# MPI/OPENMP COURSE – TOOLS

MARCH 20, 2024  I  MICHAEL KNOBLOCH I M.KNOBLOCH@FZ-JUELICH.DE

JÜLICH | JÜLICH
Forschungszentrum | SUPERCOMPUTING CENTRE

# MUST – MPI CORRECTNESS CHECKER

# HOW MANY ISSUES CAN YOU SPOT?

```c
#include <mpi.h>
#include <stdio.h>

int main(int argc, char** argv) {
    int rank, size, buf[8];

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous(2, MPI_INTEGER, &type);

    MPI_Recv(buf, 2, MPI_INT, size - rank, 123, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
    MPI_Send (buf, 2, type, size - rank, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    return 0;
}
```

At least 8 issues in this example!

JÜLICH | JÜLICH SUPERCOMPUTING
Forschungszentrum | CENTRE

# MOTIVATION

- MPI programming is error prone

- Portabillity errors
  (just on some systems, runs, configurations)

- Bugs may manifest as

  - Crash

  - Application hanging

  - Application finishes

- Questions

  - Why crash/hang?

  - Is the result correct?

  - Will the code produce the correct result on another system?

- Tools help to pin-point these issues

Error more obvious

JÜLICH
Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# TYPES OF ERRORS

- Common syntactic errors:
  - Incorrect arguments
  - Resource usage
  - Lost/Dropped Requests
  - Buffer usage
  - Type-matching
  - Deadlocks

Tool to use:

**MUST,**

**Static analysis tool,**

**(Debugger)**

- Semantic errors that are correct in terms of MPI standard, but do not match the programmer's intention:
  - Displacement/Size/Count errors

Tool to use:

**Debugger**

JÜLICH
Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# MUST

- Next generation MPI correctness and portability checker

- https://www.i12.rwth-aachen.de/go/id/nrbe

- MUST reports
  - Errors: violations of the MPI-standard
  - Warnings: unusual behavior or possible problems
  - Notes: harmless but remarkable behavior
  - Potential deadlock detection

- Usage
  - Compile with debug information (i.e. use the -g flag)
  - Run application under the control of mustrun (requires (at least) one additional MPI process)
    - E.g. on JUSUF: mustrun --must:mpiexec srun --must:np -n -n 4 ./app
  - Open output html report (might need to copy it to your local machine)

JÜLICH
Forschungszentrum

JÜLICH
SUPERCOMPUTING
CENTRE

# FIX 0: ADD MPI_INIT/MPI_FINALIZE

```c
#include <mpi.h>
#include <stdio.h>

int main(int argc, char** argv) {
  int rank, size, buf[8];

  MPI_Init(&argc, &argv);

  MPI_Comm_rank(MPI_COMM_WORLD, &rank);
  MPI_Comm_size(MPI_COMM_WORLD, &size);

  MPI_Datatype type;
  MPI_Type_contiguous(2, MPI_INTEGER, &type);

  MPI_Recv(buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
  MPI_Send (buf, 2, type, size - rank - 1, 123, MPI_COMM_WORLD);

  printf ("Hello, I am rank %d of %d.\n", rank, size);

  MPI_Finalize();

  return 0;
}
```

JÜLICH Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# MUST DETECTS DEADLOCKS

# MUST DETECTS DEADLOCKS

# FIX 1: USE ASYNC RECV

```c
#include <mpi.h>
#include <stdio.h>

int main(int argc, char** argv) {
    int rank, size, buf[8];

    MPI_Init(&argc, &argv);

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous(2, MPI_INTEGER, &type);

    MPI_Request request;
    MPI_Irecv(buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf, 2, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    MPI_Finalize();

    return 0;
}
```

Use asynchronous receive (MPI_Irecv)

JÜLICH Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# MUST DETECTS BUFFER ERRORS

| Rank(s) | Type | |
|---|---|---|
| 2(28793) | Error | A receive operation uses a (datatype... by the send it matches! The first element of the send... |

Details:

| Message | | From | References |
|---|---|---|---|
| A receive operation uses a (datatype,count) pair that can not hold the data transfered by the send it matches! The first element of the send that did not fit into the receive operation is at (contiguous)[0](MPI_INTEGER) in the send type (consult the MUST manual for a detailed description of datatype positions). The send operation was started at reference 1, the receive operation was started at reference 2. (Information on communicator: MPI_COMM_WORLD) (Information on send of count 2 with type:Datatype created at reference 3 is for Fortran, based on the following type(s): { MPI_INTEGER}) (Information on receive of count 2 with type:MPI_INT) | | Representative location: **MPI_Send** (1st occurrence) called from: #0 main@example-fix1.c:18 | References of a representative process:<br><br>reference 1 rank 2: **MPI_Send** (1st occurrence) called from: #0 main@example-fix1.c:18<br><br>reference 2 rank 1: **MPI_Irecv** (1st occurrence) called from: #0 main@example-fix1.c:16<br><br>reference 3 rank 2: **MPI_Type_contiguous** (1st occurrence) called from: #0 main@example-fix1.c:13 |

