ADMIRE Users Day

# Metric Proxy: Enabling real-time measurement at Supercomputer Scale

**Jean-Baptiste Besnard, ParaTools SAS**

**December 12th 2023.**
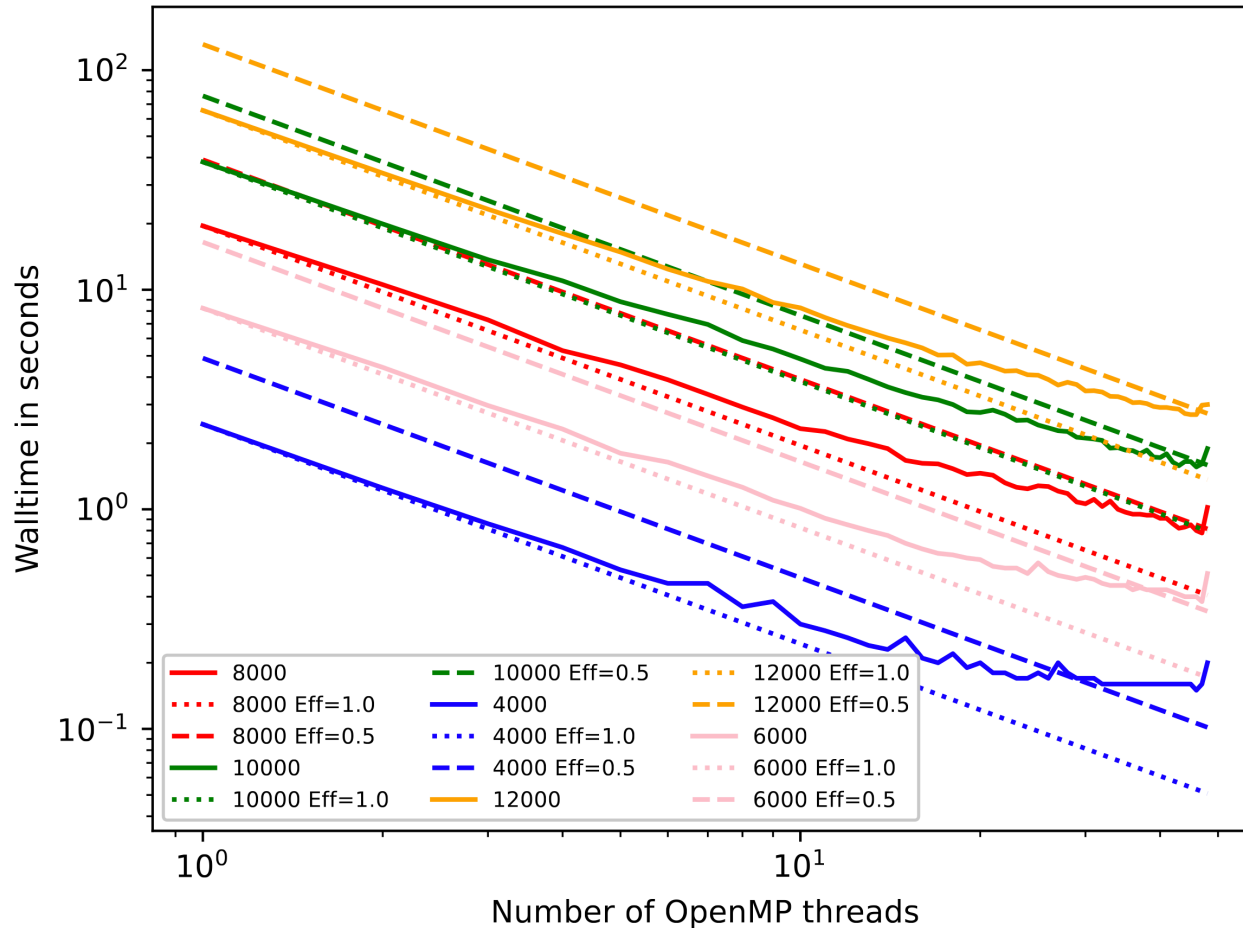
**Barcelona Supercomputing Center**

# Performance Monitoring For Malleability

- **Malleability** is about adapting the payload to external constraints to maximise machine throughput:
  - Optimize computation
  - Minimize wait-time
  - Maximize machine utilization
  - Lower Power
- It is a multi-criterion process, and therefore it requires a wide-range of monitoring capabilities to feed the various models.
- This motivated a general approach for monitoring in ADMIRE with two main challenges:
  - Need for **real-time** data (malleability is temporal)
  - Need for **machine-wide** metrics
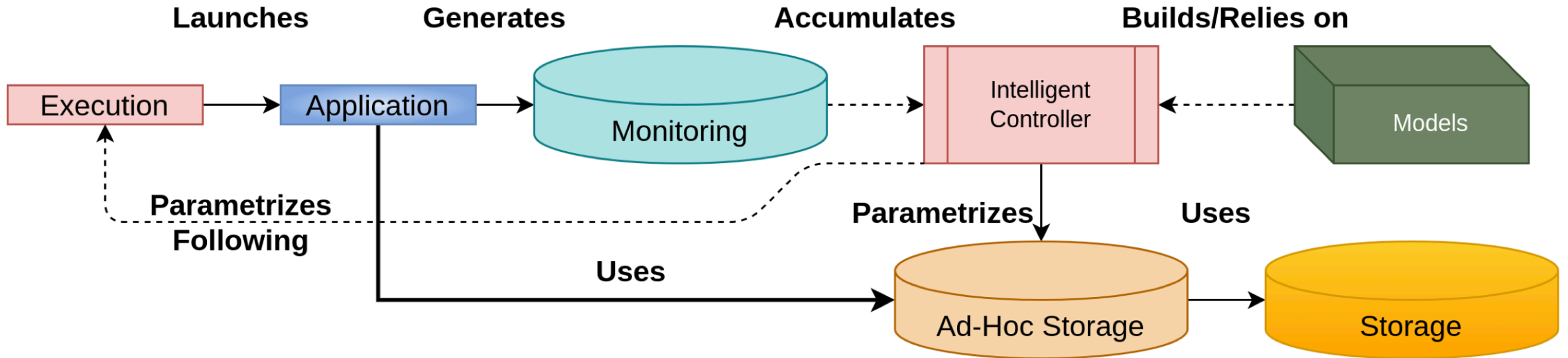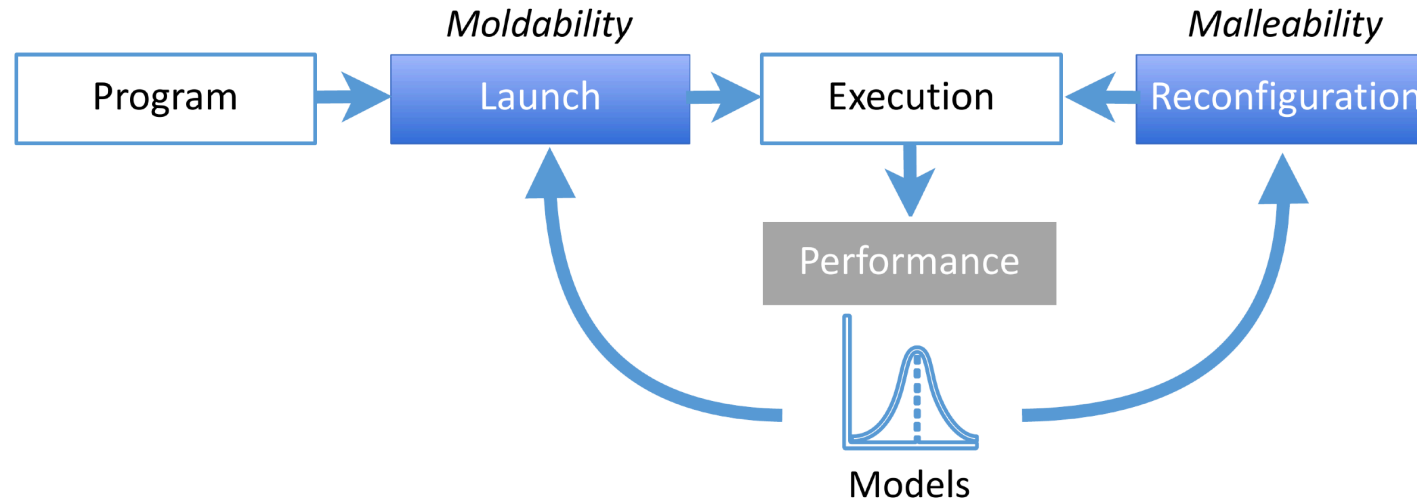  - Need for **per-program** models

