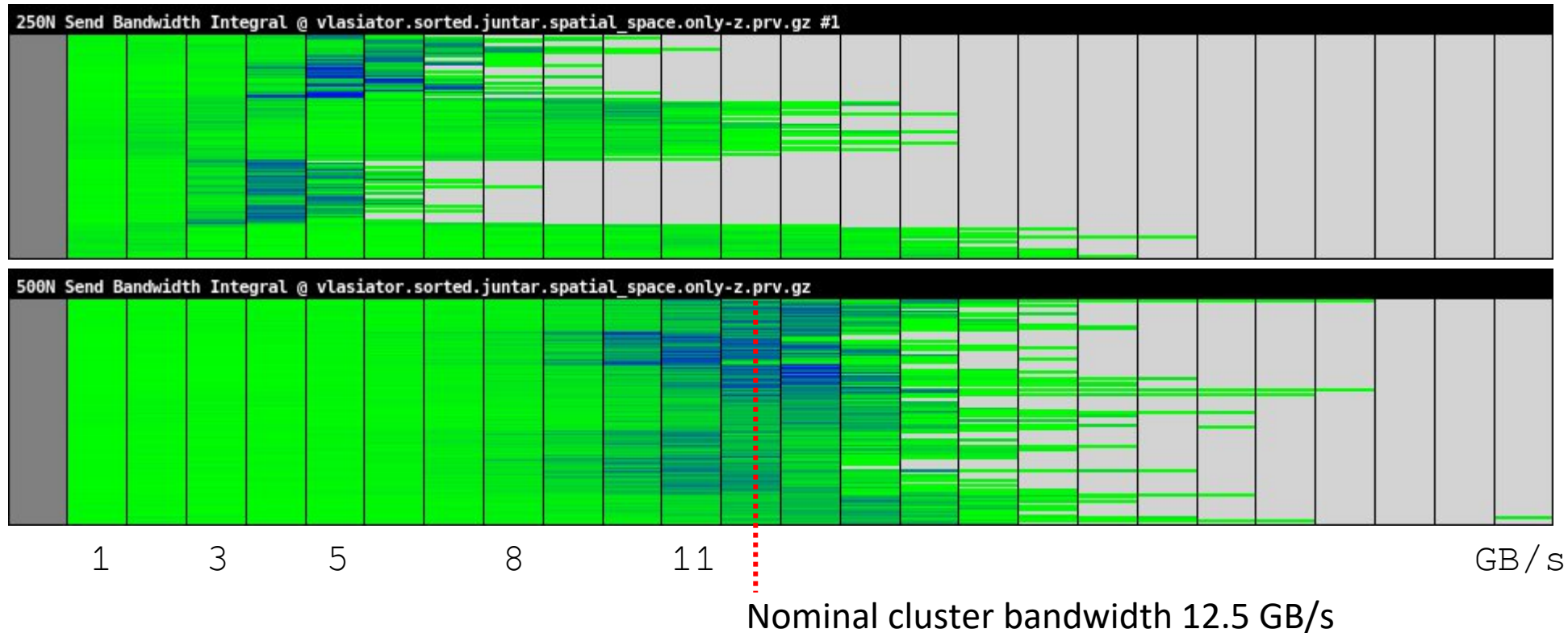Average Bytes transferred: 900MB
Maximum Bytes transferred: 4.3GB