| Rank(s) | Type | |
|---|---|---|
| 1(28792) | Error | A receive operation uses a (datatype,count) pair that can not hold the data transfered by the send it matches! The first element of the send... |
| 0-3 | Error | Argument 3 (datatype) is not commited for transfer, call MPI_Type_commit before using the type for transfer!(Information on datatypeData... |
| 2(28793) | Error | The memory regions to be transfered by this send operation overlap with regions spanned by a pending non-blocking receive operation!(In... |
| 1(28792) | Error | The memory regions to be transfered by this send operation overlap with regions spanned by a pending non-blocking receive operation!(In... |
| 3(28795) | Error | The memory regions to be transfered by this send... receive operation!(In... |
| 3(28795) | Error | A receive operation uses a (datatype,count) pair that can not... element of the send... |
| 0(28794) | Error | The memory regions to be transfered by this send operation overla... receive operation!(In... |
| 0(28794) | Error | A receive operation uses a (datatype,count) pair that can not hold t... element of the send... |

**Size of sent message larger than receive buffer**

**All detected errors are collapsed for overview - click to expand**

JÜLICH | JÜLICH SUPERCOMPUTING CENTRE
Forschungszentrum

# FIX 2: SAME MESSAGE SIZE FOR SEND/RECV

```c
#include <mpi.h>
#include <stdio.h>

int main(int argc, char** argv) {
    int rank, size, buf[8];

    MPI_Init(&argc, &argv);

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous(2, MPI_INTEGER, &type);

    MPI_Request request;
    MPI_Irecv(buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf, 1, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    MPI_Finalize();

    return 0;
}
```

Reduce the message size

JÜLICH | JÜLICH SUPERCOMPUTING CENTRE
Forschungszentrum

# MUST DETECTS DATATYPE ERRORS

# MUST DETECTS DATATYPE ERRORS



Graphical representation of the type mismatch

# FIX 3+4: C INT TYPE & USE TYPE_COMMIT

```c
#include <mpi.h>
#include <stdio.h>

int main(int argc, char** argv) {
    int rank, size, buf[8];

    MPI_Init(&argc, &argv);

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous(2, MPI_INT, &type);
    MPI_Type_commit(&type);

    MPI_Request request;
    MPI_Irecv(buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf, 1, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    MPI_Finalize();

    return 0;
}
```

Use integer datatype intended for C

Commit datatype before usage

JÜLICH | JÜLICH SUPERCOMPUTING CENTRE
Forschungszentrum

# MUST DETECTS DATARACES IN ASYNC COMM

Data race between send and asynchronous receive operation

| Rank(s) | Type | |
|---|---|---|
| 0(1605) | Error | The memory regions to be transfered by thi... receive operation!(In... |

Details:

| Message | From | References |
|---|---|---|
| The memory regions to be transfered by this send operation overlap with regions spanned by a pending non-blocking receive operation!<br><br>(Information on the request associated with the other communication:<br>Point-to-point request activated at reference 1)<br>(Information on the datatype associated with the other communication:<br>MPI_INT)<br>The other communication overlaps with this communication at position:(MPI_INT)<br><br>(Information on the datatype associated with this communication:<br>Datatype created at reference 2 is for C, commited at reference 3, based on the following type(s): { MPI_INT})<br>This communication overlaps with the other communication at position:(contiguous) [0](MPI_INT)<br><br>A graphical representation of this situation is available in a detailed overlap view (MUST_Output-files/MUST_Overlap_6893422510080_0.html). | Representative location:<br>**MPI_Send** (1st occurrence) called from:<br>#0 main@example-fix3.c:19 | References of a representative process:<br><br>reference 1 rank 0: **MPI_Irecv** (1st occurrence) called from:<br>#0 main@example-fix3.c:17<br><br>reference 2 rank 0: **MPI_Type_contiguous** (1st occurrence) called from:<br>#0 main@example-fix3.c:13<br><br>reference 3 rank 0: **MPI_Type_commit** (1st occurrence) called from:<br>#0 main@example-fix3.c:14 |
| 3(1610) | Error | The memory regions to be transfered by this send operation overlap with regions spanned by a pending non-blocking receive operation!(In... |
| 2(1608) | Error | The memory regions to be transfered by this send operation overlap with regions spanned by a pending non-blocking receive operation!(In... |
| 1(1606) | Error | The memory regions to be transfered by this send operation overlap with regions spanned by a pending non-blocking receive operation!(In... |
| 0-3 | Error | There are 1 datatypes that are not freed when MPI_Finalize was issued, a quality application should free all MPI resources before calling ... |
| 0-3 | Error | There are 1 requests that are not freed when MPI_Finalize was issued, a quality application should free all MPI resources before calling M... |

