# Multi-Application Online Profiling Tool
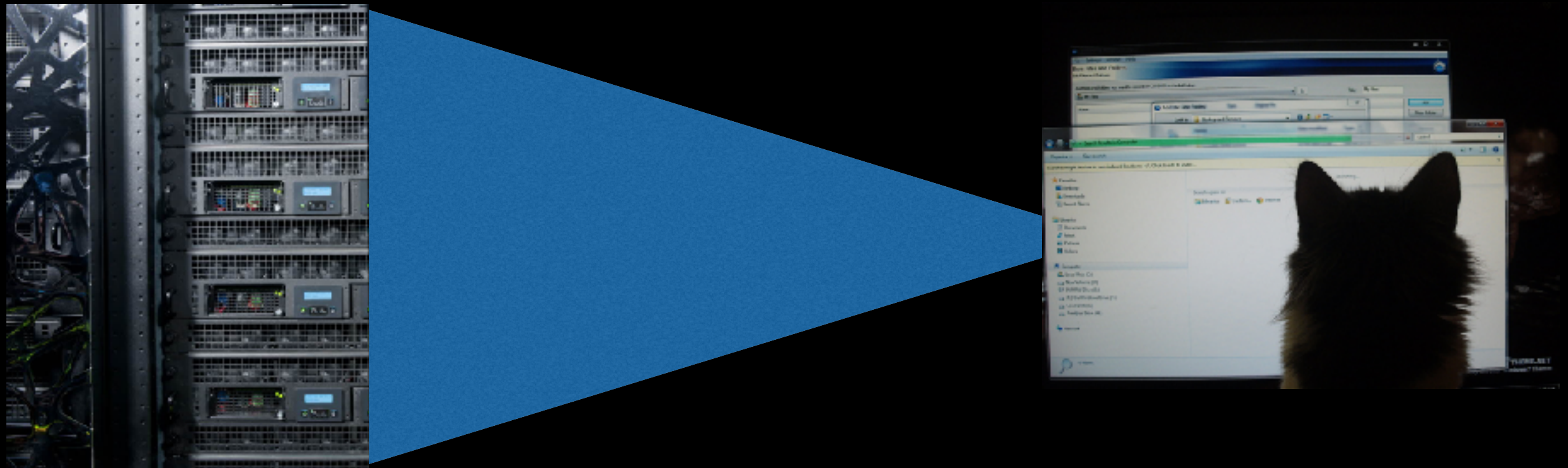
Vi-HPS

Julien ADAM, Antoine CAPRA

ParaTools

# About MALP

- MALP is a tool originally developed in CEA and in the University of Versailles (UVSQ)
- It generates rich HTML views describing parallel programs

- MALP relies on the on-line profiling approach (fine-grained events are reduced during the execution)
- MALP is currently MPI oriented with a focus on the topology and some temporal analysis