Minimum Bytes transferred: 173MB
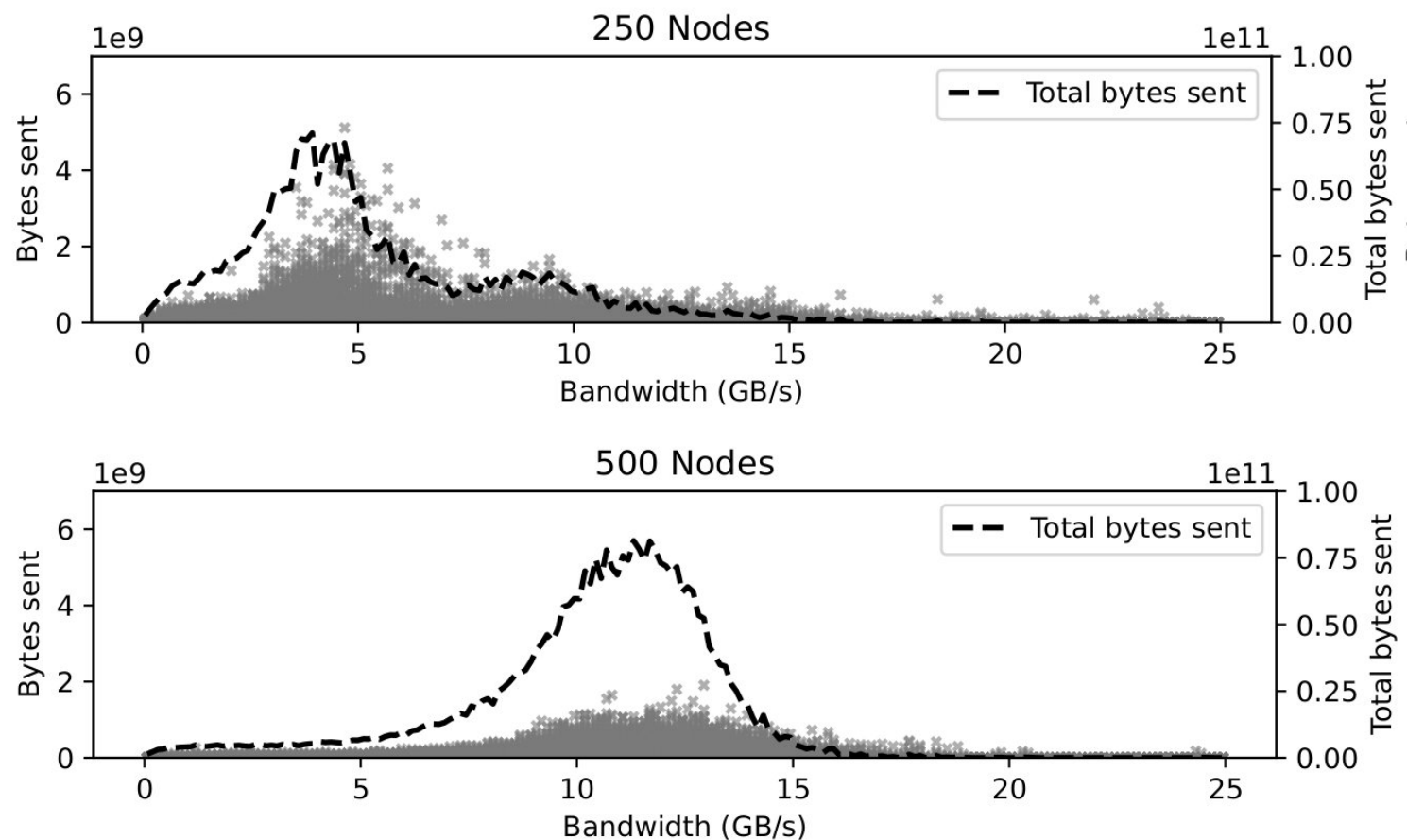Total Bytes transferred: 6.35TB