} **Challenging !**

# Example: Moldability

**Efficiency of Rodinia LU Benchmark at various scales and problem sizes**
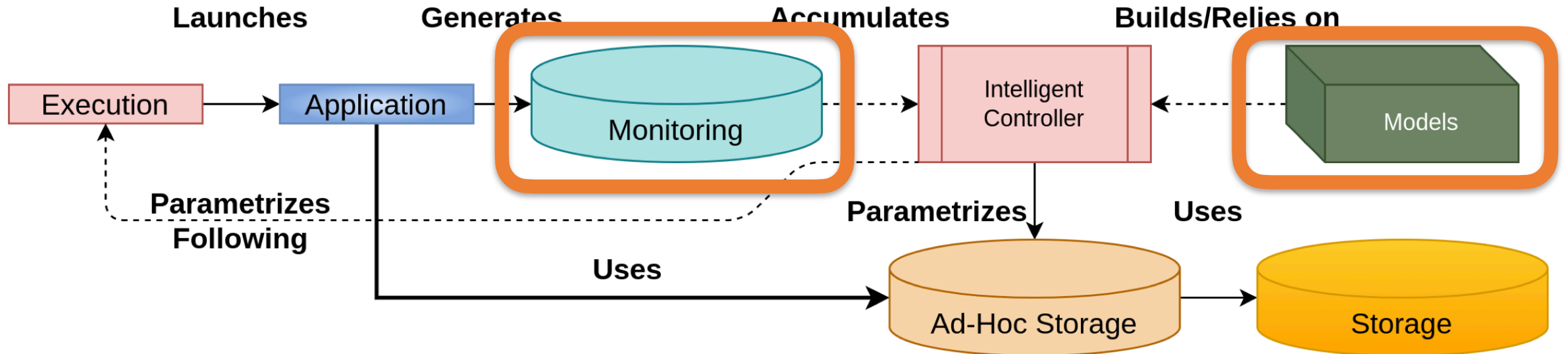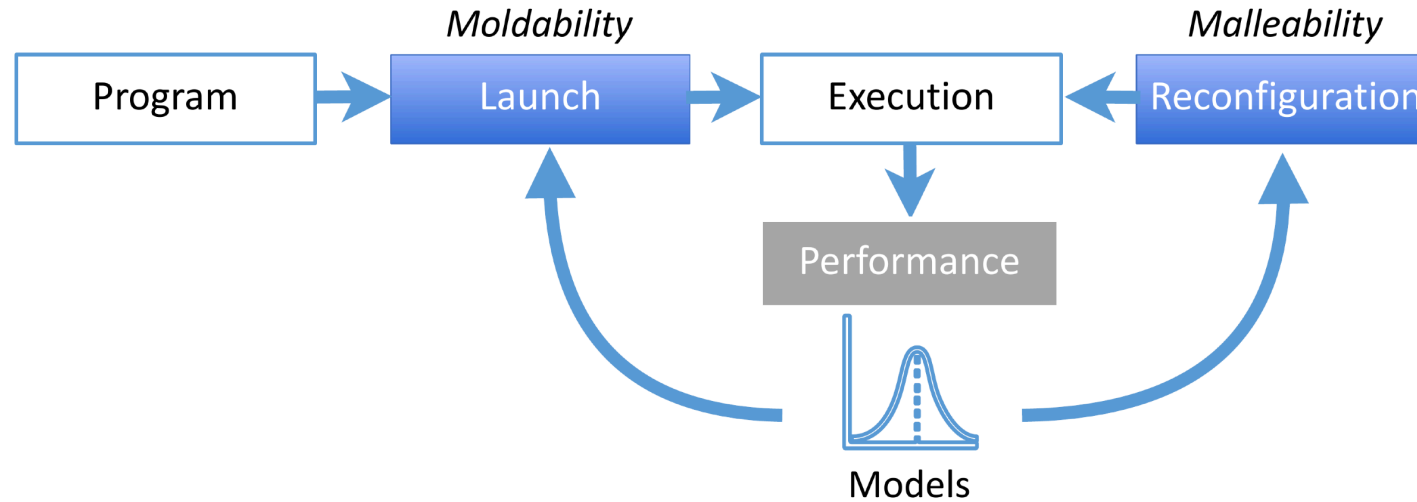