ParaTools

# MALP's trade-offs in the Measurement Process
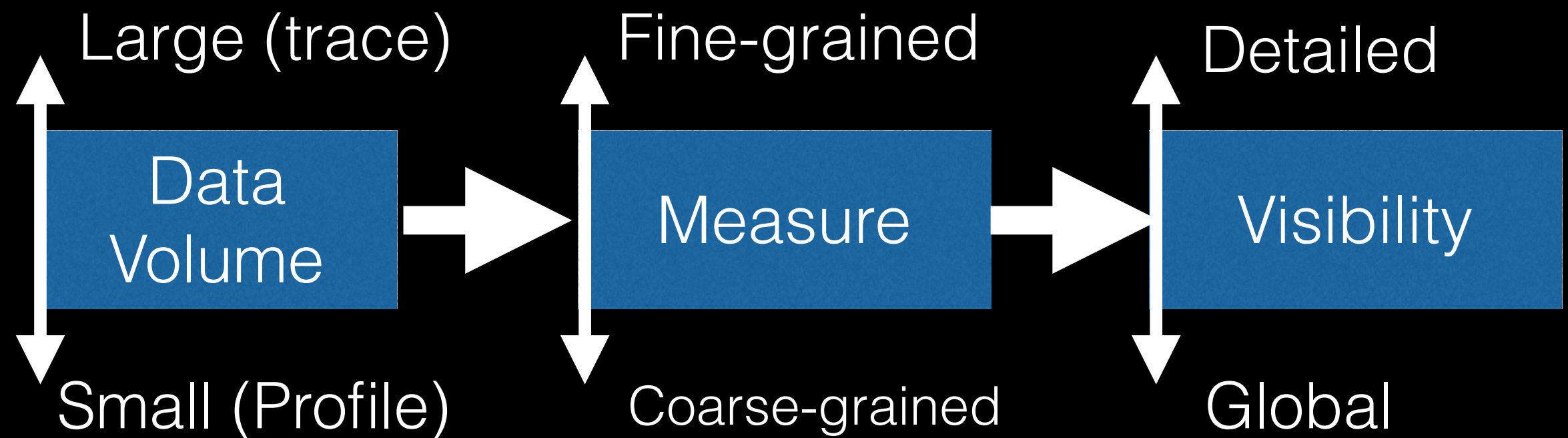
ParaTools

# The Measurement Chain



Representing the state of millions of cores necessarily require a form of **reduction**, at least to acknowledge user's **cognitive limitations**.
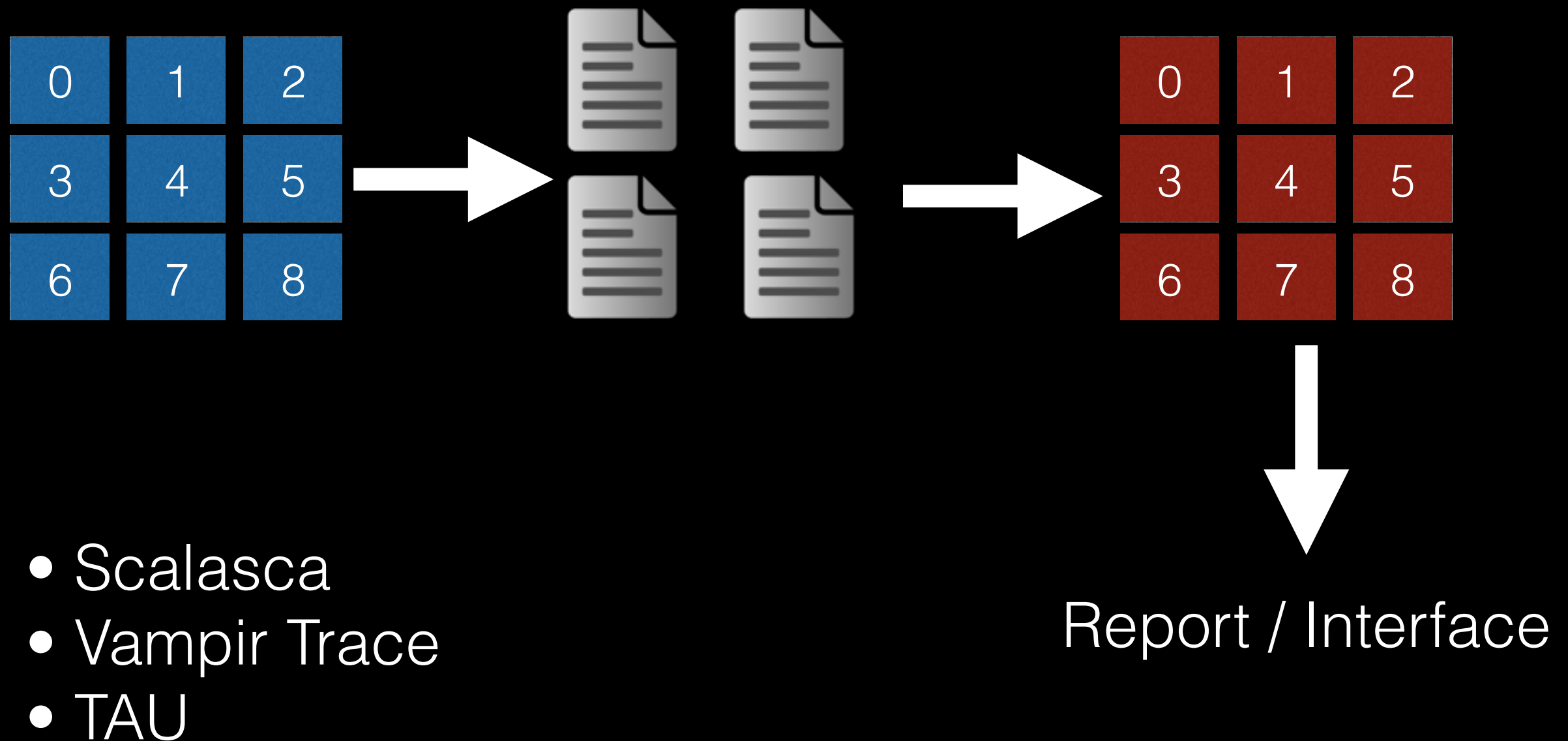
ParaTools

# The Measurement Chain

Large (trace)     Fine-grained     Detailed

**Data Volume** → **Measure** → **Visibility**

Small (Profile)     Coarse-grained     Global

A chain **collecting** and **valorizing** data
*Compromise between details and measurements' scalability*

ParaTools

# Trace-Based Approach

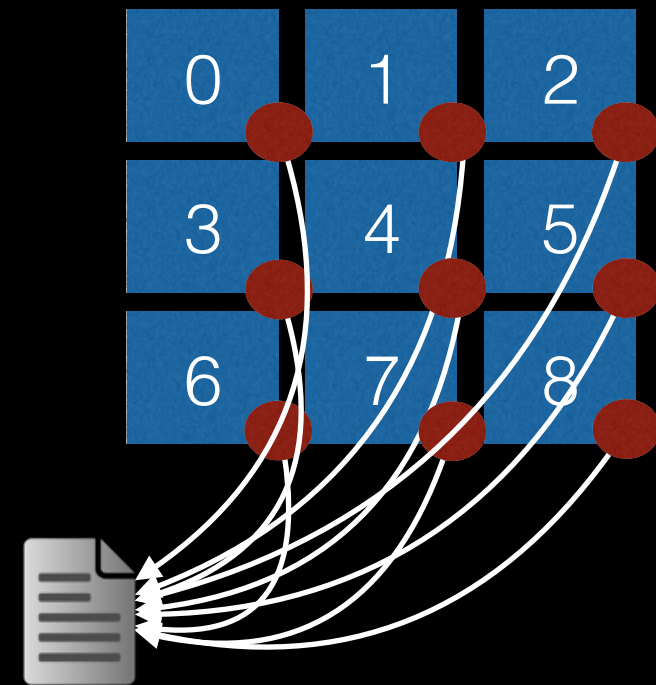| | | |
|---|---|---|
| 0 | 1 | 2 |
| 3 | 4 | 5 |
| 6 | 7 | 8 |

