



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación



**EXCELENCIA
SEVERO
OCHOA**

TALP Tracking Application Live Performance

Marta Garcia-Gasulla

10 February 2026

DLB library structure

➤ Dynamic Library

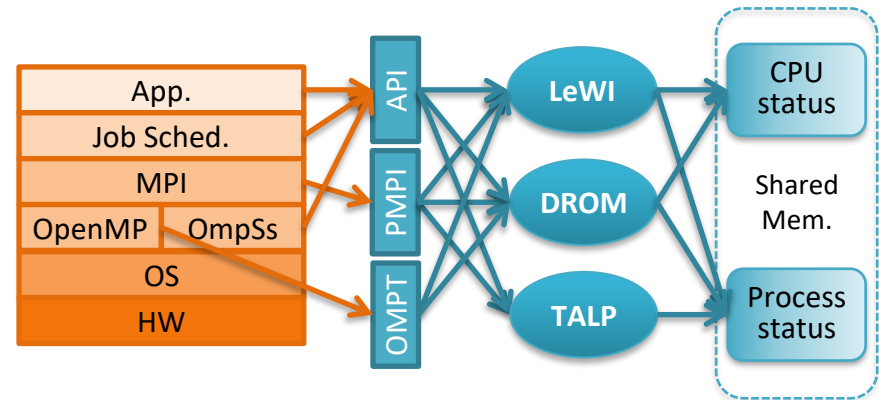
➤ Three modules:

- LeWI: For fine grain load balancing
- DROM: For coarse grain resource management
- TALP: For performance measurement

➤ Common infrastructure

- Integration with different layers of software stack
- API
- Shared memory

Three modules,
integrated but
independent

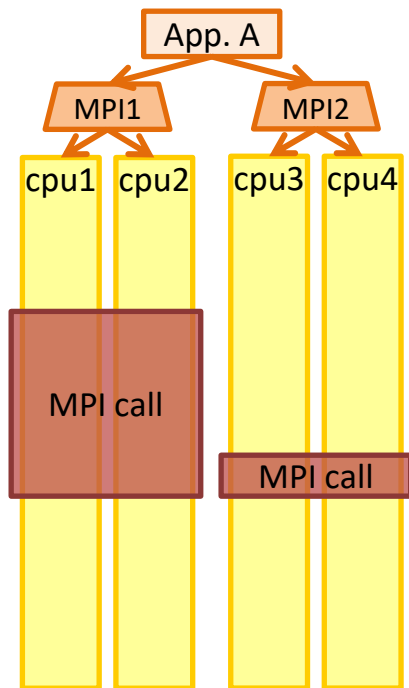


TALP: Tracking Application Live Performance

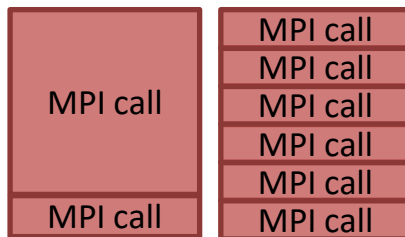
- Profiling tool with:
 - Low overhead
 - Report POP metrics
 - API to obtain metrics at runtime
 - API to instrument code and profile regions of code
- Version 3.6.0-beta1
 - MPI efficiency metrics
 - Hardware counters (cycles, instructions, IPC, and frequency)
 - OpenMP metrics
 - TALP-pages (CI/CD integration)
 - GPU efficiency metrics (NVIDIA devices)
- Work in progress
 - GPU efficiency metrics (AMD devices)
 - GPU computational metrics
 - TALP-pages support for GPU metrics

TALP a lightweight tool to Unveil Parallel Efficiency of Large Scale Executions. *In Proceedings of Performance Engineering, Modelling, Analysis, and Visualization Strategy (Permavost 2021).*

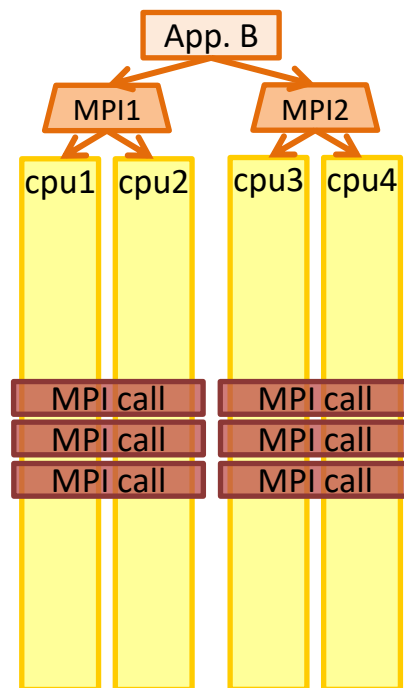
Why is more than “yet another profiling tool”?



- A profiler will report same “issue” while both cases have very different problems.



- TALP will report a low Load Balance for App. A and a low Communication efficiency for App. B