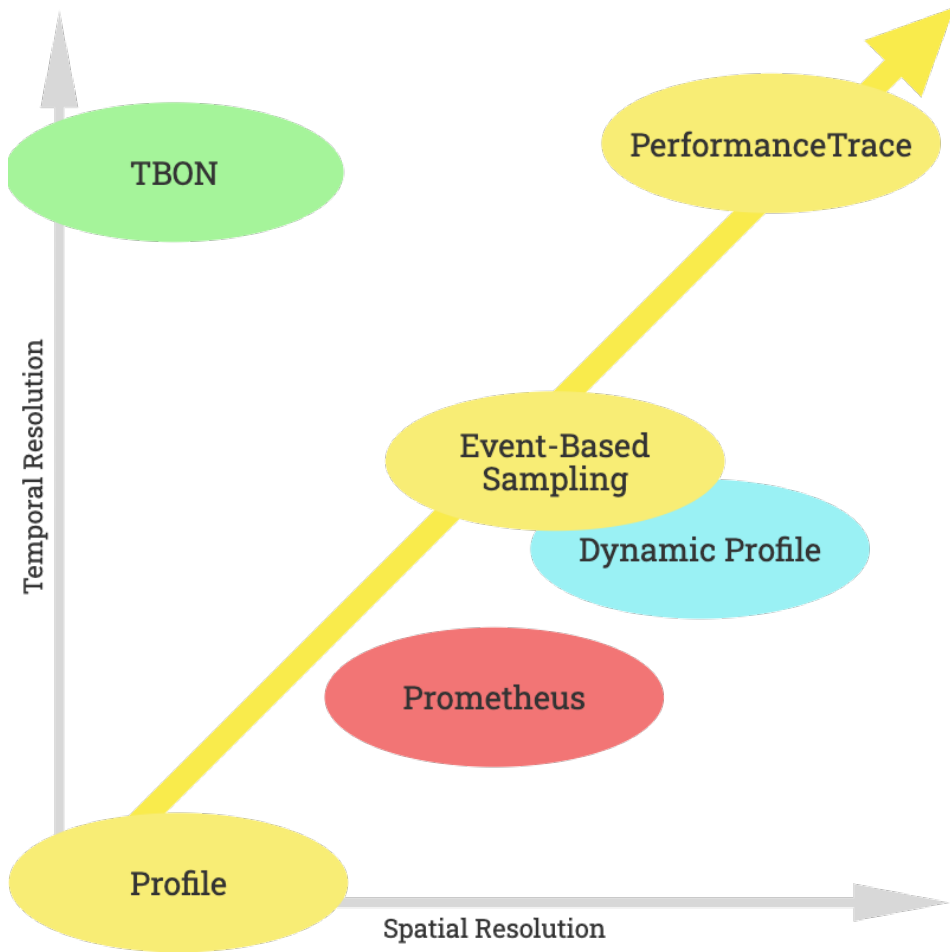It is a special case of *malleability*, called *moldability*, or more straightforwardly:

« choosing the right configuration at program start ».
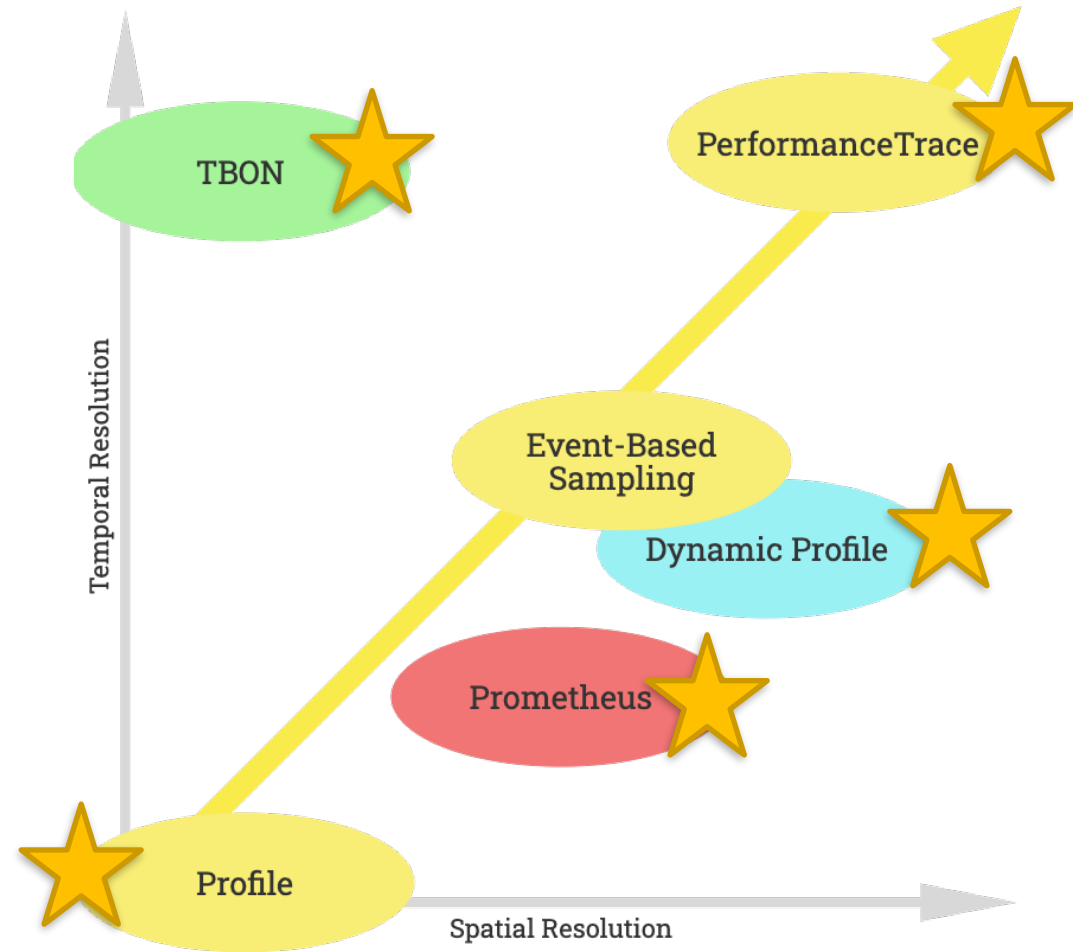
# Recall of the ADMIRE Feedback Loop
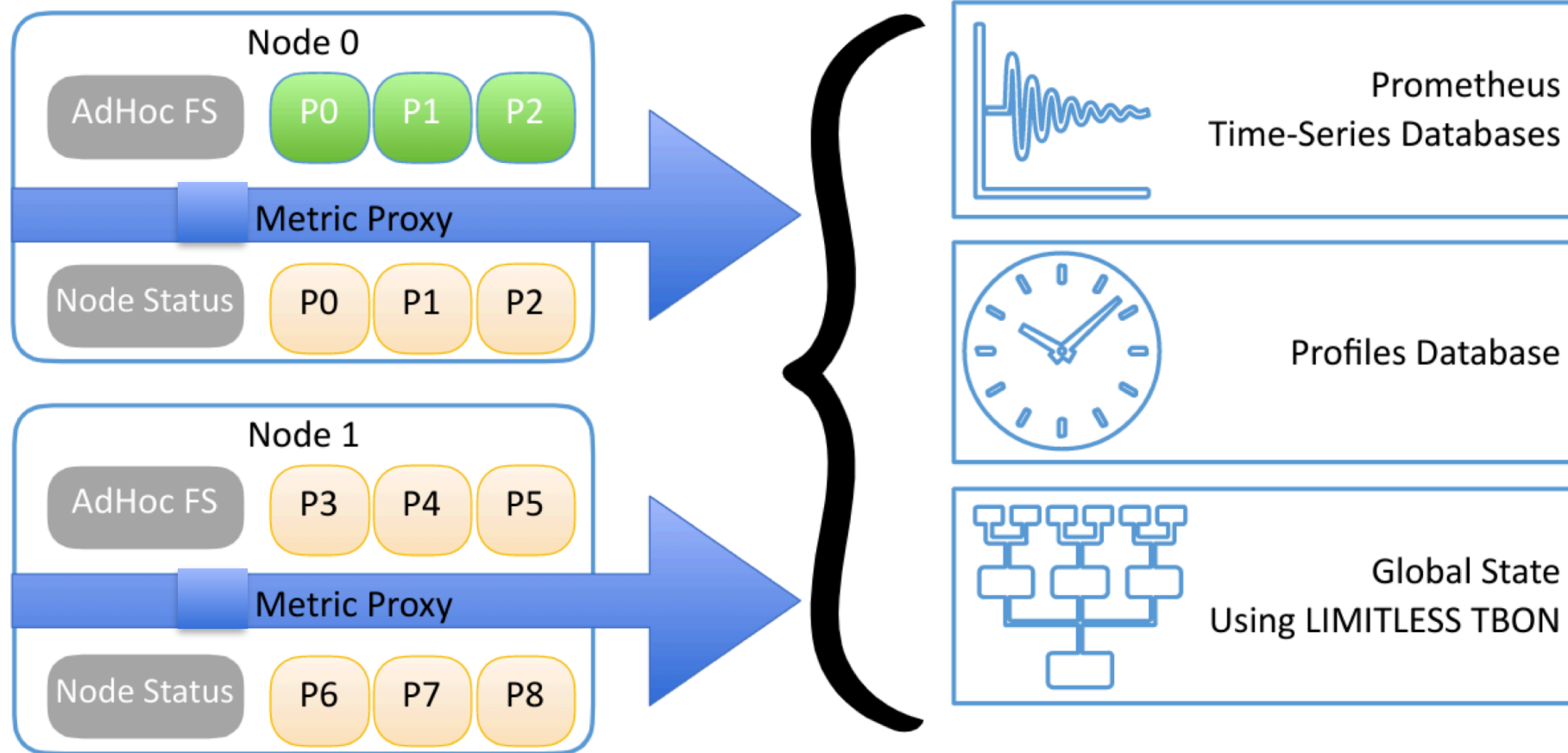
# Recall of the ADMIRE Feedback Loop