JÜLICH
Forschungszentrum

JÜLICH
SUPERCOMPUTING
CENTRE

# MUST DETECTS DATARACES IN ASYNC COMM

# FIX 5: USE INDEPENDENT MEMORY REGIONS

```c
#include <mpi.h>
#include <stdio.h>

int main(int argc, char** argv) {
    int rank, size, buf[8];

    MPI_Init(&argc, &argv);

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous(2, MPI_INT, &type);
    MPI_Type_commit(&type);

    MPI_Request request;
    MPI_Irecv(buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf + 2, 1, type, size - rank - 1, 123, MPI_COMM_WORLD);

    printf ("Hello, I am rank %d of %d.\n", rank, size);

    MPI_Finalize();

    return 0;
}
```

Offset points to independent memory

JÜLICH Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# MUST DETECTS LEAKS OF USER-DEFINED OBJECTS

| Rank(s) | Type | Message |
|---|---|---|
| 0-3 | **Error** | There are 1 datatypes that are not freed when MPI_Finalize was issued, a quality application should free all MPI resources before calling ... |

Details:

| Message | From | References |
|---|---|---|
| There are 1 datatypes that are not freed when MPI_Finalize was issued, a quality application should free all MPI resources before calling MPI_Finalize. Listing information for these datatypes:<br><br>-Datatype 1: Datatype created at reference 1 is for C, commited at reference 2, based on the following type(s): { MPI_INT} | Representative location:<br>**MPI_Type_contiguous** (1st occurrence) called from:<br>#0 main@example-fix4.c:13 | References of a representative process:<br><br>reference 1 rank 1: **MPI_Type_contiguous** (1st occurrence) called from:<br>#0 main@example-fix4.c:13<br><br>reference 2 rank 1: **MPI_Type_commit** (1st occurrence) called from:<br>#0 main@example-fix4.c:14 |

| 0-3 | **Error** | There are 1 requests that are not freed when MPI_Finalize was issued, a quality application should free all MPI resources before calling M... |

Details:

| Message | From | References |
|---|---|---|
| There are 1 requests that are not freed when MPI_Finalize was issued, a quality application should free all MPI resources before calling MPI_Finalize. Listing information for these requests:<br><br>-Request 1: Point-to-point request activated at reference 1 | Representative location:<br>**MPI_Irecv** (1st occurrence) called from:<br>#0 main@example-fix4.c:17 | References of a representative process:<br><br>reference 1 rank 1: **MPI_Irecv** (1st occurrence) called from:<br>#0 main@example-fix4.c:17 |

Unfinished non-blocking receive is resource leak and missing synchronization

Leak of user defined datatype object

- User defined objects include
  - MPI_Comms (even by MPI_Comm_dup)
  - MPI_Datatypes
  - MPI_Groups

JÜLICH
Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# FIX 6+7: USE MPI_WAIT & FREE DATATYPE

```c
#include <mpi.h>
#include <stdio.h>

int main(int argc, char** argv) {
    int rank, size, buf[8];

    MPI_Init(&argc, &argv);

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    MPI_Datatype type;
    MPI_Type_contiguous(2, MPI_INT, &type);
    MPI_Type_commit(&type);

    MPI_Request request;
    MPI_Irecv(buf, 2, MPI_INT, size - rank - 1, 123, MPI_COMM_WORLD, &request);

    MPI_Send (buf + 2, 1, type, size - rank - 1, 123, MPI_COMM_WORLD);
    MPI_Wait(&request, MPI_STATUS_IGNORE);

    printf ("Hello, I am rank %d of %d.\n", rank, size);
    MPI_Type_free(&type);

    MPI_Finalize();

    return 0;
}
```

Finish asynchronous communication

Deallocate datatype

JÜLICH | JÜLICH SUPERCOMPUTING CENTRE
Forschungszentrum

# FINALLY

| Rank(s) | Type | Message |
|---|---|---|
| | Information | MUST detected no MPI usage errors nor any suspicious behavior during this application run. |

Details:

| Message | From | References |
|---|---|---|
| MUST detected no MPI usage errors nor any suspicious behavior during this application run. | | |

No further error detected

Hopefully this message applies to many applications

JÜLICH
Forschungszentrum

JÜLICH
SUPERCOMPUTING
CENTRE

# PERFORMANCE ANALYSIS
## USING THE SCORE-P ECOSYSTEM

# MOTIVATION

- Writing parallel code is hard

- Writing **fast/efficient** parallel code is even harder

- "Parallel" (multi core/node) performance factors

  - Partitioning / decomposition

    - ☞ Load balancing

  - Communication (i.e., message passing)

  - Multithreading

  - Core binding / NUMA

  - Synchronization / locking

  - I/O

    - ☞Parallel I/O matters

**JÜLICH** Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# TUNING BASICS

- Carefully set various tuning parameters

  - The right (parallel) algorithms and libraries

  - Compiler flags and directives

  - Correct machine usage (mapping and bindings)

    ☞Get the most performance before tuning!