TALP MPI Metrics

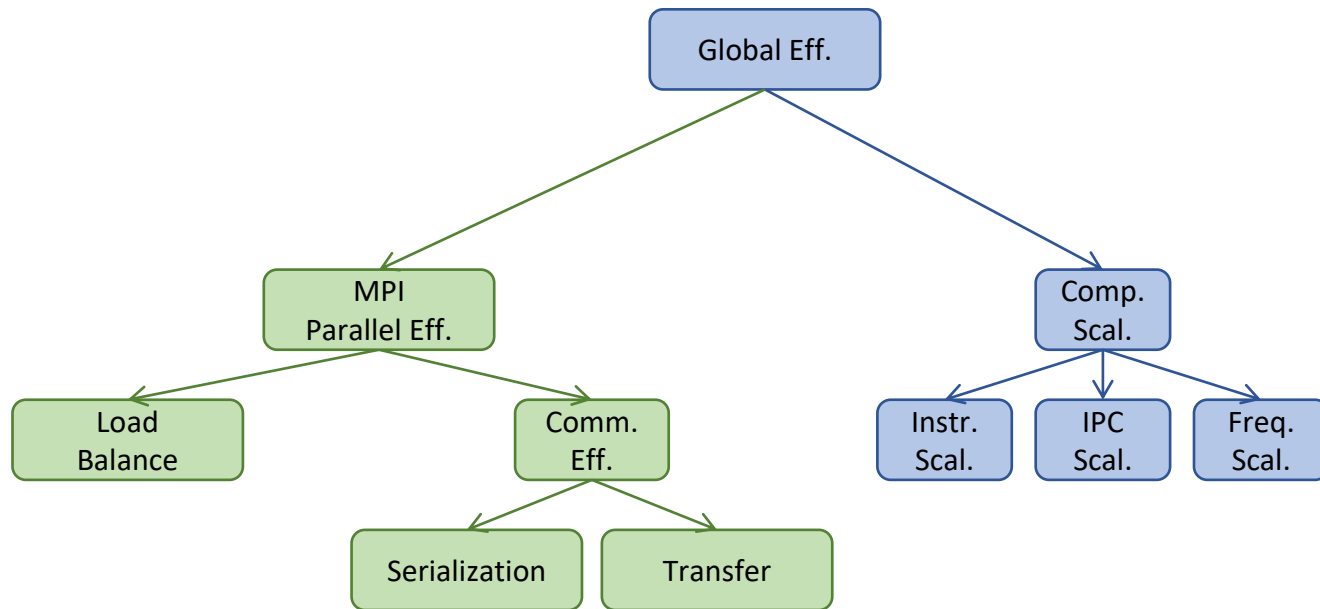


**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

MPI POP metrics

All parent metrics
are the product of
their child metrics



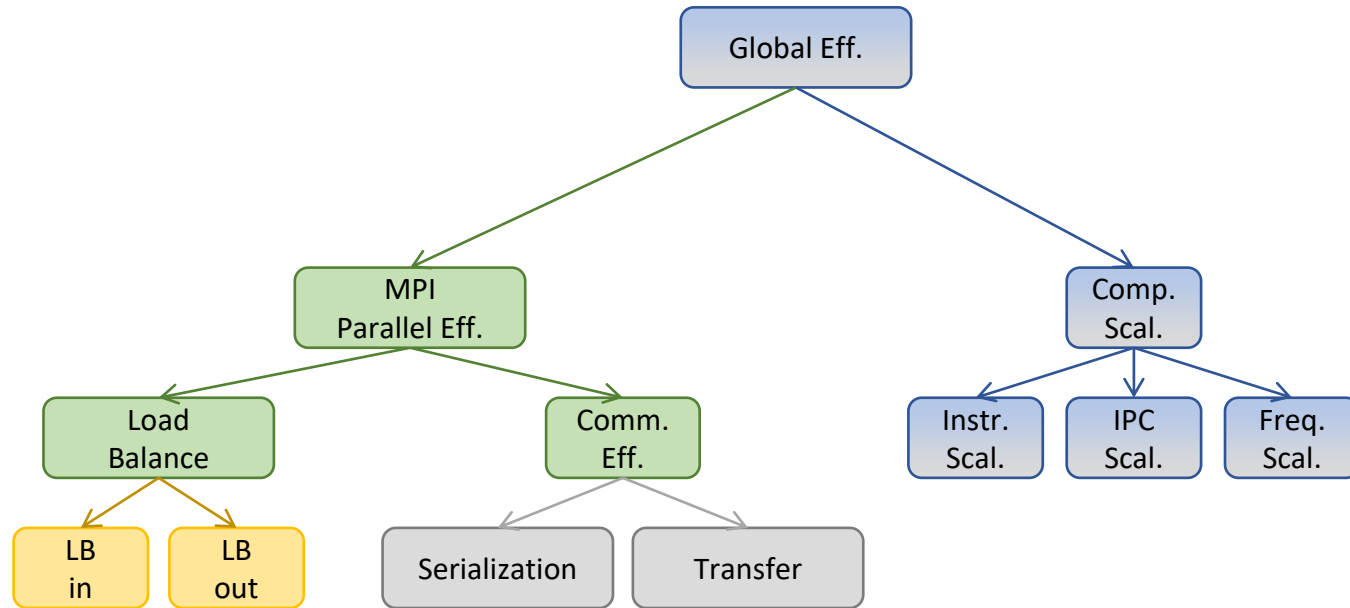
Scalability metrics
computed based on
a reference run

Efficiency metrics



TALP MPI metrics

All parent metrics
are the product of
their child metrics



Metrics not
computed by TALP

Additional metrics

Scalability metrics
computed with
several TALP runs

Efficiency metrics



Summary of metrics explanation

- **MPI parallel efficiency:** Inefficiency due to the use of parallelization. Fraction of time the resources are not being used to do useful work due to be in MPI calls.
 - **MPI Load Balance:** Inefficiency due to an uneven load distribution. Time lost due to waiting for the most loaded process.
 - **MPI Load Balance in:** Inefficiency due to an uneven load distribution inside the nodes.
 - **MPI Load Balance out:** Inefficiency due to an uneven load distribution between nodes.
 - **MPI Communication Efficiency:** Inefficiency due to the communications between processes. Time lost due to transferring data or waiting for a processes that is not the most loaded one.

TALP GPU metrics



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

TALP GPU metrics

All parent metrics
are the product of
their child metrics

Host
Global Eff.

Device
Global Eff.

Metrics not
computed by TALP

Scalability metrics
computed with
several TALP runs

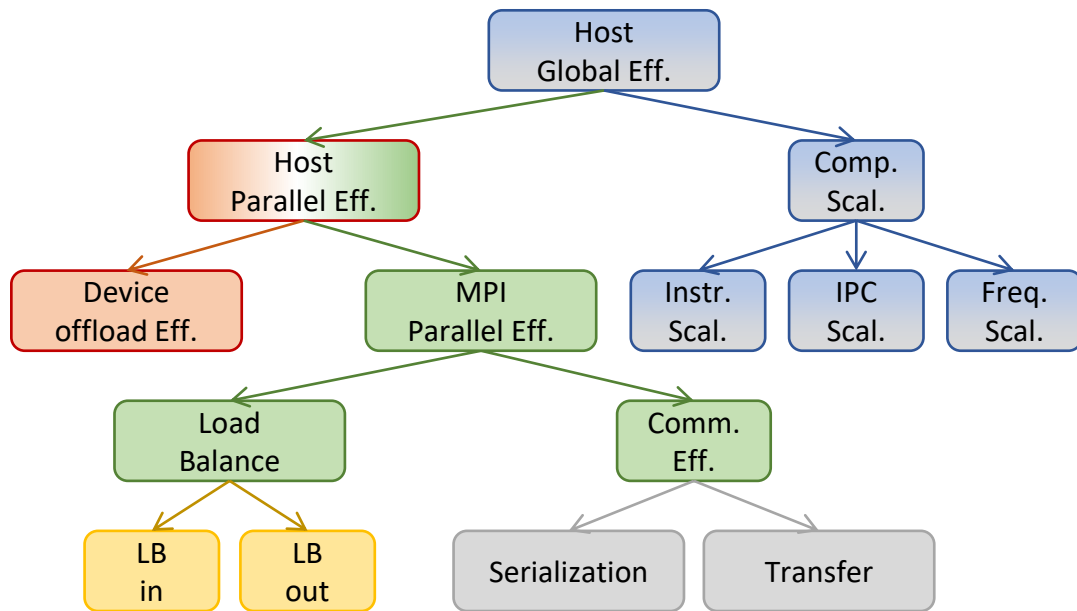
Additional metrics

Efficiency metrics



TALP GPU metrics

All parent metrics
are the product of
their child metrics



Device
Global Eff.