Performance measurement is always a compromise between Verbosity and measurement / storage overhead.
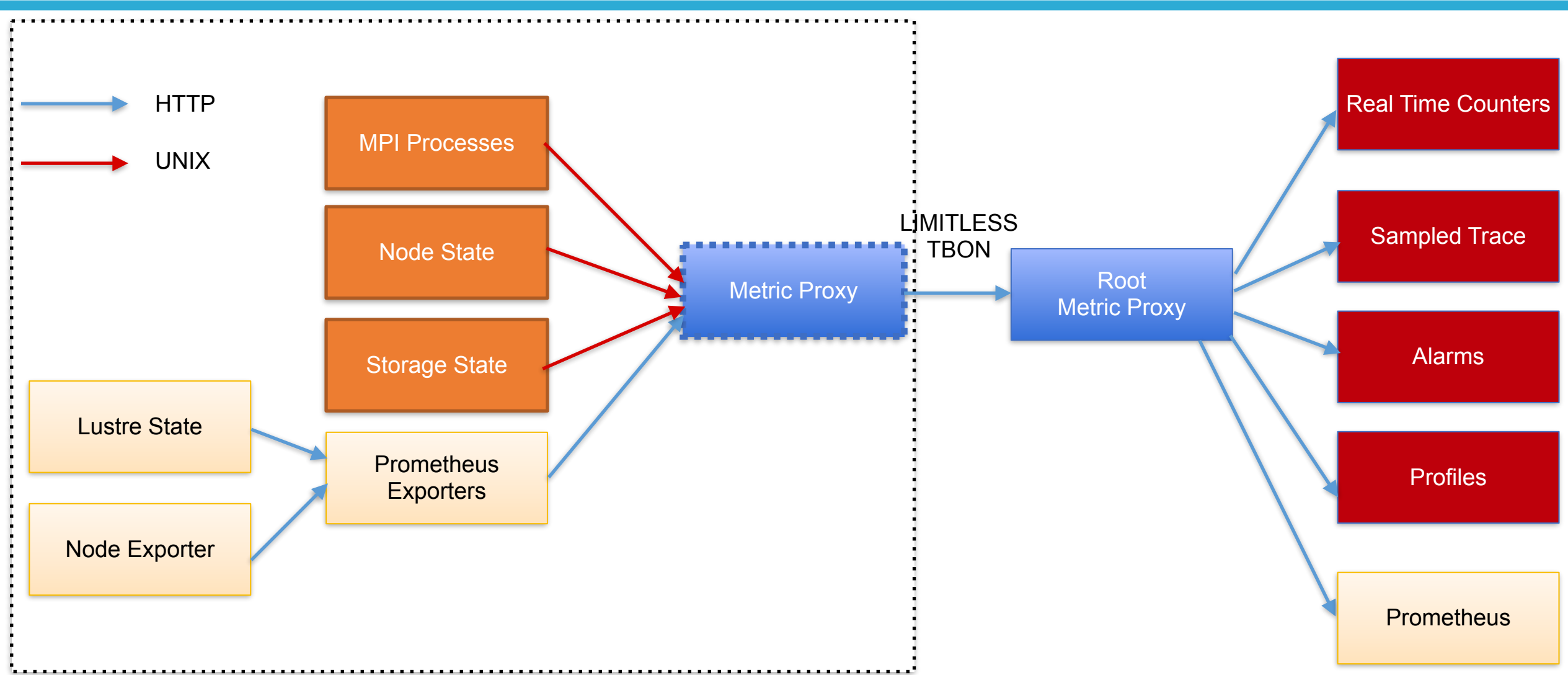
# Choosing the Right Measurement Granularity