- Measurement is better than guessing

  - To determine performance bottlenecks

  - To compare alternatives

  - To validate tuning decisions and optimizations

    ☞After each step!

JÜLICH
Forschungszentrum

JÜLICH
SUPERCOMPUTING
CENTRE

# PERFORMANCE ENGINEERING WORKFLOW

```
        ┌──────────────┐
   ┌───▶│  Preparation │────┐
   │    └──────────────┘    │
   │           │            │
   │    ┌──────────────┐    │
   │    │ Measurement  │    │
   │    └──────────────┘    │
   │           │            │
   │    ┌──────────────┐    │
   │    │   Analysis   │    │
   │    └──────────────┘    │
   │           │            │
   │    ┌──────────────┐    │
   │    │ Examination  │    │
   │    └──────────────┘    │
   │           │            │
   │    ┌──────────────┐    │
   │    │ Optimization │    │
   │    └──────────────┘    │
   └────────────────────────┘
```

- Prepare application (with symbols), insert extra code (probes/hooks)

- Collection of data relevant to execution performance analysis

- Calculation of metrics, identification of performance metrics

- Presentation of results in an intuitive/understandable form

- Modifications intended to eliminate/reduce performance problems

JÜLICH Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# THE 80/20 RULE

- Programs typically spend 80% of their time in 20% of the code

  ☞ *Know what matters!*

- Developers typically spend 20% of their effort to get 80% of the total speedup possible for the application

  ☞ *Know when to stop!*

- Don't optimize what does not matter

  ☞ *Make the common case fast!*

JÜLICH Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# PERFORMANCE MEASUREMENT

**Two dimensions**

**When** performance measurement is triggered

- **External trigger** (asynchronous)
  - **Sampling**
    - Trigger: Timer interrupt OR Hardware counters overflow

- **Internal trigger** (synchronous)
  - Code **instrumentation** (automatic or manual)

**How** performance data is recorded

- **Profile**
  - Summation of events over time

- **Trace**
  - Sequence of events over time

JÜLICH Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# Score-P

**Scalable performance measurement infrastructure for parallel codes**

- Community-developed open-source
- Replaced tool-specific instrumentation and measurement components of partners
- http://www.score-p.org

TECHNISCHE UNIVERSITÄT DRESDEN

JÜLICH Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

German Research School for Simulation Sciences

RWTH AACHEN UNIVERSITY

TUM Technische Universität München

UNIVERSITY OF OREGON

JÜLICH Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# Score-P ARCHITECTURE

JÜLICH Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

| Vampir | Scalasca | CUBE | TAU | Extra-P | TAUdb |
|--------|----------|------|-----|---------|-------|

Event traces (OTF2)

Call-path profiles (CUBE4, TAU)

**Score-P measurement infrastructure**

Hardware counter (PAPI, rusage, PERF, plugins)

Instrumentation wrapper

| Process-level parallelism (MPI, SHMEM) | Thread-level parallelism (OpenMP, Pthreads) | Accelerator-based parallelism (CUDA, OpenACC, ROCm, OpenCL, Kokkos) | I/O Activity Recording (Posix I/O, MPI-IO) | Source code instrumentation (Compiler, PDT, User) | Sampling interrupts (PAPI, PERF) |

Application

# Score-P FUNCTIONALITY

- Provide typical functionality for HPC performance tools

- **Instrumentation** (various methods)

  - Multi-process paradigms (MPI, SHMEM)

  - Thread-parallel paradigms (OpenMP, POSIX threads)

  - Accelerator-based paradigms (OpenACC, CUDA, OpenCL. Kokkos)

  - **In any combination!**

- Flexible **measurement** without re-compilation:

  - Basic and advanced **profile** generation ($\Rightarrow$ CUBE4 format)

  - Event **trace** recording ($\Rightarrow$ OTF2 format)

- Highly scalable I/O functionality

- Support all fundamental concepts of partner's tools

JÜLICH
Forschungszentrum

JÜLICH
SUPERCOMPUTING
CENTRE

# CUBE EXAMPLE

# SCORE-P: ADVANCED FEATURES

- Measurement can be extensively configured via environment variables

- Allows for targeted measurements:
  - Selective recording
  - Phase profiling
  - Parameter-based profiling
  - …

- GPU support: CUDA, OpenACC, OpenCL, HIP, Kokkos, …

- Please ask us or see the user manual for details

JÜLICH
Forschungszentrum

JÜLICH
SUPERCOMPUTING
CENTRE

# SCALASCA

- **Sc**alable **A**nalysis of **La**rge **Sc**ale **A**pplications

- Approach

  - **Instrument** C, C++, and Fortran parallel applications (**with Score-P**)