GPU metrics

Metrics not
computed by TALP

Additional metrics

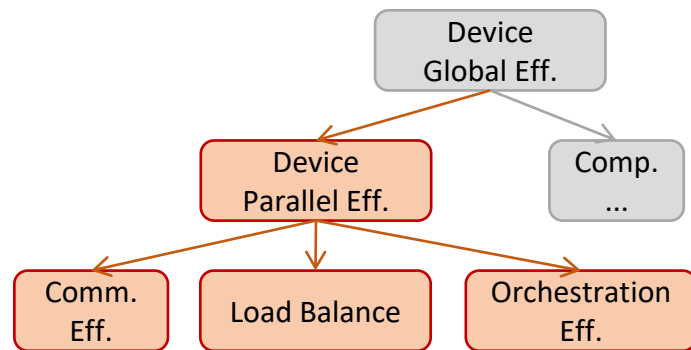
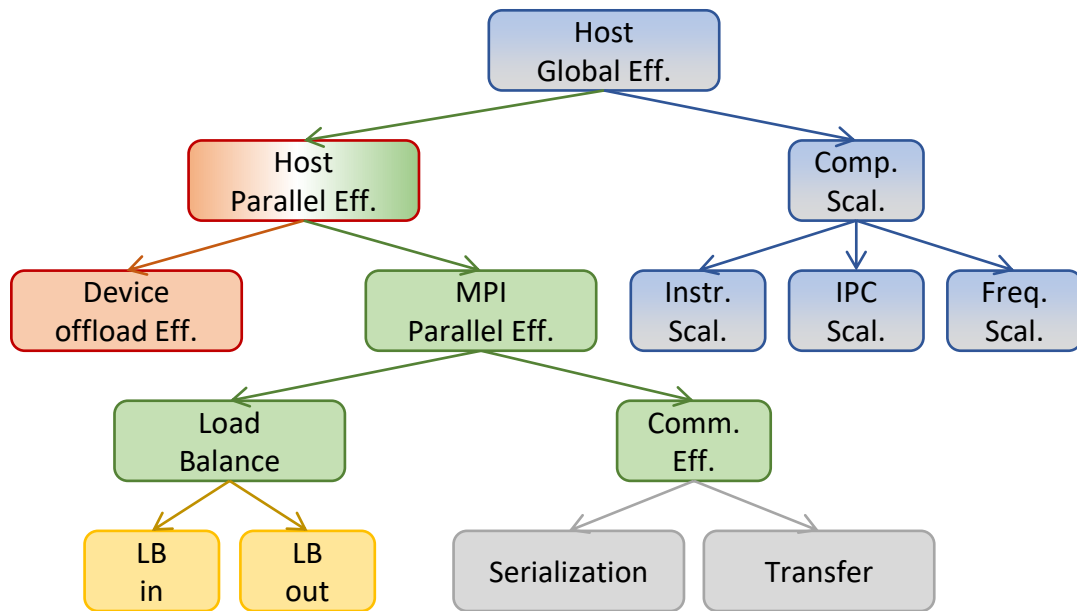
Scalability metrics
computed with
several TALP runs

Efficiency metrics



TALP GPU metrics

All parent metrics are the product of their child metrics



GPU metrics

Metrics not computed by TALP

Additional metrics

Scalability metrics computed with several TALP runs

Efficiency metrics



TALP GPU metrics

- Host and Device Efficiencies are “unrelated”
 - We can measure them separately
- Device Global Efficiency divided in
 - Device Parallel Efficiency
 - Computation (Sc./Eff.) → WiP
- We consider one GPU as a single resource
- All arrows are multiplicative

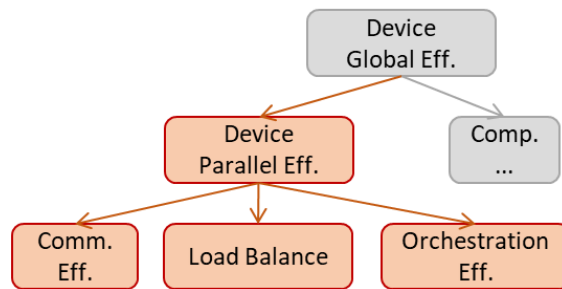
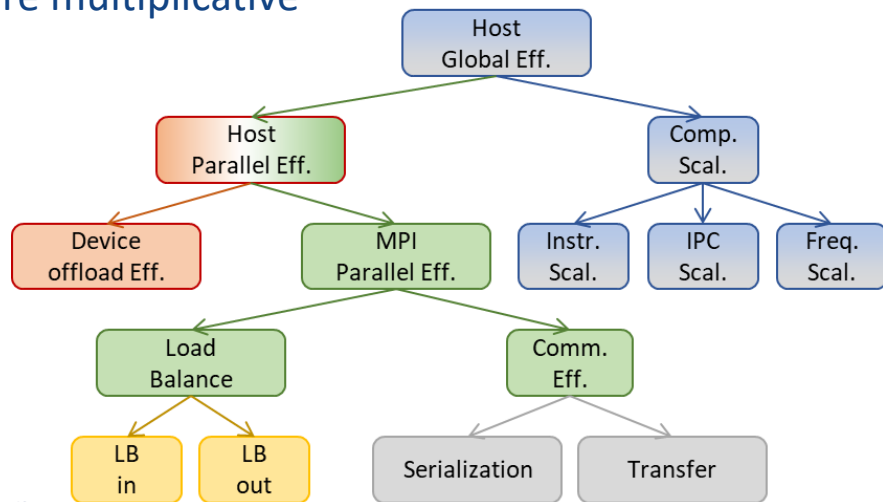
GPU metrics

Metrics not
computed by TALP

Additional metrics

Scalability metrics
computed with
several TALP runs

Efficiency metrics



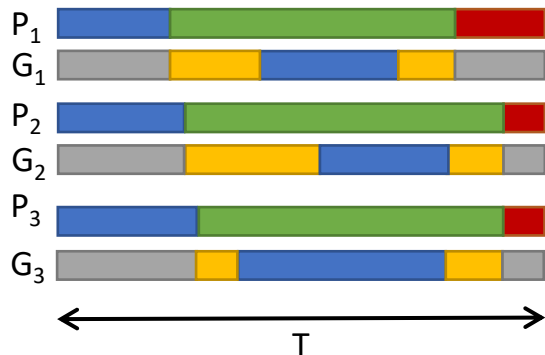
Computing TALP GPU metrics



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

States



3 states for the host (CPU):

Useful

Waiting 4 GPU

MPI

Whenever the CPU is not doing useful work because it is managing the GPU. Inside CUDA API: launching kernels, sending data to GPU, waiting for data...

P = num procs
 N = num GPUs
 T = Total elapsed time
 C_i = Useful time of CPU i
 W_i = Wait time of CPU i
 G_i = Useful time of GPU i

3 states for the device (GPU):

Useful

Data transfer

Idle

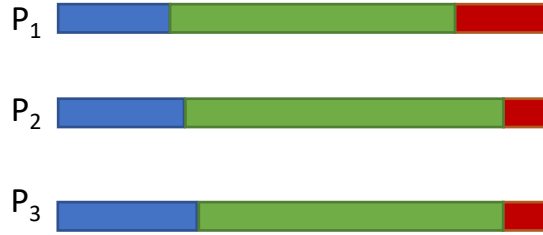
Whenever the GPU is waiting for data communication: GPU-GPU, or GPU-CPU. If communication is overlapped with computation, it is considered computation.

Host Metrics

Useful

Waiting 4 GPU

MPI



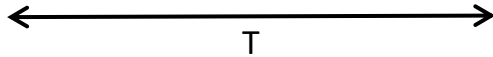
- To compute MPI metrics we consider **Waiting 4 GPU** as useful
- First, we blame MPI, then, Device offload

P = num procs

T = Total elapsed time

C_i = Useful time of CPU i

W_i = Wait time of CPU i



$$Par.Eff. = \frac{\sum_{i=1}^P c_i}{T * P} = \frac{\text{[Blue Bar]}}{T * P}$$

$$MPI Par.Eff. = \frac{\sum_{i=1}^P c_i + \sum_{i=1}^P w_i}{T * P} = \frac{\text{[Blue Bar]} + \text{[Green Bar]}}{T * P}$$

$$Device offload Eff. = \frac{\sum_{i=1}^P c_i}{\sum_{i=1}^P c_i + \sum_{i=1}^P w_i} = \frac{\text{[Blue Bar]}}{\text{[Blue Bar]} + \text{[Green Bar]}}$$



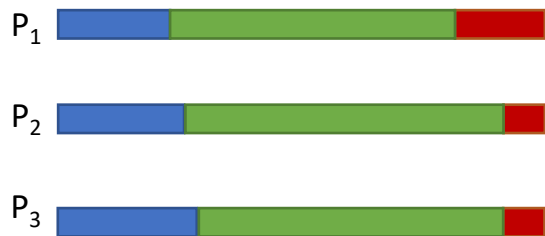
Quantifies how much time the resources are idle due the use of devices

Host Metrics (MPI branch)