- TBON: for real-time reduction of performance data using LIMITLESS
- Resampled performance traces: for temporal series
- Profiles to describe each run
- Prometheus storage for historization
- Real-time summative profiles (a.k.a snapshots) for current state
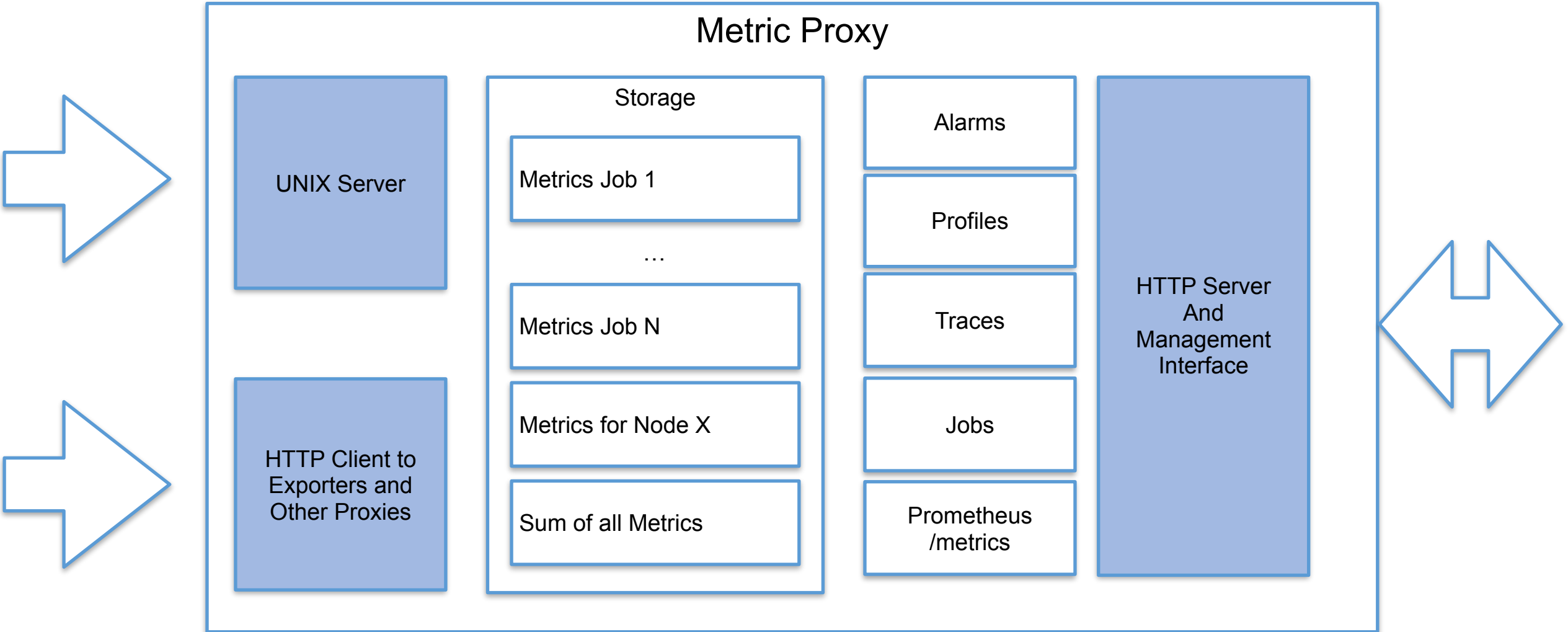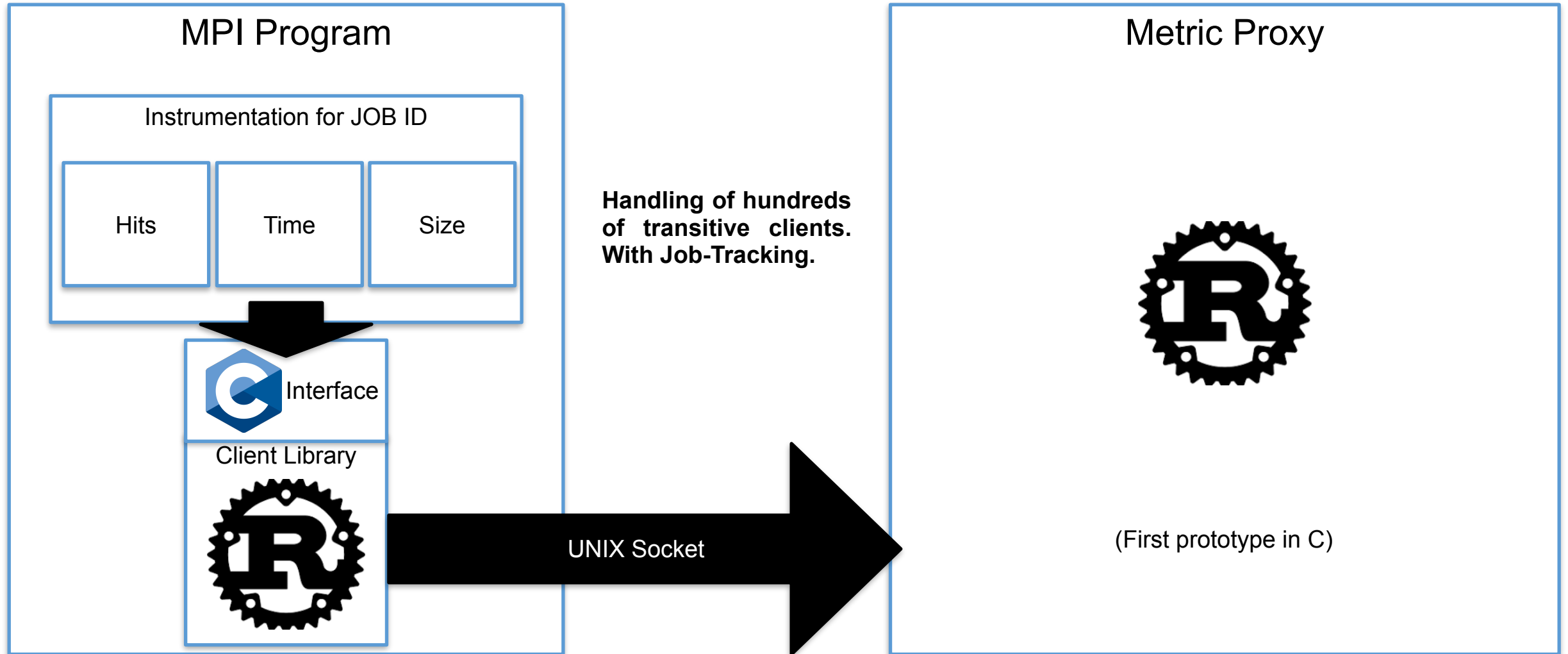
# Metric Proxy Architecture