  - Option 1: **scalable call-path profiling**

  - Option 2: **scalable event trace analysis**

    - **Collect** event traces

    - **Process trace in parallel**

      - Wait-state analysis

      - Delay and root-cause analysis

      - Critical path analysis

    - **Categorize and rank** results

# AUTOMATIC TRACE ANALYSIS

- Automatic search for patterns of inefficient behaviour

- Classification of behaviour & quantification of significance

- Identification of delays as root causes of inefficiencies



- Guaranteed to cover the entire event trace

- Quicker than manual/visual trace analysis

- Parallel replay analysis exploits available memory & processors to deliver scalability

# EXAMPLE MPI WAIT STATES



(a) Late Sender

(b) Late Receiver

(c) Late Sender / Wrong Order

(d) Wait at N x N

ENTER   EXIT   SEND   RECV   COLLEXIT

# SCALASCA ROOT CAUSE ANALYSIS

- **Root-cause analysis**

  - Wait states typically caused by load or communication imbalances earlier in the program

  - Waiting time can also propagate (e.g., indirect waiting time)

  - Enhanced performance analysis to find the root cause of wait states

- **Approach**

  - Distinguish between direct and indirect waiting time

  - Identify call path/process combinations delaying other processes and causing first order waiting time

  - Identify original **delay**

# SCALASCA TRACE ANALYSIS EXAMPLE



Additional wait-state metrics from the trace analysis

Delay / root-cause metrics

Critical-path profile

JÜLICH | JÜLICH SUPERCOMPUTING CENTRE
Forschungszentrum

# VAMPIR EVENT TRACE VISUALIZER

- Offline trace visualization for Score-Ps OTF2 trace files

- Visualization of MPI, OpenMP and application events:

  - All diagrams highly customizable (through context menus)

  - Large variety of displays for ANY part of the trace

- http://www.vampir.eu


- Advantage:

  - Detailed view of dynamic application behavior

- Disadvantage:

  - Completely manual analysis

  - Too many details can hide the relevant parts

JÜLICH
Forschungszentrum

JÜLICH
SUPERCOMPUTING
CENTRE

# EVENT TRACE VISUALIZATION WITH VAMPIR

- Visualization of dynamic runtime behaviour at any level of detail along with statistics and performance metrics

- Alternative and supplement to automatic analysis

- **Typical questions that Vampir helps to answer**

  - What happens in my application execution during a given time in a given process or thread?

  - How do the communication patterns of my application execute on a real system?

  - Are there any imbalances in computation, I/O or memory usage and how do they affect the parallel execution of my application?

▪ **Timeline charts**
  - Application activities and communication along a time axis



▪ **Summary charts**
  - Quantitative results for the currently selected time interval

JÜLICH
Forschungszentrum

JÜLICH
SUPERCOMPUTING
CENTRE

# VAMPIR PERFORMANCE CHARTS

## Timeline Charts

| | | |
|---|---|---|
| Master Timeline | ➡ | *all threads' activities* |
| Process Timeline | ➡ | *single thread's activities* |
| Summary Timeline | ➡ | *all threads' function call statistics* |
| Performance Radar | ➡ | *all threads' performance metrics* |
| Counter Data Timeline | ➡ | *single threads' performance metrics* |
| I/O Timeline | ➡ | *all threads' I/O activities* |

## Summary Charts

| | |
|---|---|
| Function Summary | Process Summary |
| Message Summary | Communication Matrix View |
| I/O Summary | Call Tree |

JÜLICH | JÜLICH SUPERCOMPUTING CENTRE
Forschungszentrum

# VAMPIR DISPLAYS

# TOOLS DEMO: BT-MZ WITH SCORE-P

# TYPICAL PERFORMANCE ANALYSIS PROCEDURE

- Do I have a performance problem at all?

    - Time / speedup / scalability measurements

- *What* is the key bottleneck (computation / communication)?

    - MPI / OpenMP / flat profiling

- *Where* is the key bottleneck?

    - Call-path profiling, detailed basic block profiling

- *Why* is it there?

    - Hardware counter analysis

    - Trace selected parts (to keep trace size manageable)

- Does the code have scalability problems?

    - Load imbalance analysis, compare profiles at various sizes function-by-function, performance modeling

JÜLICH
Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# WHAT IS THE KEY BOTTLENECK?

- Generate flat MPI profile using Score-P/Scalasca

  - Only requires re-linking

  - Low runtime overhead

- Provides detailed information on MPI usage

  - How much time is spent in which operation?

  - How often is each operation called?

  - How much data was transferred?