Useful

Waiting 4 GPU

MPI



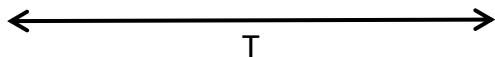
- Added for completeness
- Metrics under MPI branch are computed the same way as the MPI only metrics considering **Waiting 4 GPU** as useful

P = num procs

T = Total elapsed time

C_i = Useful time of CPU i

W_i = Wait time of CPU i



$$LB = \frac{\sum_{i=1}^P c_i + w_i}{\max(c_i + w_i) * P} = \frac{\text{blue} + \text{green}}{\max(\text{blue}_i + \text{green}_i) * P}$$

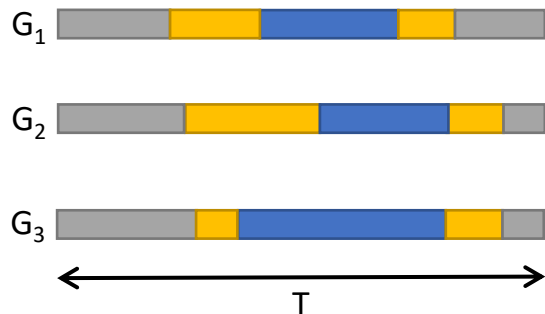
$$Comm = \frac{\max(c_i + w_i)}{T} = \frac{\max(\text{blue}_i + \text{green}_i) * P}{\text{blue} + \text{green} + \text{red}}$$

Device Parallel Efficiency

Useful

Data move

Idle



N = num GPUs

T = Total elapsed time

G_i = Useful time of GPU i

D_i = Time waiting for data in GPU i

$$Par.Eff. = \frac{\sum_{i=1}^N G_i}{T * N} = \frac{\text{Blue}}{T * N}$$

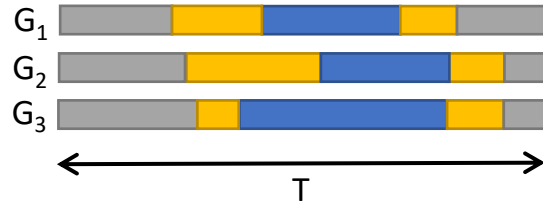
Quantifies how much time the devices are used to do useful work

Device Load Balance

Useful

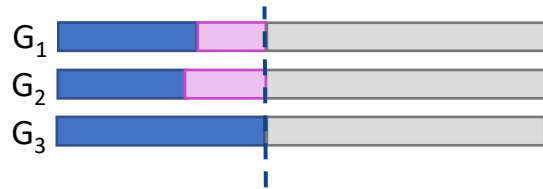
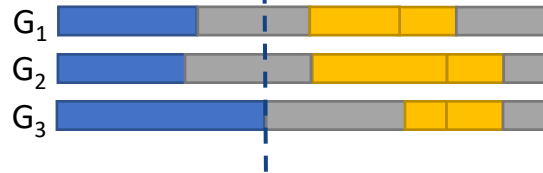
Data move

Idle



$$LB = \frac{\sum_{i=1}^N G_i}{\max(G_i) * N} = \frac{\text{Useful}}{\text{Useful} + \text{Data move}}$$

Aggregate useful time per device



N = num GPUs

T = Total elapsed time

G_i = Useful time of GPU i

D_i = Time waiting for data in GPU i

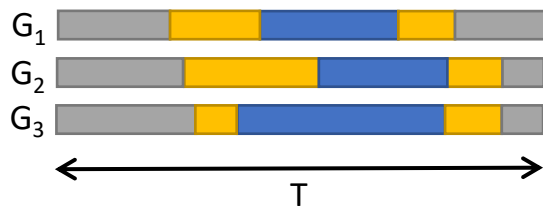
Quantifies how much time the devices are idle due to one device spending more time in useful work than others

Device Communication Efficiency

Useful

Data move

Idle



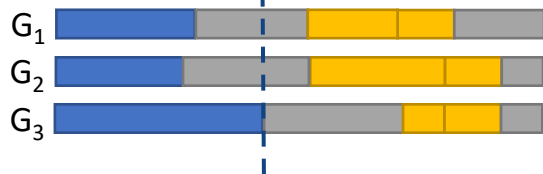
N = num GPUs

T = Total elapsed time

G_i = Useful time of GPU i

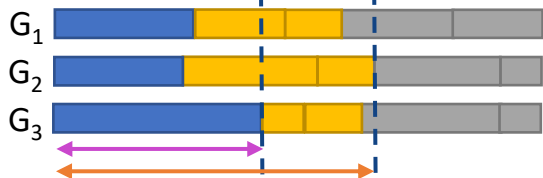
D_i = Time waiting for data in GPU i

Aggregate useful time per device



$$\text{Comm. Eff.} = \frac{\max(G_i)}{\max(G_j + D_j)} = \frac{\text{Useful time of fastest GPU}}{\text{Total time of slowest GPU}} = \frac{\text{Useful time}}{\text{Total time}}$$

Aggregate useful data waiting time per device



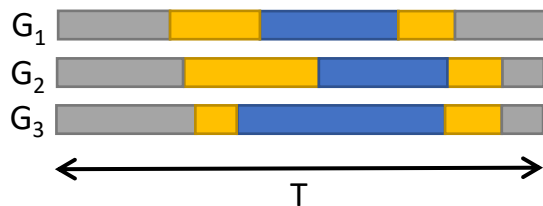
Quantifies how much time the devices are busy due to data movements with respect to useful computation

Device Orchestration Efficiency

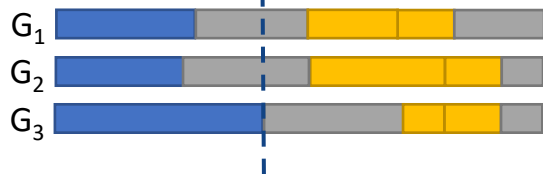
Useful

Data move

Idle



Aggregate useful time per device



$$Orch. Eff. = \frac{\max(G_i + D_i)}{T} = \frac{\text{orange arrow}}{T}$$

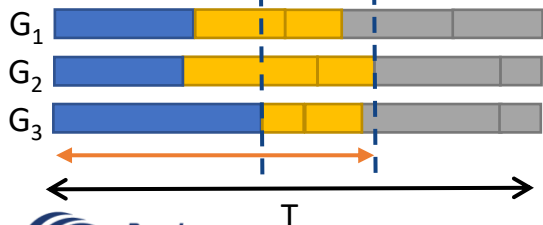
N = num GPUs

T = Total elapsed time

G_i = Useful time of GPU i

D_i = Time waiting for data in GPU i

Aggregate useful data waiting time per device



Quantifies how much time the devices are idle because there is no pending work to do

Summary of metrics explanation

➤ Host:

- **Device offload Eff.:** Inefficiency due to use of accelerator.
 - Includes: offloading work, waiting for kernels, waiting for data or sending data to accelerator.
 - Will be very low for applications that “only” use the device.

➤ Device:

- **Orchestration Eff.:** Inefficiency due to lack of available work to do.
 - Includes: waiting for work from host, dependencies between kernels...
 - Low value indicates the GPU is not efficiently used because lack of work to do
- **Communication Eff.:** Inefficiency due to data movement not instantaneous.
 - Includes: waiting for data from host, sending data to host, sending data to other accelerator, NCCL comm, MPI CUDA aware communication.
- **Load Balance:** Inefficiency due to not all the GPUs computing the same amount of time
 - Does not differentiate between GPUs from the same process or different processes
 - Maybe in the future we can add child metrics similar to LB_in and LB_out
- **Computation:** How well the resources inside the accelerator are being used.
 - TBD: streams, warps, occupancy, instructions, tensor core use....