# Metric Proxy Architecture

# Metric Proxy Architecture

**Metric Proxy**

UNIX Server

HTTP Client to Exporters and Other Proxies

**Storage**

Metrics Job 1

…

Metrics Job N

Metrics for Node X

Sum of all Metrics

Alarms

Profiles

Traces

Jobs

Prometheus /metrics

HTTP Server And Management Interface

# Handling of Transitive Clients and Jobs

# Metric Proxy : Interfaces

**Each metric proxy on each node provides an HTTP endpoint on port 1337 by default.**

# Proxy Topology



The proxy reduction tree is built automatically by « pivoting » the nodes on a root server which then returns the address of one of the proxy. Here an example with 32 nodes, seen from the root.

**} Scrapes**

# Real-Time Job Tracking

# Real-Time Monitoring at Scale

**METRIC PROXY**

Home | Jobs | Alarms | Proxy Topology | Trace | Profiles | API Documentation

## JOB DETAILS

### Job Description

| Key | Value |
| --- | --- |
| jobid | 535363585 |
| command | ./IMB-MPI1 |
| size | 3 |
| nodelist | |
| partition | |
| cluster | |
| run_dir | /tmp/IMB/src |
| start_time | 1702039102 |
| end_time | 0 |

### Counters

| Name | Documentation | Value | | |
| --- | --- | --- | --- | --- |
| mpi___hits___mpi_cartdim_get | Number of function calls for MPI_Cartdim_get | 0 | | |
| mpi___hits___mpi_publish_name | Number of function calls for MPI_Publish_name | 0 | | |
| mpi___time___mpi_type_ub | Total seconds spent for MPI_Type_ub | 0 | | |
| mpi___hits___mpi_send_init | Number of function calls for MPI_Send_init | 0 | | |
| mpi___hits___mpi_status_set_elements_x | Number of function calls for MPI_Status_set_elements_x | 0 | | |
| mpi___time___mpi_comm_set_attr | Total seconds spent for MPI_Comm_set_attr | 0 | | |
| mpi___time___mpi_cart_create | Total seconds spent for MPI_Cart_create | 0 | | |
| mpi___hits___mpi_raccumulate | Number of function calls for MPI_Raccumulate | 0 | | |
| mpi___hits___mpi_type_create_subarray | Number of function calls for MPI_Type_create_subarray | 0 | | |
| mpi___time___mpi_win_delete_attr | Total seconds spent for MPI_Win_delete_attr | 0 | | |
| proxy_memory_swap_used_percent | Total swap usage on the system in percent | AVG: 23.154901660001453 | Min: 23.154901660001453 | Max: 23.154901660001453 |
| mpi___hits___mpi_win_fence | Number of function calls for MPI_Win_fence | 0 | | |
| mpi___hits___mpi_allreduce | Number of function calls for MPI_Allreduce | 147047 | | |