- Limitations:

  - Computation on non-master threads and outside of MPI_Init/MPI_Finalize scope ignored

JÜLICH Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# FLAT MPI PROFILE: RECIPE

1. Prefix your *link command* with
   "`scorep --nocompiler`"

2. Prefix your MPI *launch command* with
   "`scalasca -analyze`"

3. After execution, examine analysis results using
   "`scalasca -examine scorep_<title>`"

JÜLICH
Forschungszentrum

JÜLICH
SUPERCOMPUTING
CENTRE

# FLAT MPI PROFILE: EXAMPLE (CONT.)

# WHERE IS THE KEY BOTTLENECK?

- Generate call-path profile using Score-P/Scalasca

  - Requires re-compilation

  - Runtime overhead depends on application characteristics

  - Typically needs some care setting up a good measurement configuration

    - Filtering

    - Selective instrumentation

- Option 1 (recommended for beginners):
  Automatic compiler-based instrumentation

- Option 2 (for in-depth analysis):
  Manual instrumentation of interesting phases, routines, loops

**JÜLICH** Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# CALL-PATH PROFILE: RECIPE

1. Prefix your ***compile & link commands*** with
   "`scorep`"

2. Prefix your MPI ***launch command*** with
   "`scalasca -analyze`"

3. After execution, compare overall runtime with uninstrumented run to determine overhead

4. If overhead is too high

   1. Score measurement using
      "`scalasca -examine -s scorep_<title>`"

   2. Prepare filter file

   3. Re-run measurement with filter applied using prefix
      "`scalasca -analyze -f <filter_file>`"

5. After execution, examine analysis results using
   "`scalasca -examine scorep_<title>`"

JÜLICH
Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# CALL-PATH PROFILE: EXAMPLE (CONT.)

```
% scalasca -examine -s scorep_myprog_Ppnxt_sum
scorep-score -r ./scorep_myprog_Ppnxt_sum/profile.cubex
INFO: Score report written to ./scorep_myprog_Ppnxt_sum/scorep.score
```

- Estimates trace buffer requirements

- Allows to identify canditate functions for filtering

  ☞Computational routines with high visit count
  and low time-per-visit ratio

- Region/call-path classification

  - MPI (pure MPI library functions)

  - OMP (pure OpenMP functions/regions)

  - USR (user-level source local computation

  - COM ("combined" USR + OpeMP/MPI)

  - ANY/ALL (aggregate of all region types)

JÜLICH
Forschungszentrum

JÜLICH
SUPERCOMPUTING
CENTRE

# CALL-PATH PROFILE: EXAMPLE (CONT.)

```
% less scorep_myprog_Ppnxt_sum/scorep.score
Estimated aggregate size of event trace:                          162GB
Estimated requirements for largest trace buffer (max_buf): 2758MB
Estimated memory requirements (SCOREP_TOTAL_MEMORY):         2822MB
(hint: When tracing set SCOREP_TOTAL_MEMORY=2822MB to avoid
 intermediate flushes or reduce requirements using USR regions
 filters.)

flt type      max_buf[B]          visits  time[s] time[%] time/        region
                                                          visit[us]

     ALL 2,891,417,902 6,662,521,083 36581.51   100.0      5.49   ALL
     USR 2,858,189,854 6,574,882,113 13618.14    37.2      2.07   USR
     OMP    54,327,600    86,353,920 22719.78    62.1    263.10   OMP
     MPI       676,342       550,010   208.98     0.6    379.96   MPI
     COM       371,930       735,040    34.61     0.1     47.09   COM

     USR   921,918,660 2,110,313,472  3290.11     9.0      1.56   matmul_sub
     USR   921,918,660 2,110,313,472  5914.98    16.2      2.80   binvcrhs
     USR   921,918,660 2,110,313,472  3822.64    10.4      1.81   matvec_sub
     USR    41,071,134    87,475,200   358.56     1.0      4.10   lhsinit
     USR    41,071,134    87,475,200   145.42     0.4      1.66   binvrhs
     USR    29,194,256    68,892,672    86.15     0.2      1.25   exact_solution
     OMP     3,280,320     3,293,184    15.81     0.0      4.80   !$omp parallel
     [...]
```