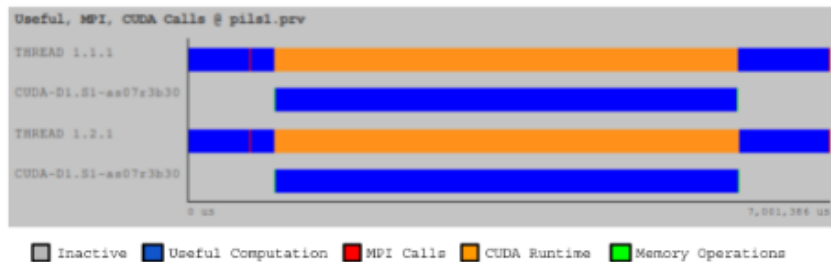
TALP GPU metrics examples



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

Most work offloaded to GPU. CPUs well balanced, GPUs well balanced

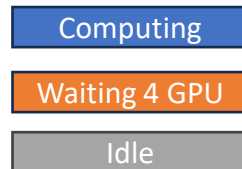
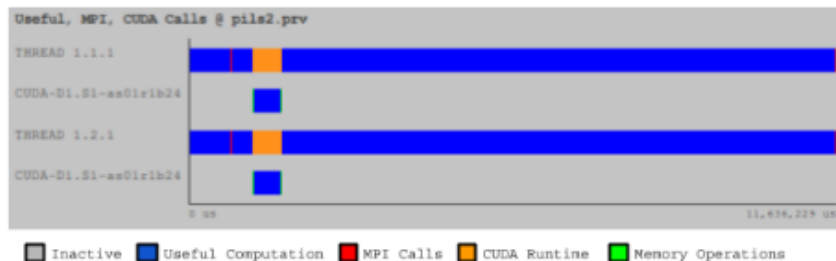


Metrics

Host	Parallel efficiency	0.16
	└ MPI Parallel eff.	1.00
	└ └ Comm. eff.	1.00
	└ └ Load Balance	1.00
	└ Device Offload eff.	0.16
Device	Parallel efficiency	0.82
	└ Load Balance	1.00
	└ Communication eff.	1.00
	└ Orchestration eff.	0.82

- Low Offload efficiency indicates CPUs are only used to offload work to GPUs
- Good device efficiency, GPUs are used efficiently

Few work offloaded to GPU. CPUs well balanced, GPUs well balanced

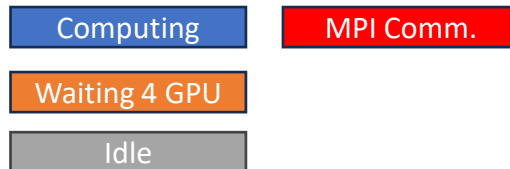


Metrics

Host	Parallel efficiency	0.94
	└ MPI Parallel eff.	1.00
	└ └ Comm. eff.	1.00
	└ └ Load Balance	1.00
	└ Device Offload eff.	0.94
Device	Parallel efficiency	0.05
	└ Load Balance	1.00
	└ Communication eff.	1.00
	└ Orchestration eff.	0.05

- Good host efficiency
 - CPUs are used efficiently
- Low orchestration efficiency indicates GPUs are not used efficiently because not enough work is offloaded to them

Load imbalance between GPUs. CPUs well balanced

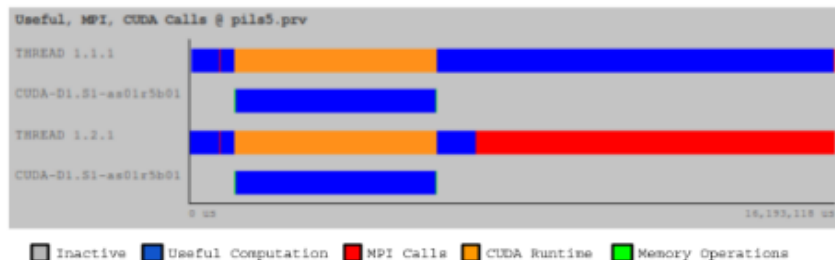


Metrics

Host	Parallel efficiency	0.16
	MPI Parallel eff.	0.63
	Comm. eff.	1.00
	Load Balance	0.63
	Device Offload eff.	0.26
Device	Parallel efficiency	0.45
	Load Balance	0.55
	Communication eff.	1.00
	Orchestration eff.	0.82

- Low MPI Load Balance
 - Indicates one process has more work to do than the other
- Low Offload efficiency
 - Indicates CPUs are not being used while waiting for GPUs
- Low device Load Balance
 - Indicates GPUs are not well balanced

GPUs well balanced CPUs unbalanced



Computing

MPI Comm.

Waiting 4 GPU

Idle

Metrics

Host	Parallel efficiency	0.36
	└ MPI Parallel eff.	0.70
	└ └ Comm. eff.	1.00
	└ └ Load Balance	0.70
	└ Device Offload eff.	0.51
Device	Parallel efficiency	0.33
	└ Load Balance	1.00
	└ Communication eff.	1.00
	└ Orchestration eff.	0.33

➤ Low MPI Load Balance

- Indicates there is load imbalance between processes

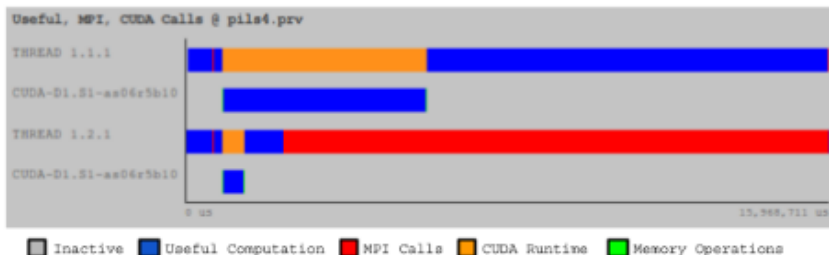
➤ Low offload efficiency

- Indicates CPUs are not working while waiting for GPUs

➤ Low orchestration efficiency

- Indicates not enough work is being offloaded to GPUs

GPUs unbalanced. CPUs unbalanced



Computing

MPI Comm.

Waiting 4 GPU

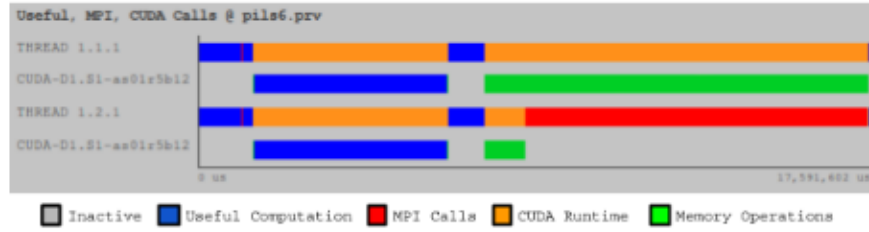
Idle

Metrics

Host	Parallel efficiency	0.36
	MPI Parallel eff.	0.55
	Comm. eff.	1.00
	Load Balance	0.55
	Device Offload eff.	0.65
Device	Parallel efficiency	0.18
	Load Balance	0.55
	Communication eff.	1.00
	Orchestration eff.	0.33

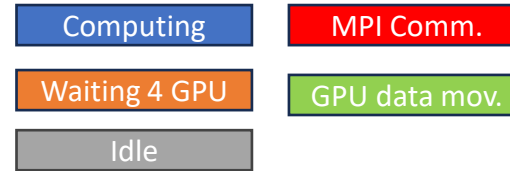
- Low MPI Load Balance
 - Indicates there is load imbalance between processes
- Low offload efficiency
 - Indicates CPUs are not working while waiting for GPUs
- Low device Load Balance
 - Indicates GPUs are not well balanced
- Low orchestration efficiency
 - Indicates not enough work is being offloaded to GPUs