➡️

| | | |
|---|---|---|
| 0 | 1 | 2 |
| 3 | 4 | 5 |
| 6 | 7 | 8 |

- Scalasca
- Vampir Trace
- TAU

Report / Interface

ParaTools

# On-Line



**Tree-Based Overlay Network (TBON)**
- Periscope
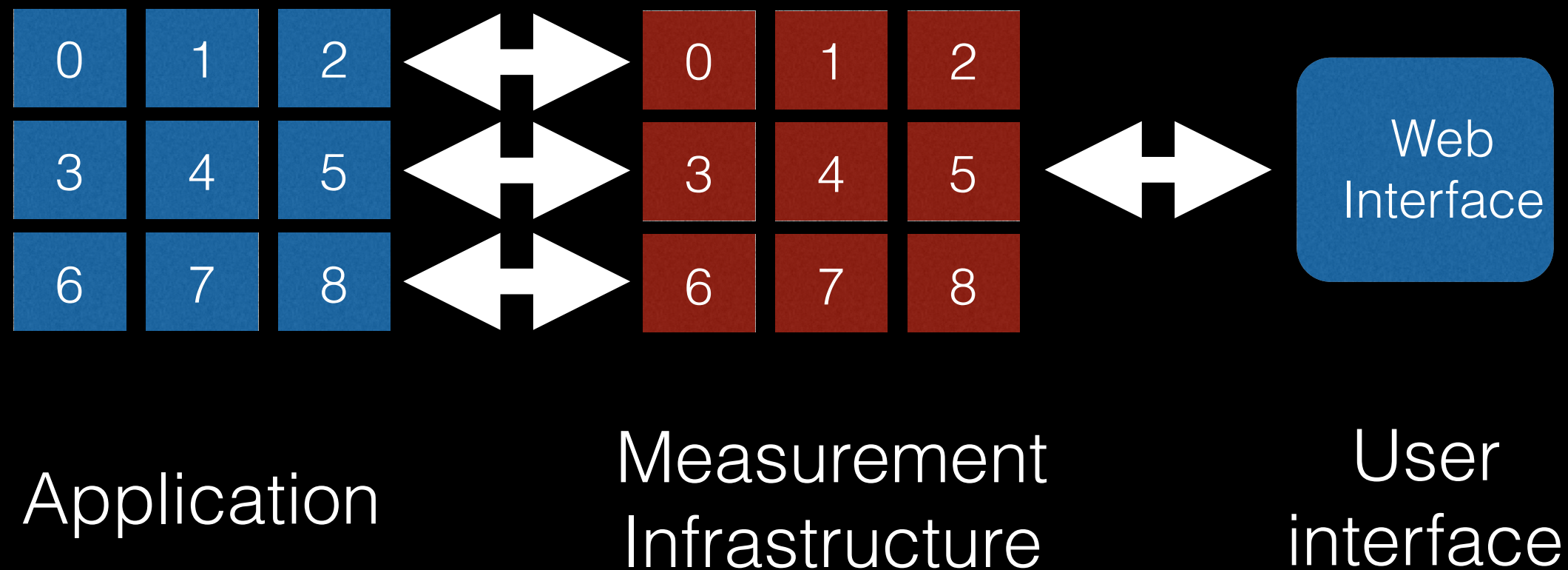- MrNet/Stat (Debug)
- MUST (Validation)
- DDT/Totalview/MAP