Legend:
- ■ Outside MPI
- ■ MPI_Isend
- ■ MPI_Irecv
- ■ MPI_Waitall
- ■ MPI_Barrier

Bytes sent(color) between the processes

# Spatial-space::transfer-stencil-data-z



```
250N Send Bandwidth Integral @ vlasiator.sorted.juntar.spatial_space.only-z.prv.gz #1
```

```
500N Send Bandwidth Integral @ vlasiator.sorted.juntar.spatial_space.only-z.prv.gz
```

1    3    5    8    11    GB/s

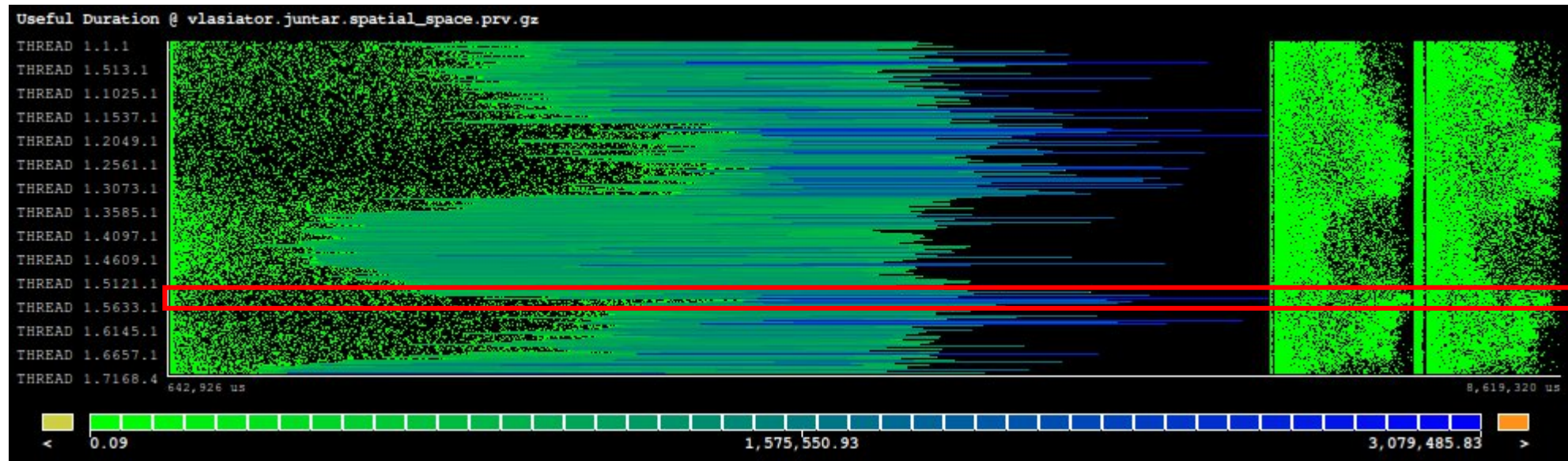Nominal cluster bandwidth 12.5 GB/s

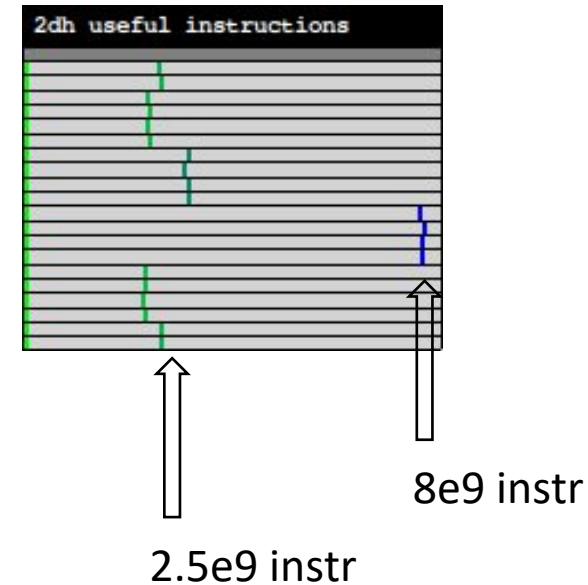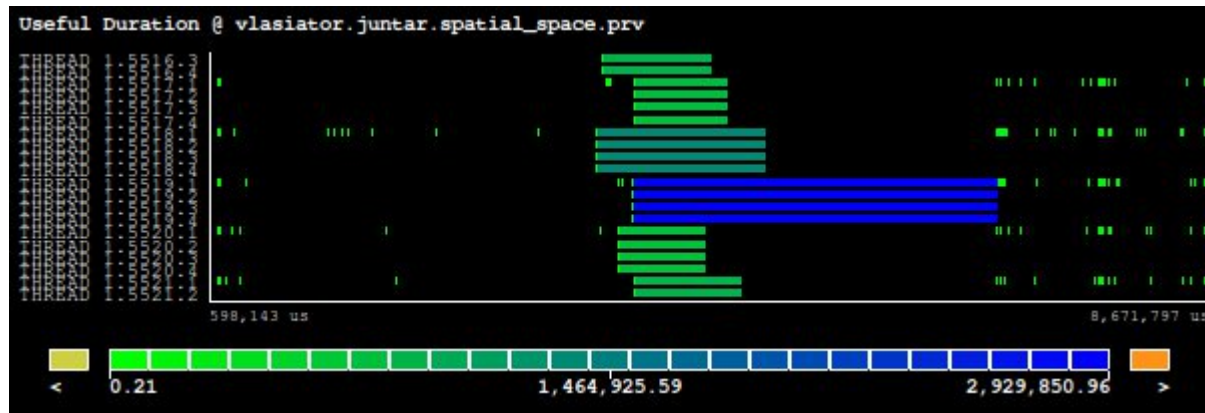-> Under utilization of network bandwidth for 250 node execution by the MPI stack.