CPU's well balanced. GPU's well balanced. Data movement in one of the processes



Metrics

Host	Parallel efficiency	0.06
	└ MPI Parallel eff.	0.72
	└ └ Comm. eff.	1.00
	└ └ Load Balance	0.72
	└ Device Offload eff.	0.09
Device	Parallel efficiency	0.31
	└ Load Balance	1.00
	└ Communication eff.	0.36
	└ Orchestration eff.	0.86



- Low MPI load balance
 - Indicates imbalance between processes
- Low offload efficiency
 - Indicates CPUs are not being used while waiting for GPUs
- Low device communication efficiency
 - Indicates data movement is limiting the use of the GPUs

TALP Use cases



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

TALP use cases

➤ Use cases examples:

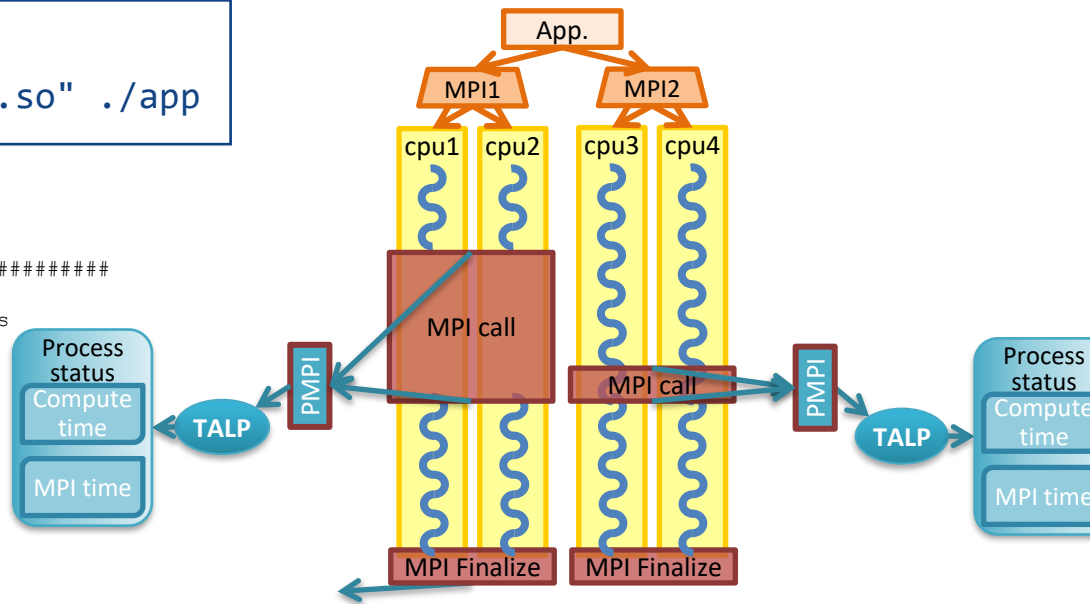
- Transparent use for the user
- With user defined regions
- Getting metrics at runtime
- TALP pages: Continuous Performance Monitoring

TALP: Transparent use for the user

```
DLB_ARGS=" --talp"
```

```
env LD_PRELOAD="$DLB_LIBS/libdlb_mpi.so" ./app
```

```
DLB[...]: ##### Monitoring Region POP Metrics #####
DLB[...]: ### Name: Global
DLB[...]: ### Elapsed Time: 31.76 s
DLB[...]: ### Parallel efficiency: 0.70
DLB[...]: ### - MPI Parallel efficiency: 0.70
DLB[...]: ### - Communication efficiency: 1.00
DLB[...]: ### - Load Balance: 0.70
DLB[...]: ### - In: 0.70
DLB[...]: ### - Out: 1.00
DLB[...]: ### - OpenMP Parallel efficiency: 0.84
DLB[...]: ### - Load Balance: 1.00
DLB[...]: ### - Scheduling efficiency: 1.00
DLB[...]: ### - Serialization efficiency: 0.84
DLB[...]: ### Computational metrics:
DLB[...]: ### - Average useful IPC: 0.59
DLB[...]: ### - Average useful frequency: 2.95 GHz
DLB[...]: ### - Number of instructions: 1.55E+11
```



No knowledge from the user needed
Metrics reported transparently at finalization

TALP: With user defined regions

```
# include < dlb_talp.h >
...
// Register a new region or obtain an existing handler
dlb_monitor_t * monitor = DLB_MonitoringRegionRegister ("Name");
// Start TALP monitoring region
DLB_MonitoringRegionStart( monitor );
...
// Stop TALP monitoring region
DLB_MonitoringRegionStop( monitor );
...
```

Code modification needed
Metrics reported by region
Any kind of nesting is allowed

regions

Main region

me”);

Region without MPI code

Region with MPI code

```
DLB[...]: ##### Monitoring Region POP Metrics #####
DLB[...]: ### Name: Global
DLB[...]: ### Elapsed Time: 25 s
DLB[...]: ### Parallel efficiency: 0.70
DLB[...]: ### - MPI Parallel efficiency: 0.70
DLB[...]: ### - Communication efficiency: 1.00
DLB[...]: ### - Load Balance: 0.70
DLB[...]: ### - In: 0.70
DLB[...]: ### - Out: 1.00
DLB[...]: ### Computational metrics:
DLB[...]: ### - Average useful IPC: 1.15
DLB[...]: ### - Average useful frequency: 2.99 GHz
DLB[...]: ### - Number of instructions: 1.20E+11
DLB[...]: ##### Monitoring Region POP Metrics #####
DLB[...]: ### Name: Kernel computation
DLB[...]: ### Elapsed Time: 25 s
DLB[...]: ### Parallel efficiency: 1.00
DLB[...]: ### Computational metrics:
DLB[...]: ### - Average useful IPC: 1.15
DLB[...]: ### - Average useful frequency: 2.99 GHz
DLB[...]: ### - Number of instructions: 1.20E+11
DLB[...]: ##### Monitoring Region POP Metrics #####
DLB[...]: ### Name: Main loop
DLB[...]: ### Elapsed Time: 25 s
DLB[...]: ### Parallel efficiency: 0.70
DLB[...]: ### - MPI Parallel efficiency: 0.70
DLB[...]: ### - Communication efficiency: 1.00
DLB[...]: ### - Load Balance: 0.70
DLB[...]: ### - In: 0.70
DLB[...]: ### - Out: 1.00
DLB[...]: ### Computational metrics:
DLB[...]: ### - Average useful IPC: 1.15
DLB[...]: ### - Average useful frequency: 2.99 GHz
DLB[...]: ### - Number of instructions: 1.20E+11
```

incl

...

// Reg

dlb_mc

// Sta

DLB_Mc

...