**Profiles:**
- Mpi-P
- Scalasca (partly)

# MALP
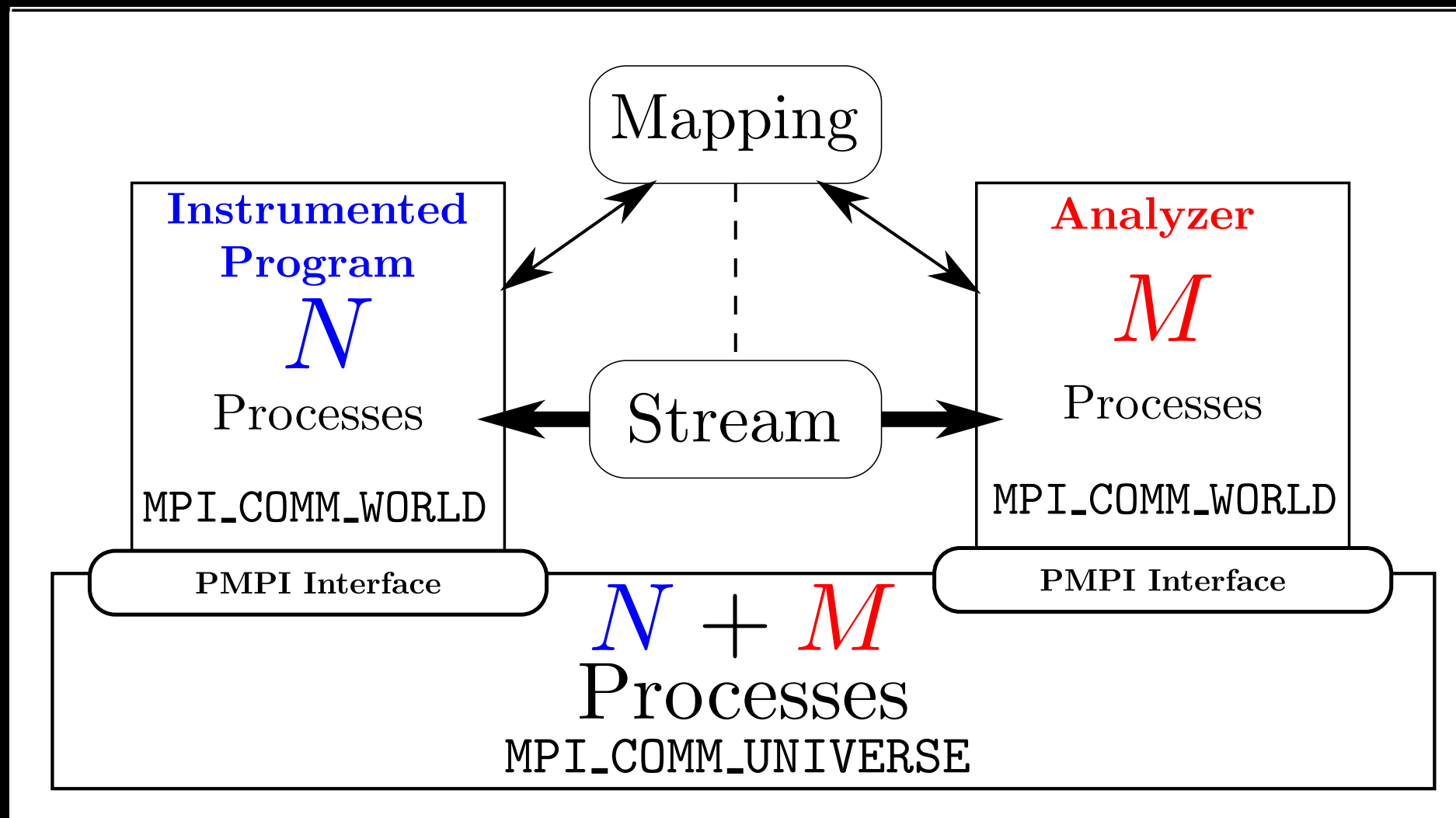


Application

Measurement Infrastructure

User interface

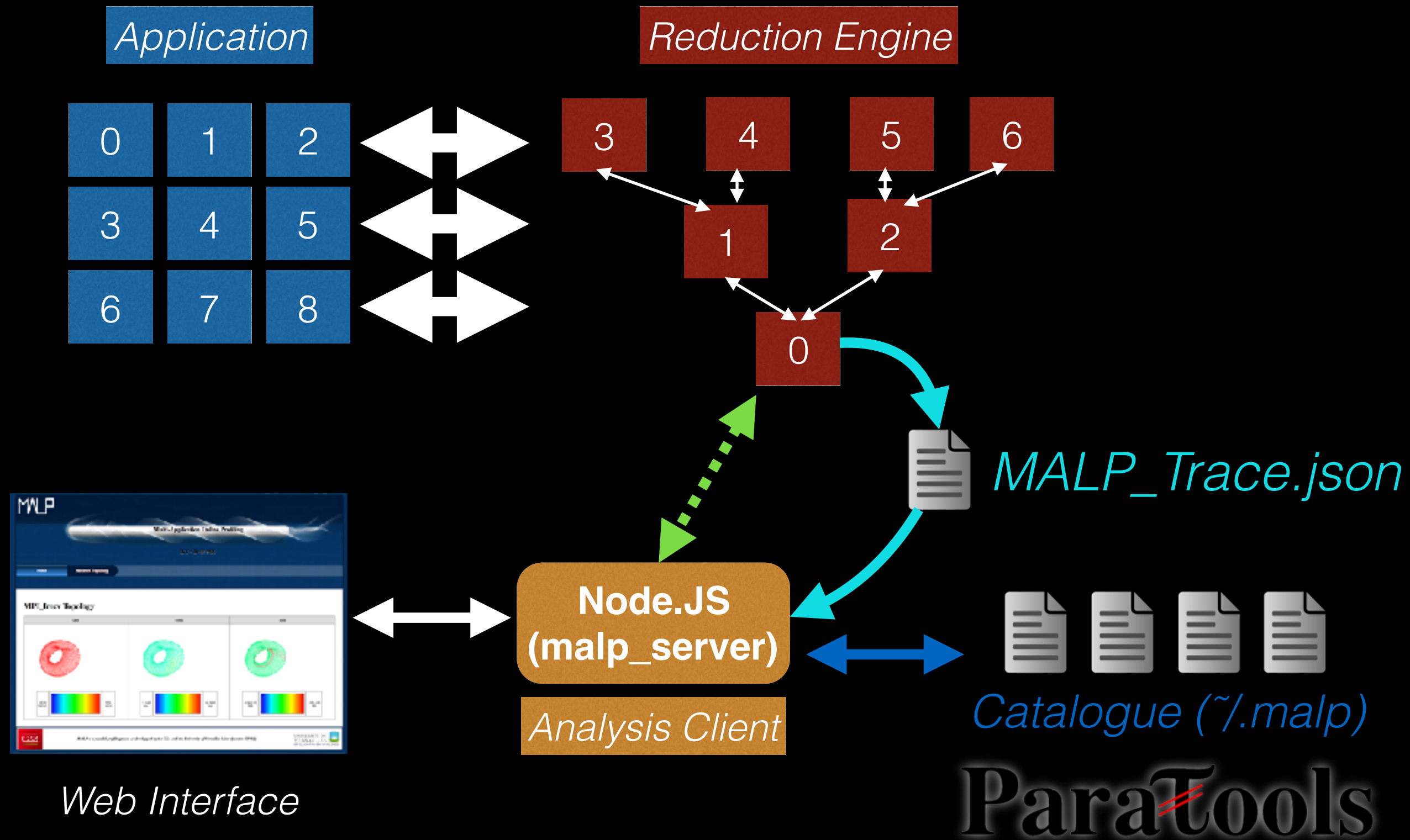Idea: allocate cores to analyze and valorize data (reduction).

# Instrumentation Analysis Coupling
## Overview



MPI Virtualization: an original idea from the P^NMPI paper

# Coupling with the Interface



*Application*

| 0 | 1 | 2 |
| 3 | 4 | 5 |
| 6 | 7 | 8 |

*Reduction Engine*

3  4  5  6

1  2

0

*MALP_Trace.json*

**Node.JS (malp_server)**

*Analysis Client*

*Catalogue (~/.malp)*

*Web Interface*

ParaTools

# Closer look on the Interface Side



TCP or File

Analysis 1

Analysis 2

**Node.JS (malp_server)**

*Analysis Client*

HTTP

Dynamic Views

ParaTools

# MALP vs Others

**Pros**

- MALP is easy to use (preload only)
- Analysis is done online (small data footprint)
- Performance data rendering is done in a browser
- Some original topological analysis
- Overhead is relatively small

**Cons**

- Due to its approach MALP presents reduced data
- Exhaustive program behavior is not provided by MALP for this use a trace-based approach
- Extra processing power is required during the execution

ParaTools

# Outputs from MALP

# Closer look on the Interface Side

# The Upload Form



**MALP Trace Upload form**

Please select a trace file "MALP_trace.json" in order to store it in the MALP server.