**Send Bandwidth integral:** How much data (the more blue -> more GB) is sent at which effective bandwidth.
Colors scales are different for the two subplots

Taken from: https://doi.org/10.1016/j.procs.2025.08.237

# Spatial-space



Useful Duration @ vlasiator.juntar.spatial_space.prv.gz

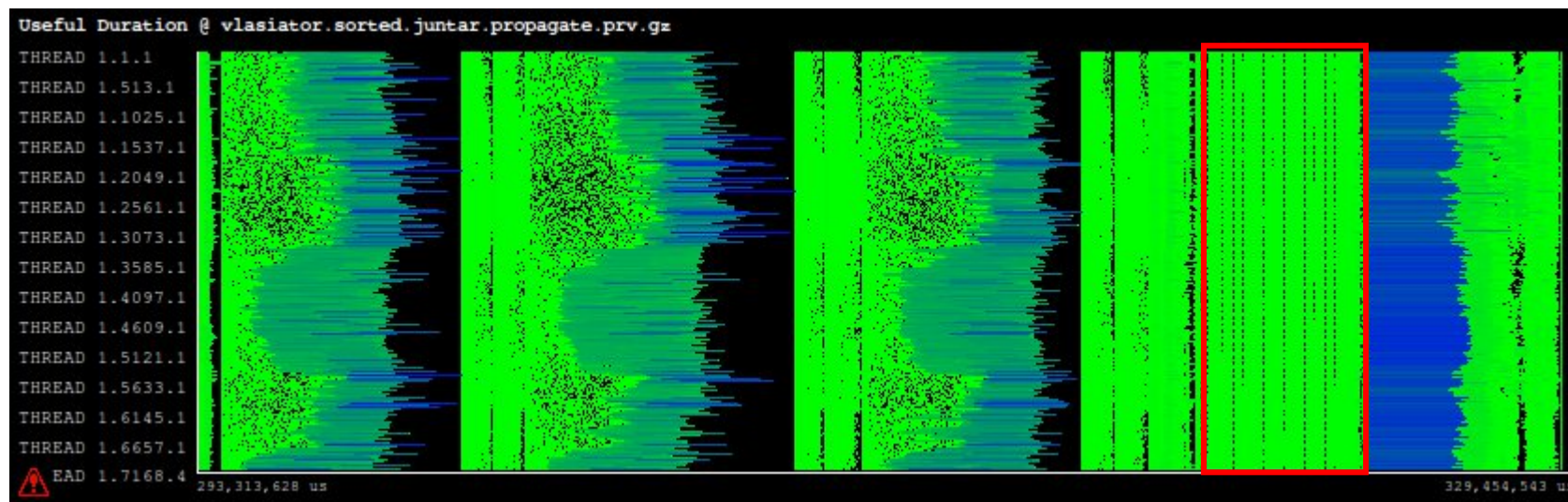# Spatial-space::compute-mapping-z



8e9 instr

2.5e9 instr

6 different processes arriving similarly after communication phase -> Then encountering load imbalance caused by instructions imbalance between the MPI processes