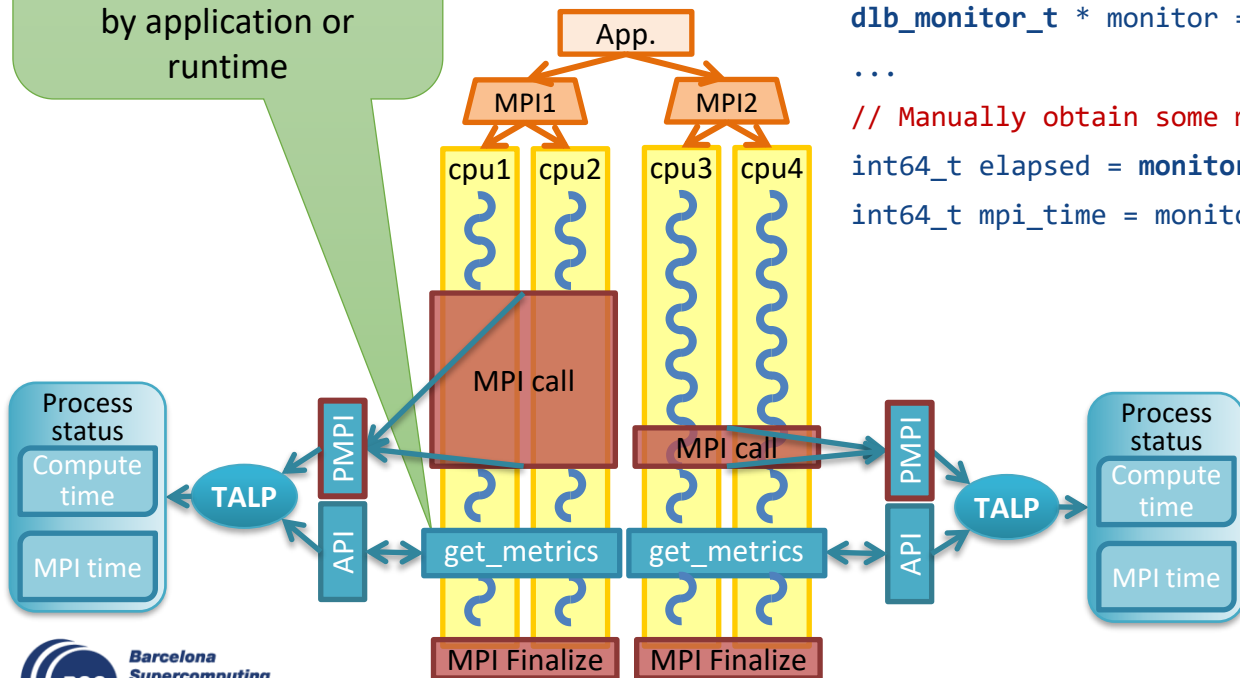
// Stc

DLB_Mc

...

TALP: getting metrics at runtime

At this point, informed decisions to improve efficiency can be taken by application or runtime



```
# include < dlb_talp.h >
```

```
...
```

```
// Register a new region or obtain an existing handler
```

```
dlb_monitor_t * monitor = DLB_MonitoringRegionRegister ("Name");
```

```
...
```

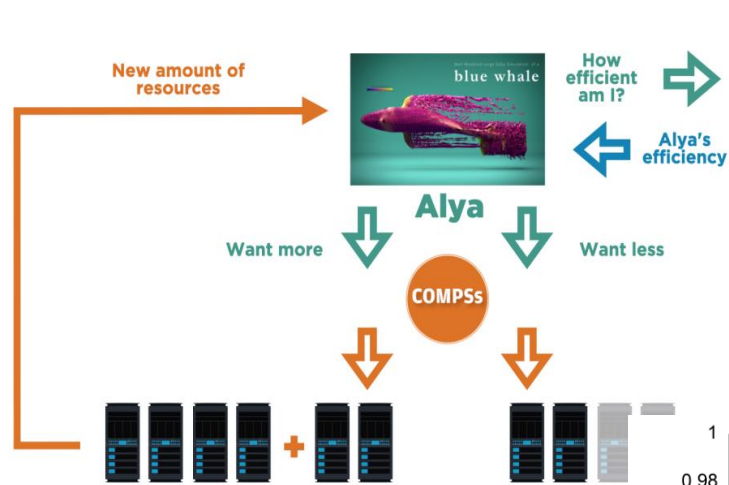
```
// Manually obtain some metrics from the monitor
```

```
int64_t elapsed = monitor->elapsed_time;
```

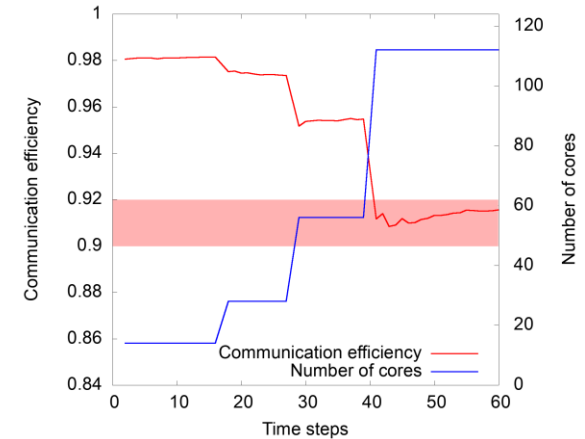
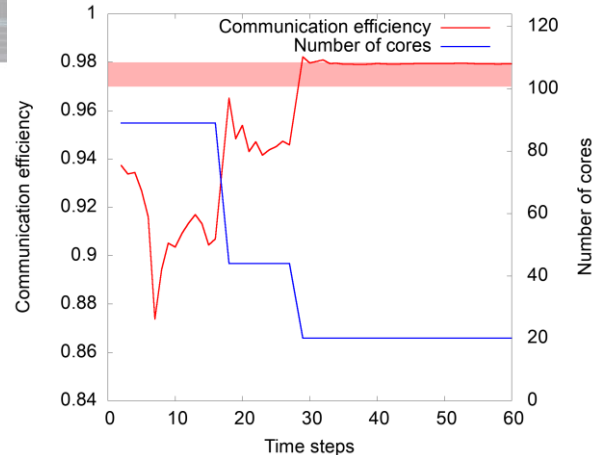
```
int64_t mpi_time = monitor->mpi_time
```

Enable dynamic resource management, load balancing or malleability

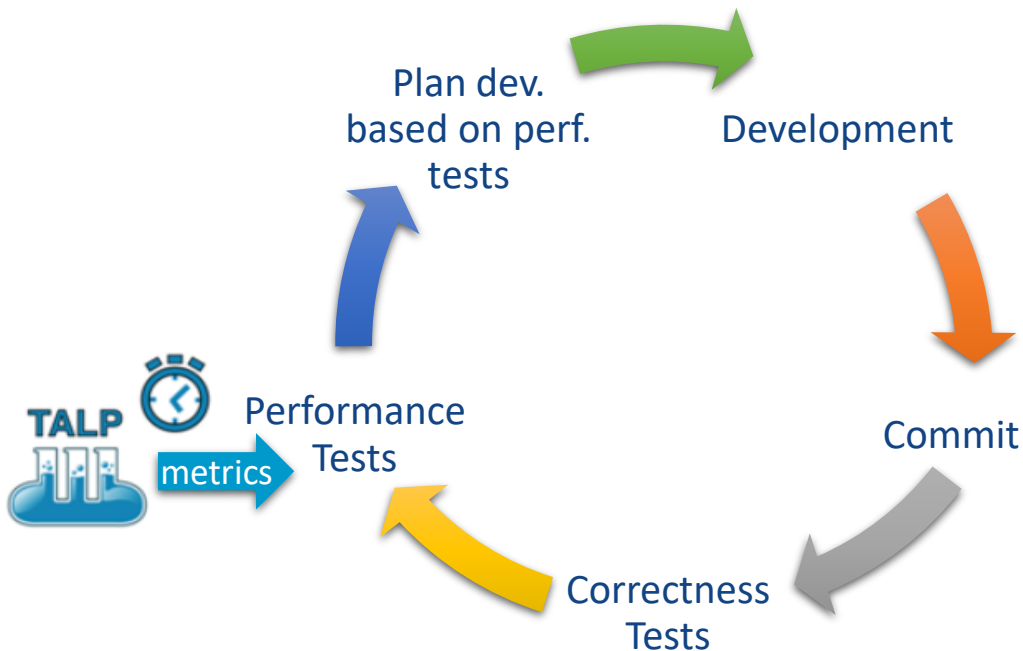
Success story 3: Malleable simulation



- Application obtain efficiency metrics through TALP
- Adjust resources accordingly using COMPSs
- Target efficiency reached after some iterations