Choisir le fichier    aucun fichier sél.        Upload Trace

ParaTools

# The Meta-Data Form

# Trace Uploaded !

**Trace Loaded Sucessfully !**

Content:

- ./a.out

ParaTools

# The Performance Catalog
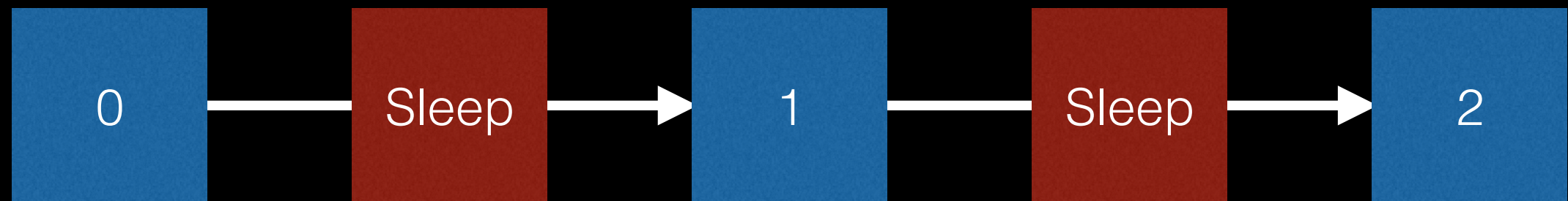
# MPI Profile

# « Per Rank » Analysis
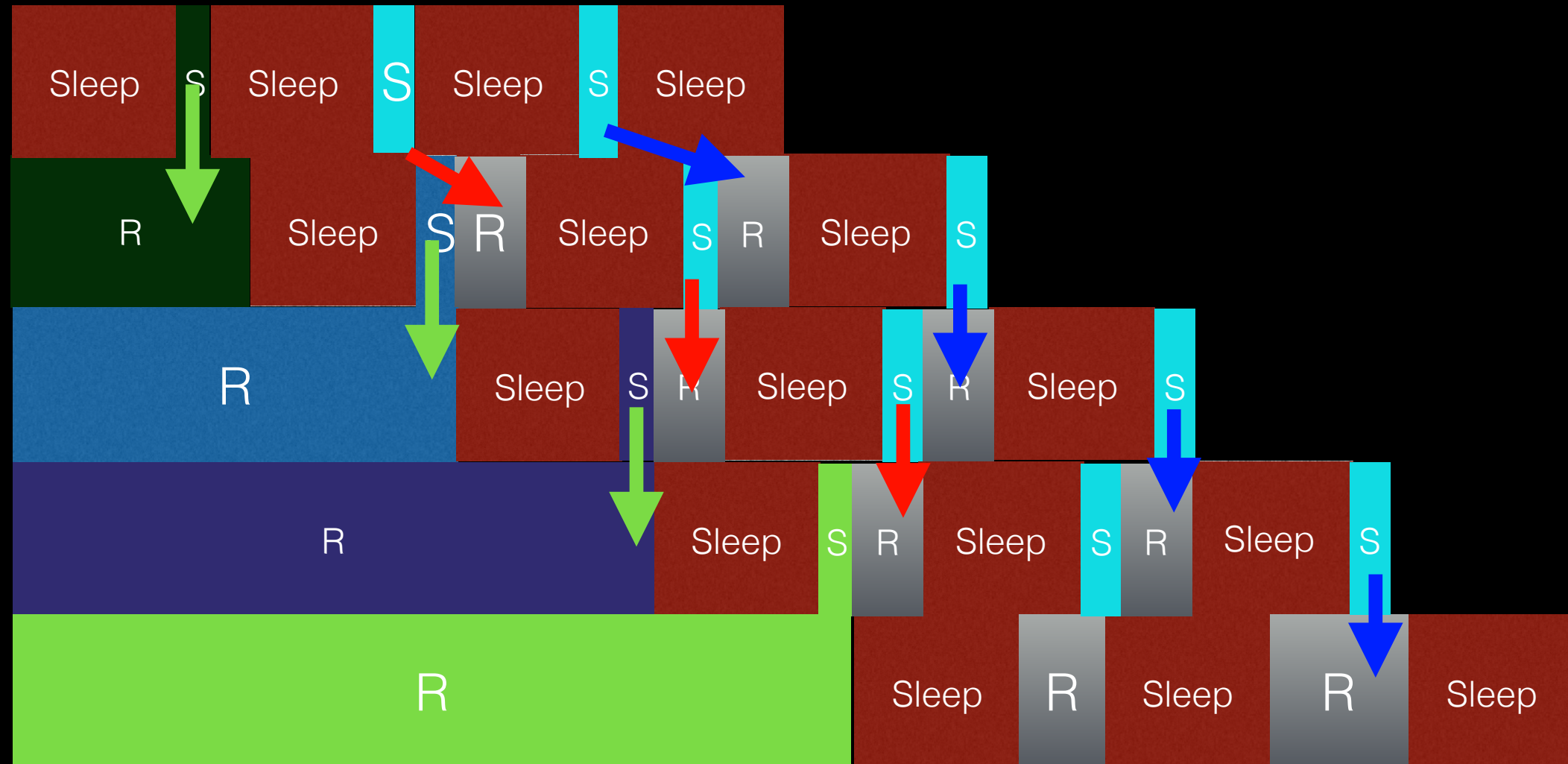