Forschungszentrum | CENTRE

# CALL-PATH PROFILE: FILTERING

- In this example, the 6 most frequently called routines are of type USR
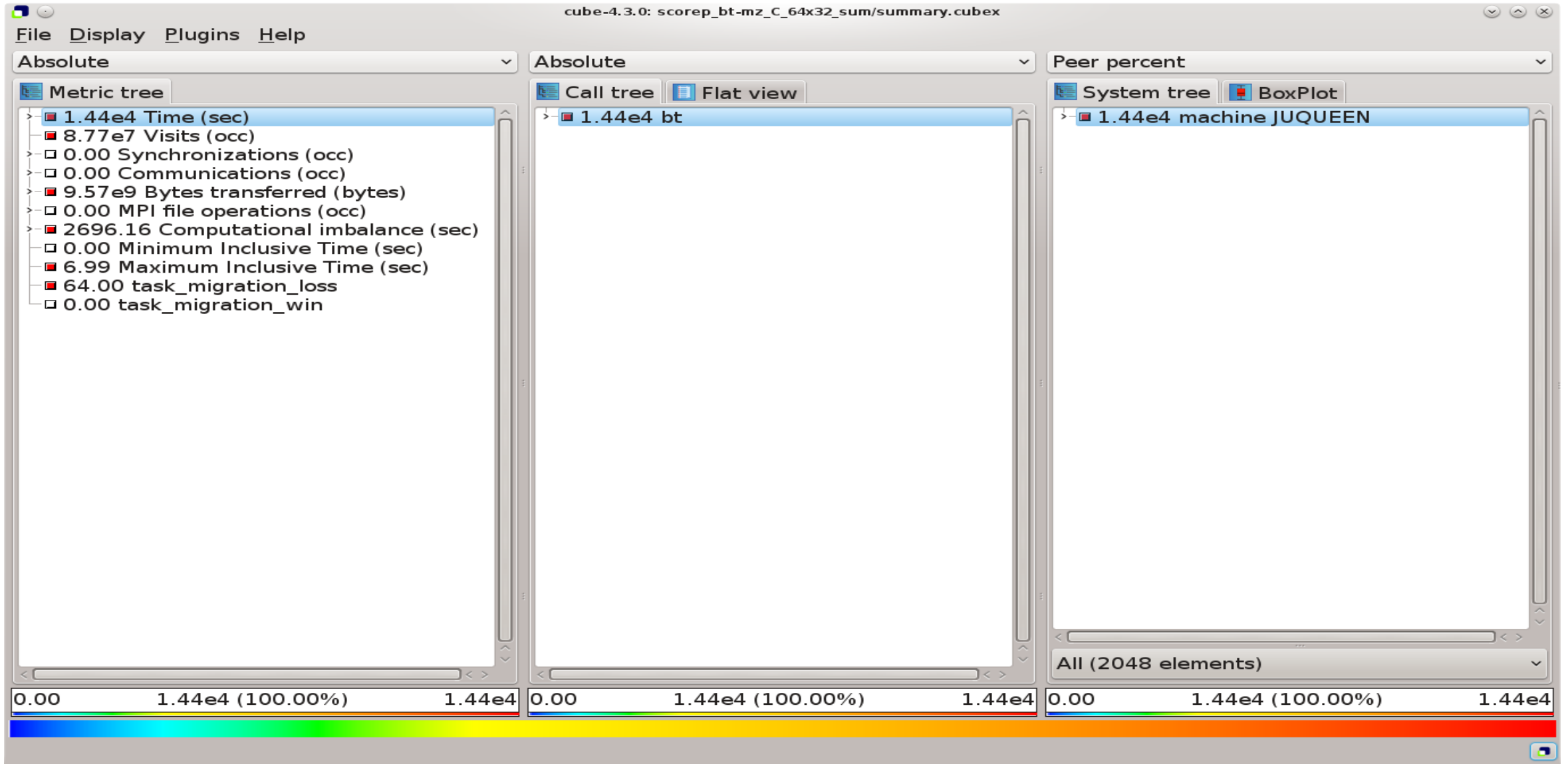
- These routines contribute around 35% of total time
  - However, much of that is most likely measurement overhead
    - Frequently executed
    - Time-per-visit ratio in the order of a few microseconds

☞ Avoid measurements to reduce the overhead

☞ List routines to be filtered in simple text file

**JÜLICH** Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# FILTERING: EXAMPLE

```
% cat filter.txt
SCOREP_REGION_NAMES_BEGIN
    EXCLUDE
        binvcrhs
        matmul_sub
        matvec_sub
        binvrhs
        lhsinit
        exact_solution
SCOREP_REGION_NAMES_END
```

- Score-P filtering files support

  - Wildcards (shell globs)

  - Blacklisting

  - Whitelisting

  - Filtering based on filenames

Mitglied der Helmholtz-Gemeinschaft

**JÜLICH** Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# CALL-PATH PROFILE: EXAMPLE (CONT.)

# CALL-PATH PROFILE: EXAMPLE (CONT.)

# CALL-PATH PROFILE: EXAMPLE (CONT.)



Split base metrics into more specific metrics

# WHY IS THE BOTTLENECK THERE?

- This is highly application dependent!