TALP pages: Continuous Performance Monitoring



*Integrate in CI/CD platform to
detect performance issues added*

➤ Two visualization modes available:

- **Scaling efficiency** tables with POP-like metrics (to gain **insight**)
- **Metrics evolution** plots of POP-like metrics (to track **regression**)

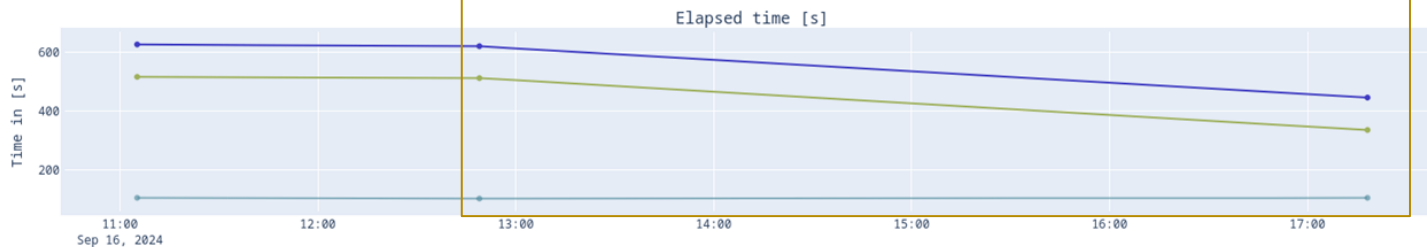
TALP pages: Scaling efficiency

Region: timestep

Metrics	4xMPI 56xOpenMP	8xMPI 56xOpenMP	4xMPI 112xOpenMP
Global Efficiency	0.88	0.86	0.61
- Parallel efficiency	0.88	0.85	0.75
-- MPI Parallel efficiency	0.98	0.97	0.98
--- MPI Communication efficiency	1	1	1
--- MPI Load balance	0.99	0.97	0.98
---- MPI In-node load balance	0.99	0.98	1
---- MPI Inter-node load balance	1	0.99	0.98
-- OpenMP Parallel efficiency	0.88	0.85	0.75
--- OpenMP Scheduling efficiency	1	1	0.99
--- OpenMP Load balance	0.99	0.98	0.97
--- OpenMP Serialization efficiency	0.9	0.87	0.78
- Computation Scalability	1	1.01	0.81
-- Instructions scaling	1	0.99	0.93
-- IPC scaling	1	0.95	0.82
-- Frequency scaling	1	1.07	1.06
Useful IPC	2.65	2.52	2.18
Frequency [GHz]	2.63	2.83	2.8
Elapsed time [s]	106.35	54.57	77.35

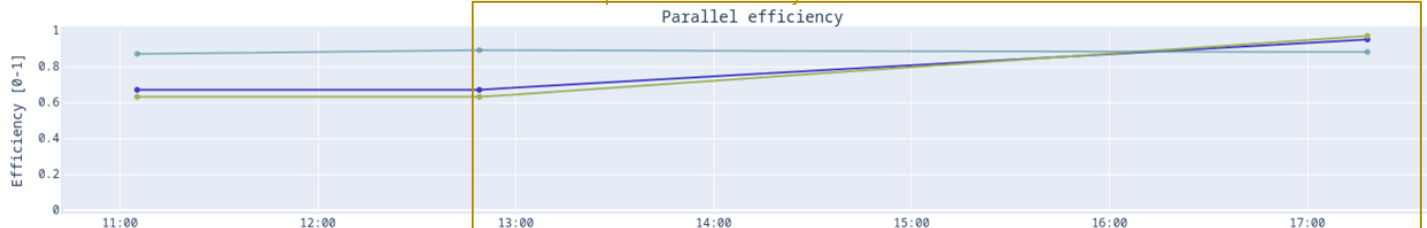
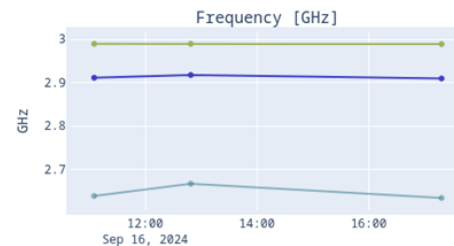
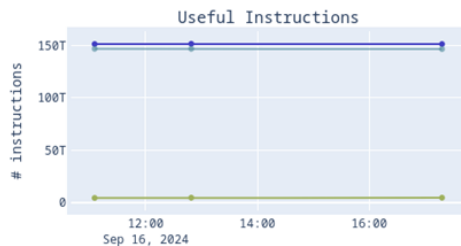
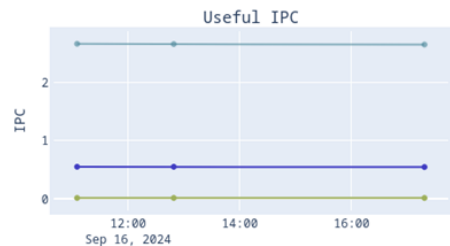
TALP Pages: Metrics evolution

Performance metrics evolution



selectable regions

- Regions
- Application
 - genex
 - initialize
 - initialize_timestep
 - initialize_mesh
 - timestep
 - calc_rhs_static
 - op_rhs_vlasov_eq_static
 - solve_maxwells_equation
 - solve_ohms_low
 - op_solve_ohms_low
 - op_solve_gh_eq
 - calc_rhs_dynamic
 - op_rhs_vlasov_eq_dynamic
 - op_mom_ohms_low
 - exchange
 - op_solve_amps_low
 - mpi_allreduce
 - op_mom_maxwells_eq
 - add
 - lin_comb
 - file_io
 - receive_2d
 - pack_5d
 - unpack_5d



Summary



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

Summary

- TALP is a lightweight tool to gather efficiency metrics
 - At finalization
 - Allows continuous performance monitoring
 - Integration with CI/CD systems (TALP-pages)
 - At runtime
 - Allows dynamic adjustment of execution
 - Enable load balancing
 - Dynamic resource management
 - Its low overhead allows its use in production runs
 - Does not store TB of data
 - Provides API to annotate regions and maximize the information gathered
 - Will indicate when detailed analysis using traces is needed

Summary

➤ Version 3.6.0-beta1 (2025-9)

- MPI metrics – Fully supported
- Hardware counters (Instructions, cycles, IPC) – Fully supported
- GPU metrics
 - NVIDIA (CUDA and OpenACC) – available
 - AMD (HIP) – Under development
 - Computational metrics – Under development
- OpenMP metrics – Under testing

➤ Download DLB (Free Download under LGPL-v3 license):

- <https://pm.bsc.es/dlb-downloads>
- <https://github.com/bsc-pm/dlb/releases/tag/v3.6.0-beta1>
- User Guide: <https://pm.bsc.es/ftp/dlb/doc/user-guide/>
- Hands-on: <https://gitlab.pm.bsc.es/dlb/dlb-training>



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación



EXCELENCIA
SEVERO
OCHOA

Thank you

marta.garcia@bsc.es

valentin.seitz@bsc.es

<https://pm.bsc.es/dlb>