# Propagate



neSmiK: Level 4 @ vlasiator.sorted.juntar.propagate.prv.gz

1% — 69% — 2% — 13% — 10% — 4%

Legend:
- Update system boundaries (Vlasov pre-translation)
- Velocity-space
- Propagate Fields
- Update system boundaries (Vlasov post-translation)
- Spatial-space
- Bailout-allreduce
- Update system boundaries (Vlasov post-acceleration)
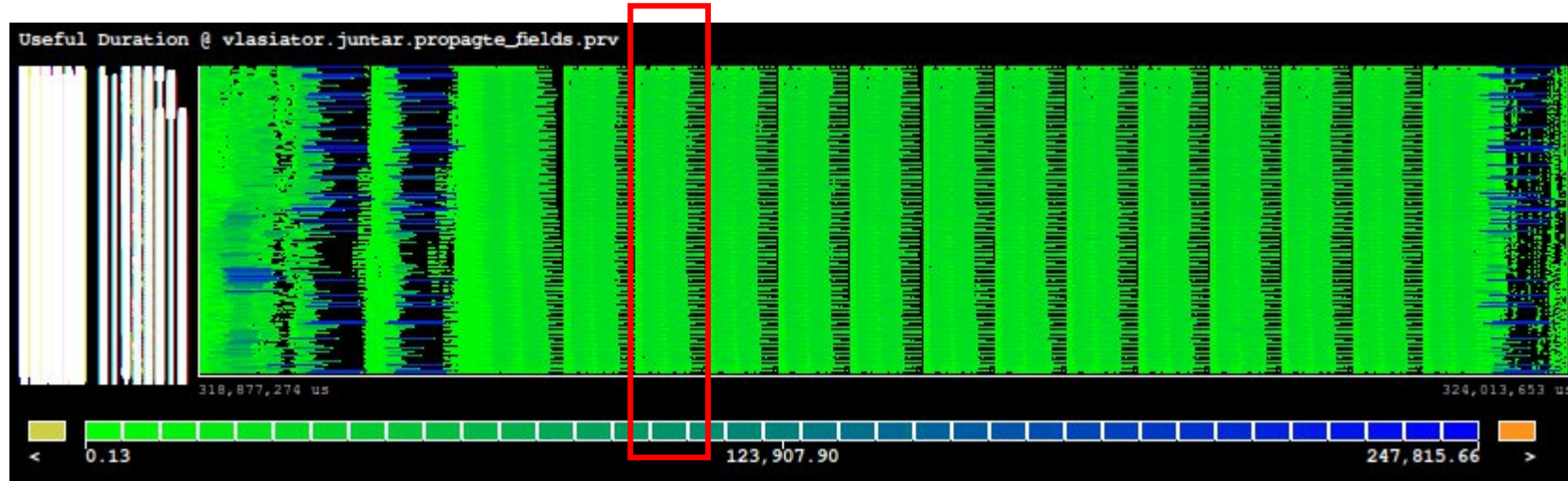- logfile-io
- write-restart

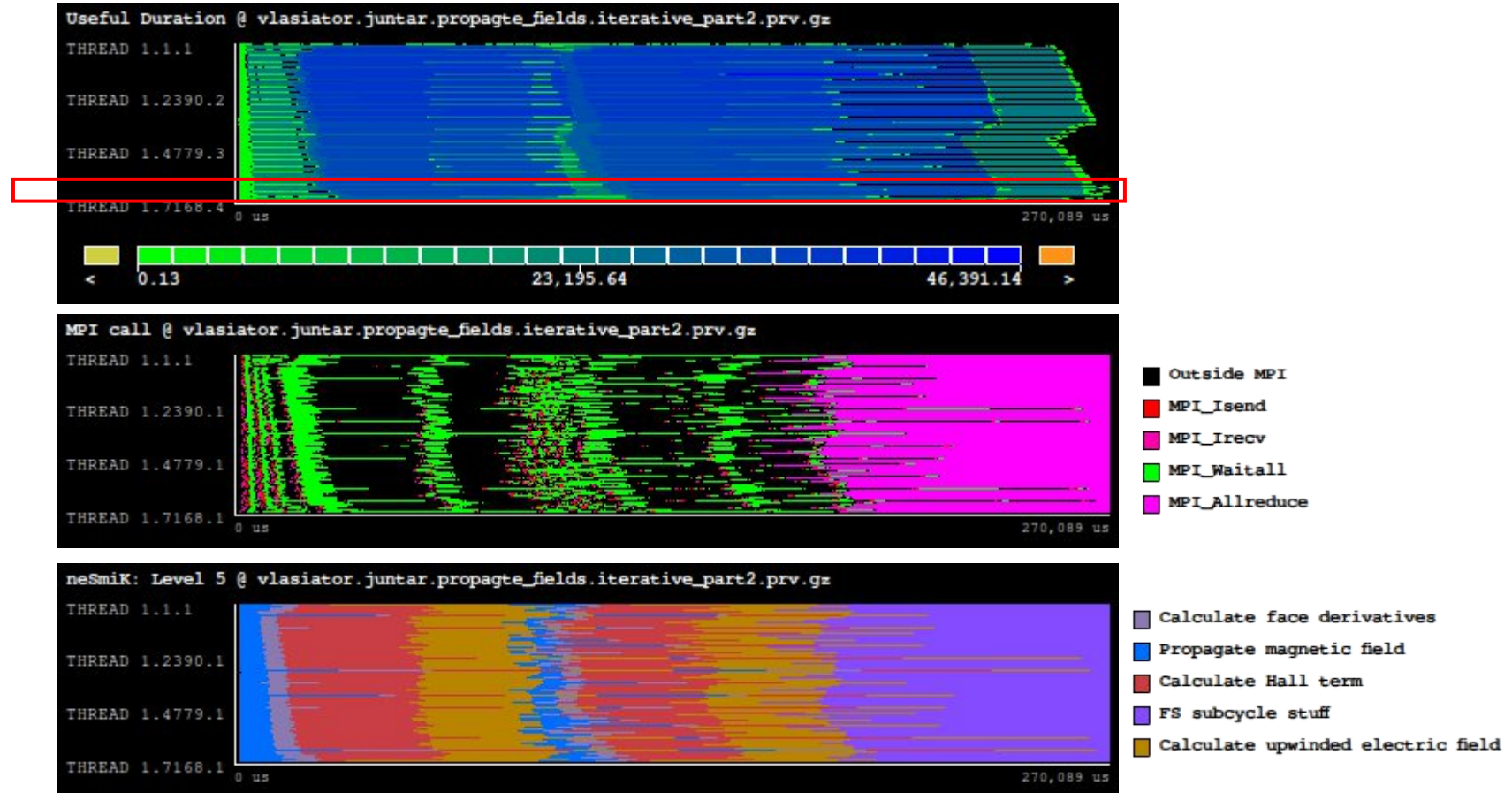Useful Duration @ vlasiator.sorted.juntar.propagate.prv.gz