# « Rank Matrix » Analysis

# « Time Matrix » Analysis

Consider this example where you send a message to your right rank after waiting



```
MPI_Recv( &value, 1, MPI_INT, rank - 1, 0, MPI_COMM_WORLD, &status );
usleep(XXX);
if (rank < size - 1)
    MPI_Send( &value, 1, MPI_INT, rank + 1, 0, MPI_COMM_WORLD );
```
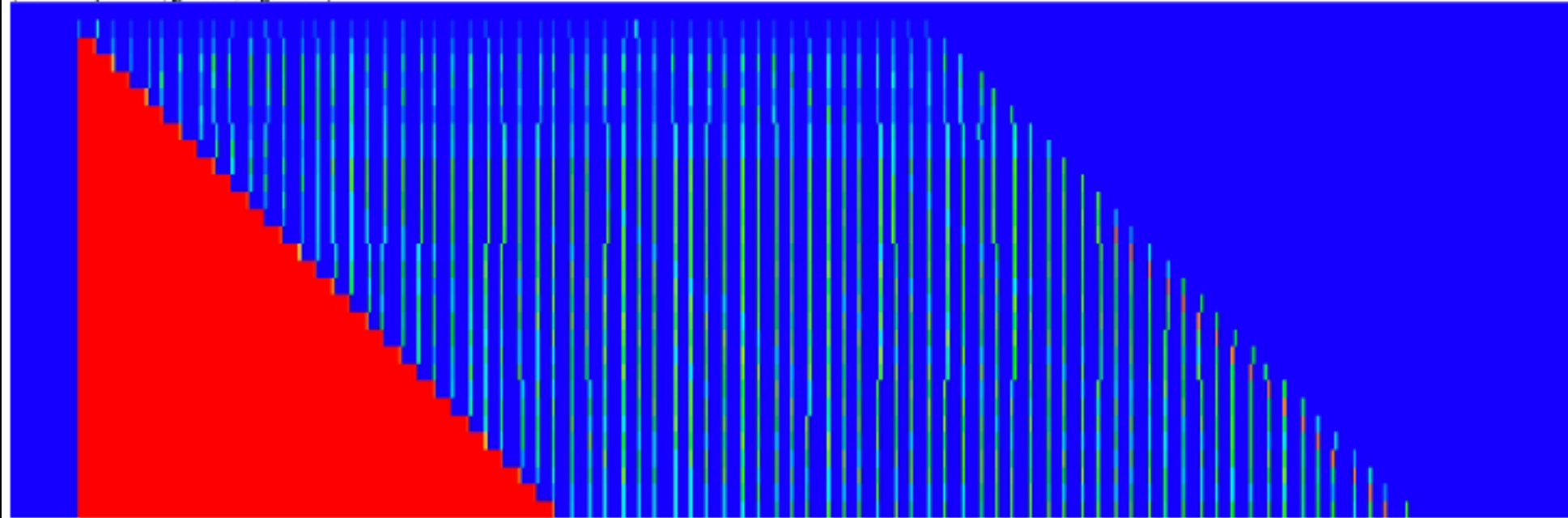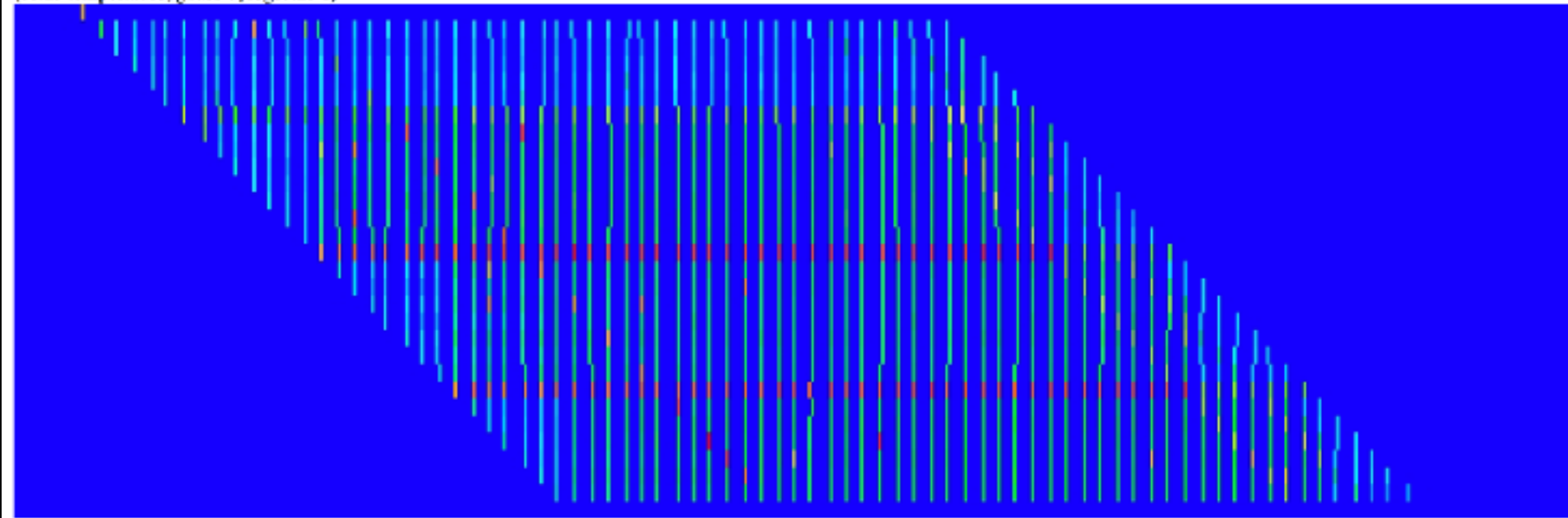
ParaTools