EuroHPC
Joint Undertaking

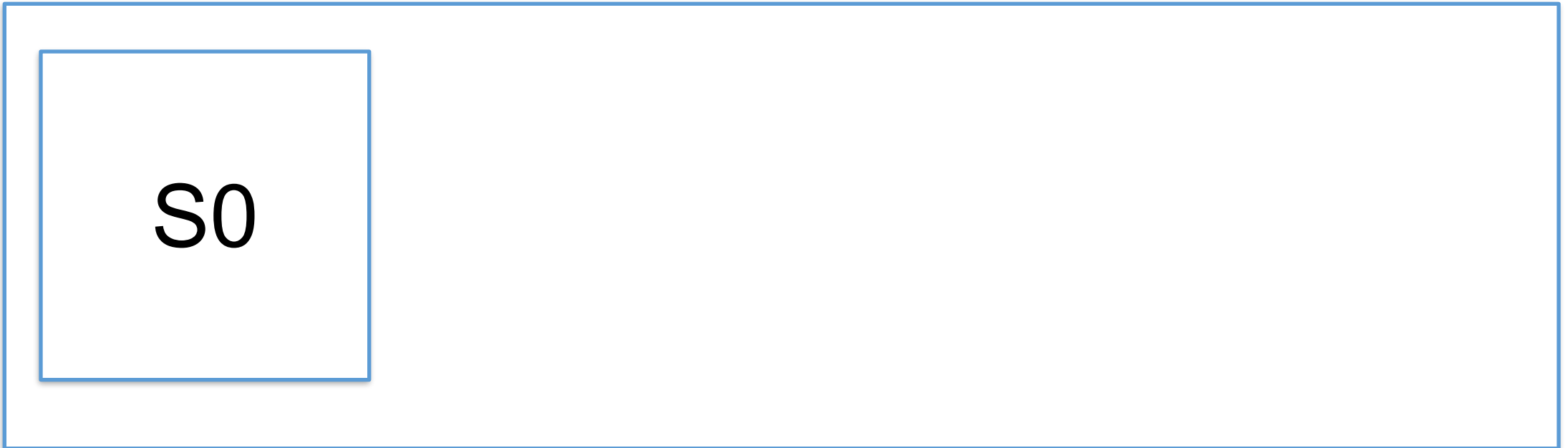# Real-Time Job Tracing



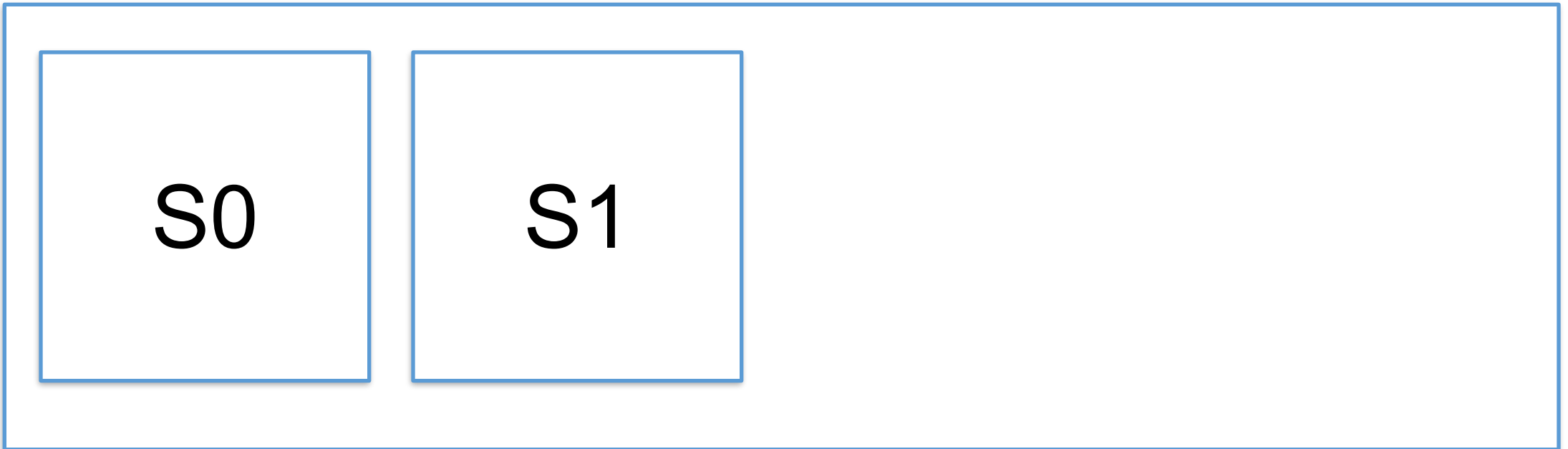Fixed size of 32MB maximum per job, filled with sample every one second and slowing down by a factor 2 on resampling.

This maintains full-range traces with dynamically decreasing resolution and bounded size.
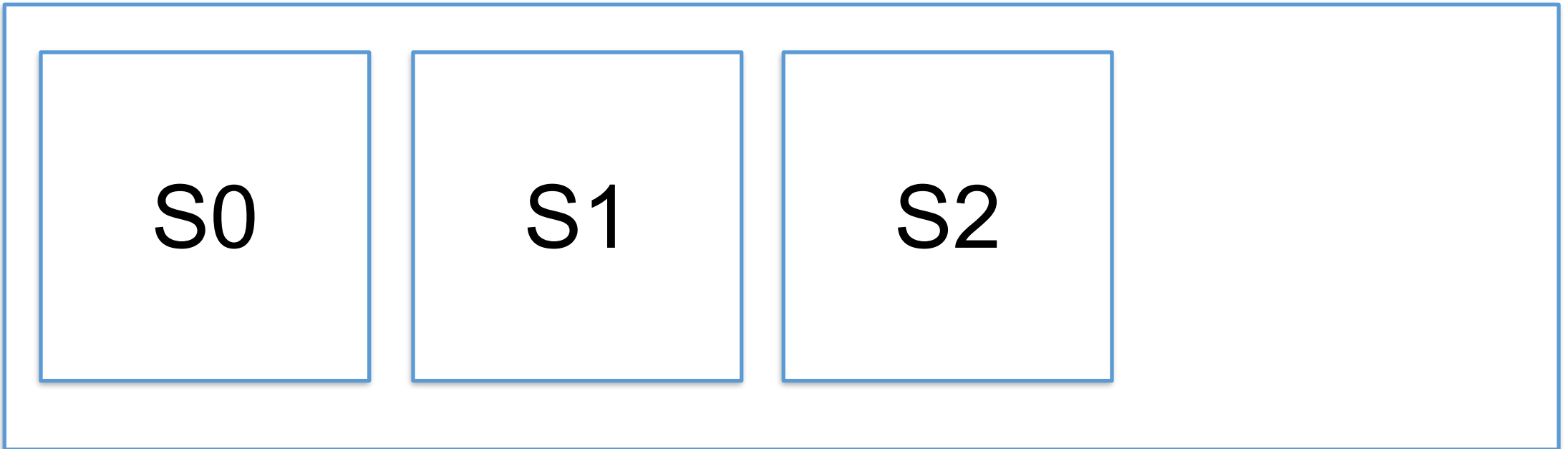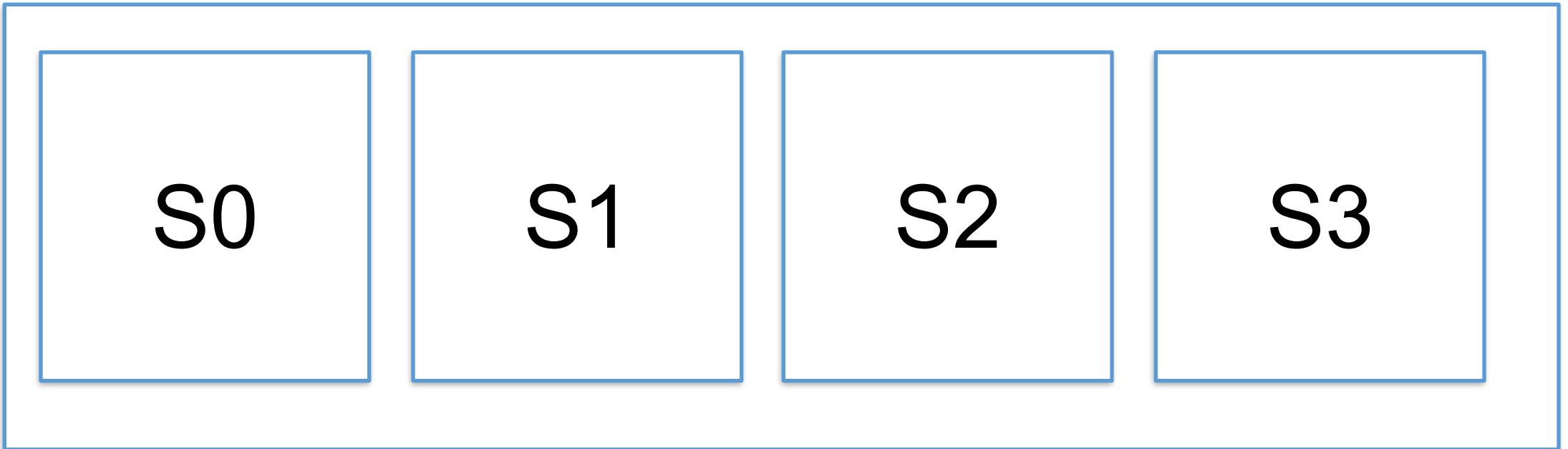
# Real-Time Job Tracing