48th VI-HPS Tuning Workshop

# POP metrics Propagate-fields

| Metrics | Propagate Fields |
|---|---|
| Global Effiency | 0.15 |
| - Parallel efficiency | 0.15 |
| -- MPI Parallel efficiency | 0.27 |
| --- MPI Communication efficiency | 0.85 |
| --- MPI Load balance | 0.32 |
| ---- MPI In-node load balance | 0.41 |
| ---- MPI Inter-node load balance | 0.77 |
| -- OpenMP Parallel efficiency | 0.36 |
| --- OpenMP Scheduling efficiency | 1 |
| --- OpenMP Load balance | 1 |
| --- OpenMP Serialization efficiency | 0.36 |
| - Computation Scalability | 1 |
| -- Instructions scaling | 1 |
| -- IPC scaling | 1 |
| -- Frequency scaling | 1 |
| Useful IPC | 2.84 |
| Frequency [GHz] | 2.93 |
| Elapsed time [s] | 56.27 |

- Metrics for the Propagate-fields
- Main issues:
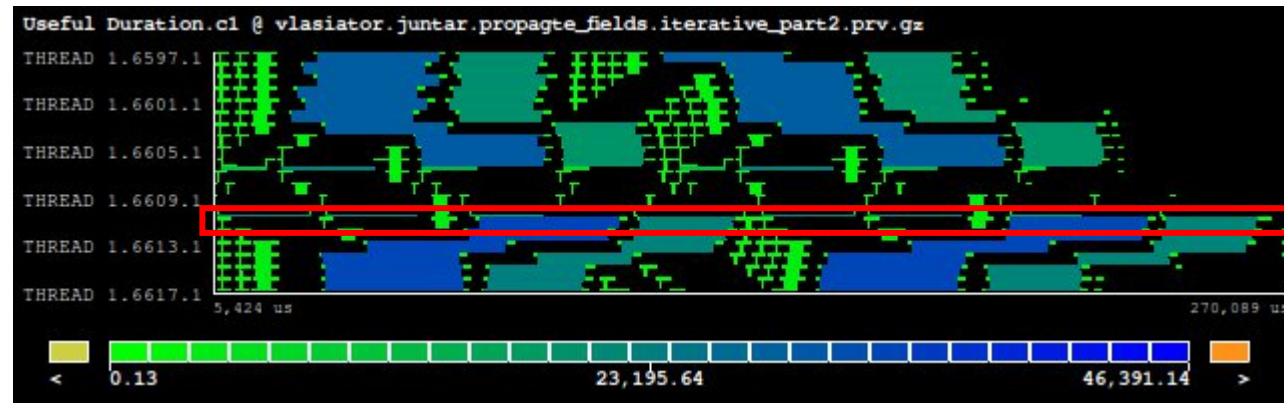  - MPI Load imbalance
  - OpenMP serialization