# « Time Matrix » Analysis

# « Time Matrix » Analysis

# « Interactive Topology » Analysis

# « Interactive Topology » Analysis



MG



BT



LU



CG

ParaTools

# Allinea MAP Support

ParaTools

# MALP can consume MAP events



28

# MALP can consume MAP events

Polling Thread

MAP Plugin Array

Handler

MALP MAP_Events

MAP Plugin

MAP Plugin

Metric

Metric

Metric

Metric

ParaTools

# Sample MAP IO stat Collector

```c
int rchar_collect(metric_id_t metric_id,
                  struct timespec *in_out_sample_time,
                  double *out_value)
{
    struct io_stat st;

    if( read_io_stat( &st ) )
    {
        allinea_set_metric_error_messagef(metric_id, errno,
        "Error opening /proc/self/io\n");
        return;
    }

    *out_value = ( st.rchar - prev_values.rchar) / (1024.0 * 1024.0);
    prev_values.rchar = st.rchar;

}
```
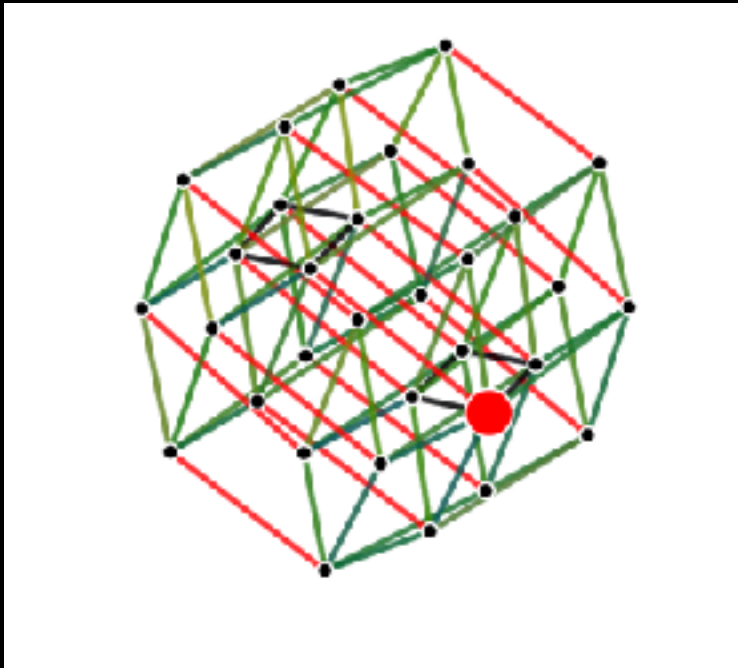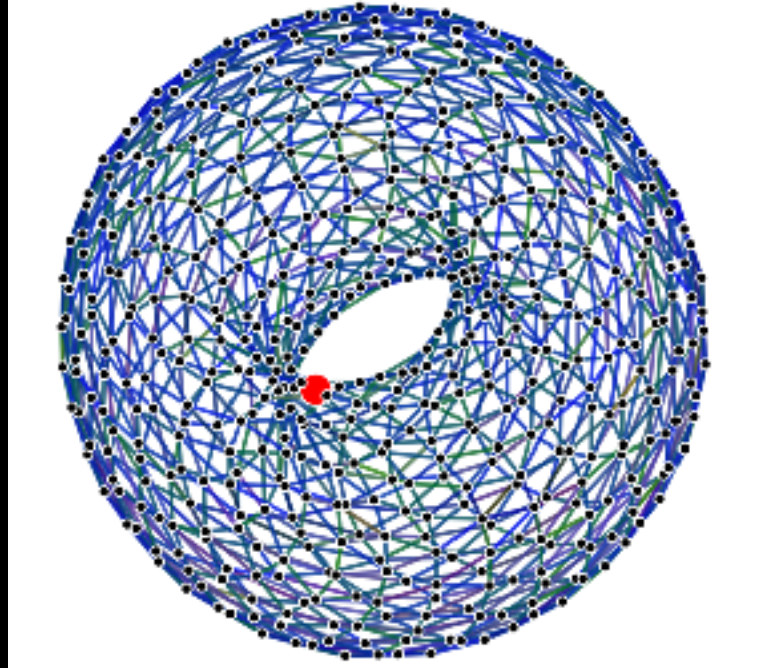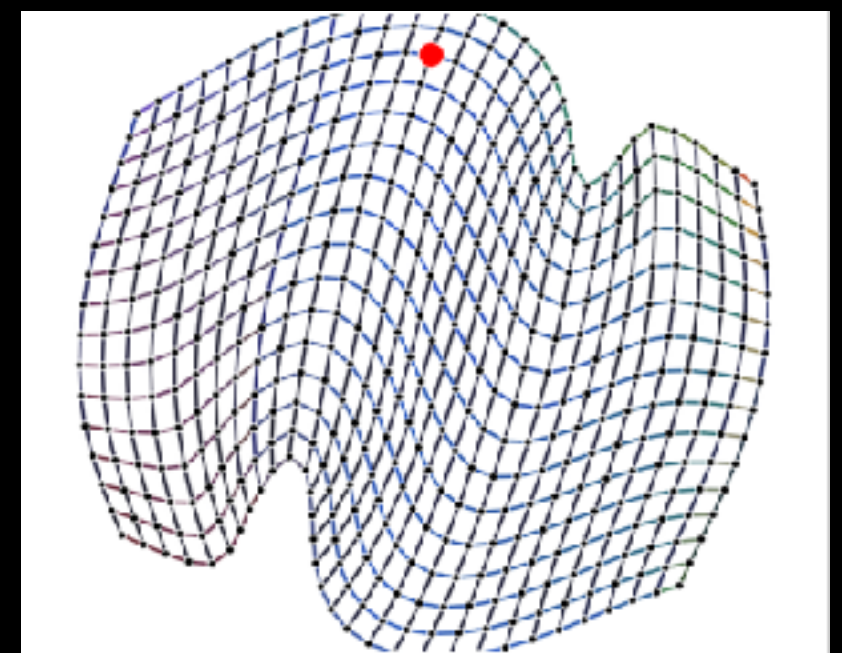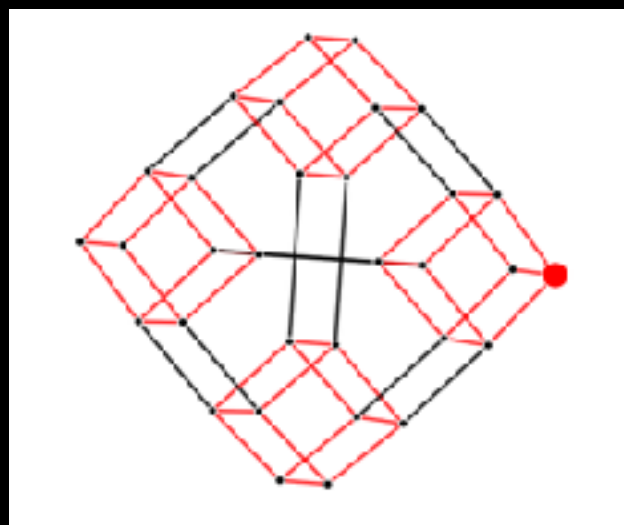
ParaTools