Period 1

# Real-Time Job Tracing
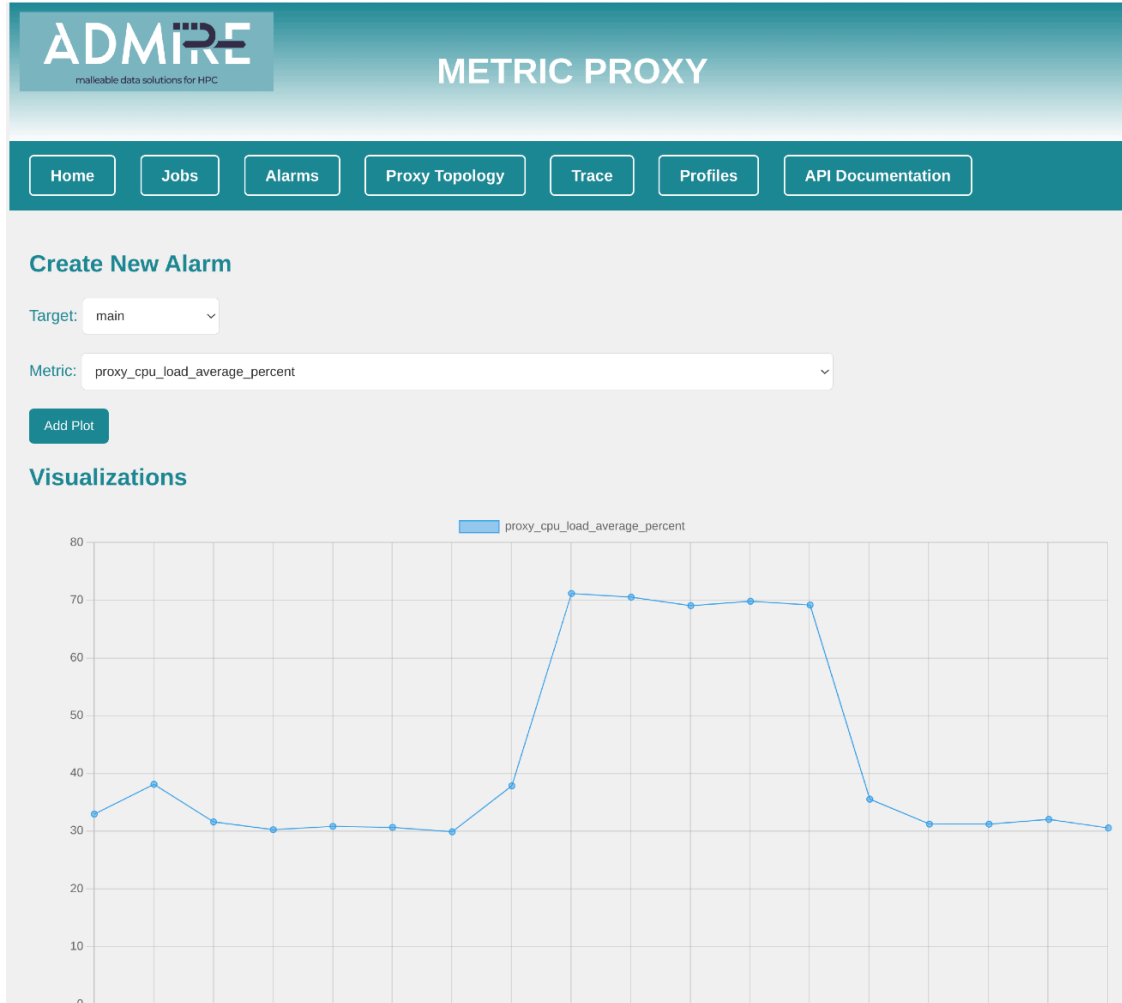


S0    S1    S2

Period 2

S0 + S1    S2 + S3    **Resampling**

# Real-Time Job Tracing

**32 KB Limit**

# Real-Time Job Tracing

**32 KB Limit**

# Job Profile Store



## METRIC PROXY

Home | Jobs | Alarms | Proxy Topology | Trace | Profiles | API Documentation

### Profile Data

#### ./lulesh2.0 -s 30 -i 50 -p

| jobid | command | size | nodelist | partition | cluster | run_dir | start_time | end_time |
|---|---|---|---|---|---|---|---|---|
| 2475294721 (JSON) | ./lulesh2.0 -s 30 -i 50 -p | 1 | | | | /tmp/lulesh-2.0.3 | 1702042634 | 1702042637 |
| 2476998657 (JSON) | ./lulesh2.0 -s 30 -i 50 -p | 8 | | | | /tmp/lulesh-2.0.3 | 1702042637 | 1702042643 |

#### ./lulesh2.0 -s 20 -i 50 -p

| jobid | command | size | nodelist | partition | cluster | run_dir | start_time | end_time |
|---|---|---|---|---|---|---|---|---|
| 2466643969 (JSON) | ./lulesh2.0 -s 20 -i 50 -p | 8 | | | | /tmp/lulesh-2.0.3 | 1702042630 | 1702042634 |
| 2468282369 (JSON) | ./lulesh2.0 -s 20 -i 50 -p | 1 | | | | /tmp/lulesh-2.0.3 | 1702042628 | 1702042630 |

#### ./lulesh2.0 -s 10 -i 50 -p

| jobid | command | size | nodelist | partition | cluster | run_dir | start_time | end_time |
|---|---|---|---|---|---|---|---|---|
| 2423128065 (JSON) | ./lulesh2.0 -s 10 -i 50 -p | 1 | | | | /tmp/lulesh-2.0.3 | 1702042623 | 1702042625 |
| 2425749505 (JSON) | ./lulesh2.0 -s 10 -i 50 -p | 8 | | | | /tmp/lulesh-2.0.3 | 1702042625 | 1702042628 |

# Profile Modelling with Extra-P

| JOB | Profile |
|-----|---------|
| JOB | Profile |
| JOB | Profile |
| JOB | Profile |

Extra-P Model File

Extra-P

Projected Profile

# Complete Prometheus Integration



Spatial Aggregation

Proxy · Proxy · Proxy · Proxy

Proxy · Proxy

Proxy

**Time Series of sample values are accumulated over time:**
- For nodes
- For global state

Prometheus Database

# Sample Prometheus Outputs (Integrated View)

# Sample Grafana Output

# Overhead Assessment (on V1)

# Source Code & Demo

Try it at:

**http://github.com/besnardjb/proxy_v2/**

Demo

# Conclusion

- **We presented the ADMIRE Metric proxy**
  - **Implements an aggregating prometheus push gateway (and more)**
  - **Made in Rust**
- **We have machine wide monitoring capabilities:**
  - **Profiles (per node, per job)**
  - **Traces(per job, per node)**
  - **Machine wide state in real time (0.5 sec resolution)**
- **We work on modelling capabilities thanks to Extra-P and FTIO (WIP)**
- **We will make a first official release in the near future after more testing at scale on the ADMIRE testing supercomputer (thanks to UNITO)**
- **Code available at: http://github.com/besnardjb/proxy_v2/**

Try it at:

**http://github.com/besnardjb/proxy_v2/**