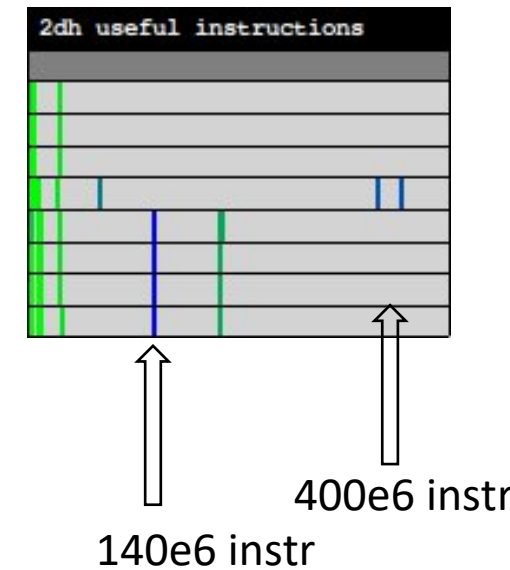# Propagate-fields



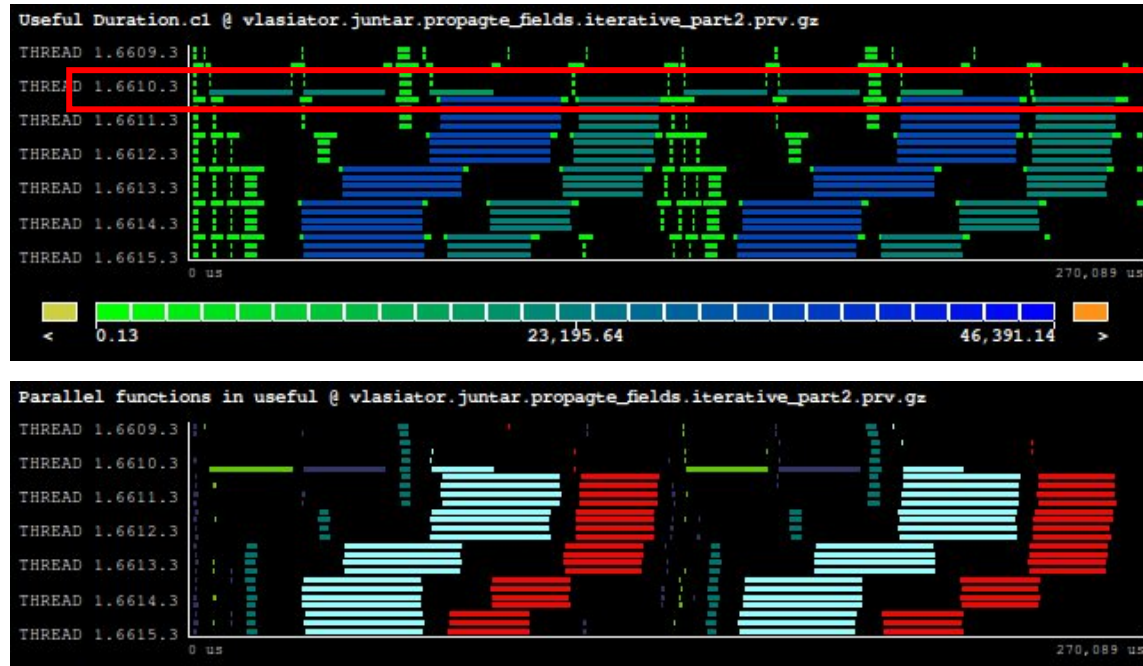Useful Duration @ vlasiator.juntar.propagte_fields.prv

318,877,274 us                                                            324,013,653 us

< 0.13                          123,907.90                          247,815.66 >

# Propagate-fields

# Propagate-fields

# Propagate-fields



Useful Duration.c1 @ vlasiator.juntar.propagte_fields.iterative_part2.prv.gz

Parallel functions in useful @ vlasiator.juntar.propagte_fields.iterative_part2.prv.gz

Parallel functions in useful

2dh useful instructions

400e6 instr

140e6 instr

- Issues in 2 parallel regions:
  - Green/purple: MPI LB, not parallelized OpenMP
  - Blue/red: MPI and OpenMP LB

# Summary

- Spatial space updates
    - Improper bandwidth utilization for 250 Node execution
    - Instruction LB problems in `compute-mapping-z`
- Field Solvers
    - LB Problems related to instructions
    - Propagate-magnetic-field: worksharing in OpenMP worsens instruction LB

# Performance Optimisation and Productivity 3
## A Centre of Excellence in HPC

EuroHPC
Joint Undertaking