# MPI-IO Write Bandwith Example



**IMB-IO**
**MALP « time matrix »**
**and MAP timeline**

# How to use MALP ?

# Measurement Process in MALP

Two applications have to be co-launched:
- The instrumented program (LD_PRELOAD)
- The analysis engine



Program

Instrum Library

Analysis

ParaTools

# Measurement Process in MALP

*malpl 16 8 CMD [ARGS]*

Will launch 18 processes
(16 processes plus an analysis process for 8 instrumented)

_____

*mpirun -np 16 malp CMD [ARGS] : -np 2 malp_an*
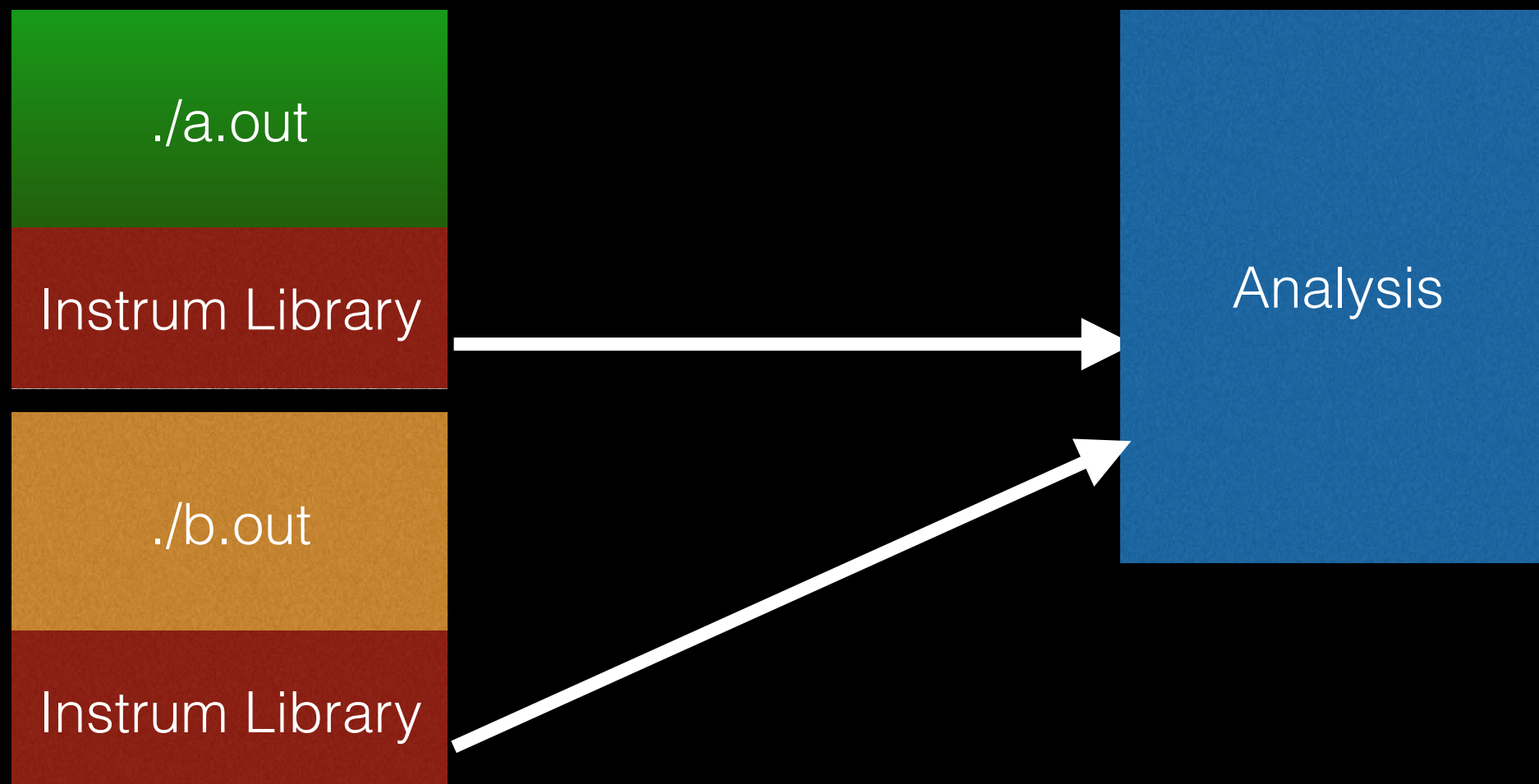
Will launch 18 processes
(same as above classical MPMD syntax)

# Instrumenting Multiple-Applications

*mpirun* *-np 16  malp ./a.out [ARGS] \*
        *: -np 16 malp ./b.out \*
        *: -np 4 malp_an*

Will launch 36 processes

# Enable MAP Plugin Support

Simply point MALP to your plugin directory:

export MAP_PLUGIN_PREFIX=$PREFIX/lib/malp_map_plugins/

You may change sampling frequency (milliseconds):

export MAP_PLUGIN_FREQ=10000

If correctly enabled the following should appear when running the target program:

## Loading MAP plugins :
Loading plugin 'interrupts.xml' ...
Loading plugin 'io.xml' ...
## Done

# Conclusion

- You may download MALP from:

    http://malp.hpcframework.com

You can now use MALP for MPI profiling on your own application or load one of the examples.

ParaTools

# Future Work

Short-term

- We are in the process of adding new MPI-T and OMP-T analysis inside MALP
- An OTF-2 consumer library is also to be released
- A Spack installer is to be released

Medium term

- We are adding call-stack context support inside MALP
- We want to provide time-line views in the browser
- The notion of phase will be used for temporal view

ParaTools

# Multi-Application Online Profiling Tool

Vi-HPS

Julien ADAM, Antoine CAPRA

ParaTools