- Might require additional measurements

  - Hardware-counter analysis

    - CPU utilization

    - Cache behavior

  - Selective instrumentation

  - Automatic/manual event trace analysis

JÜLICH
Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# HARDWARE COUNTERS

- Counters: set of registers that count processor events, e.g. floating point operations or cycles

- Number of registers, counters and simultaneously measurable events vary between platforms

-  Can be measured by:

  - perf:

    - Integrated in Linux since Kernel 2.6.31

    - Library and CLI

  - LIKWID:

    - Direct access to MSRs (requires Kernel module)

    - Consists of multiple tools and an API

  - PAPI (Performance API)

JÜLICH
Forschungszentrum

JÜLICH
SUPERCOMPUTING
CENTRE

# PAPI

- Portable API: Uses the same routines to access counters across all supported architectures
- Used by most performance analysis tools

- High-level interface:
  - Predefined standard events, e.g. PAPI_FP_OPS
  - Availability and definition of events varies between platforms
  - List of available counters: `papi_avail (-d)`
- Low-level interface:
  - Provides access to all machine specific counters
  - Non-portable
  - More flexible
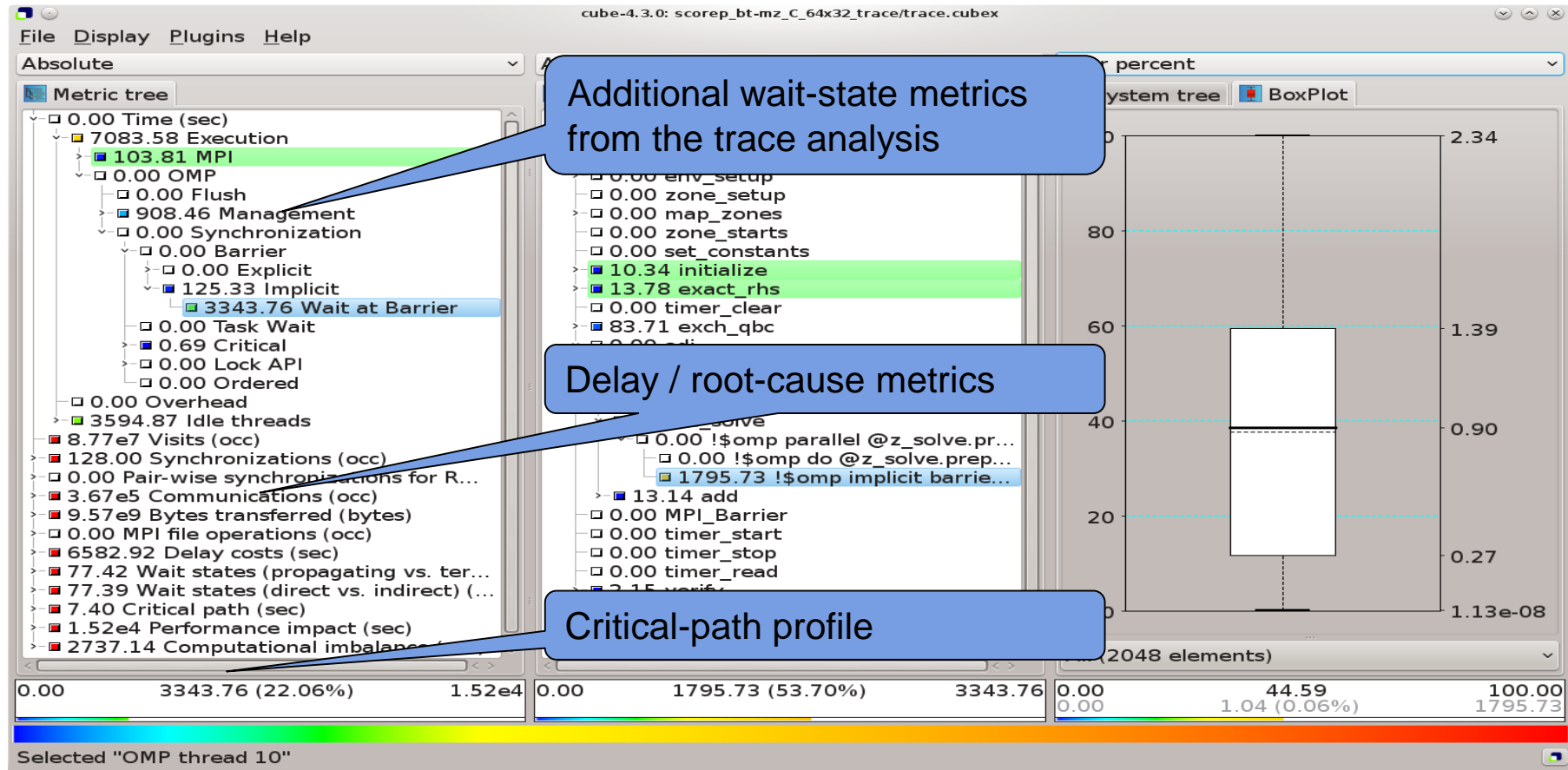  - List of available counters: `papi_native_avail`

Mitglied der Helmholtz-Gemeinschaft
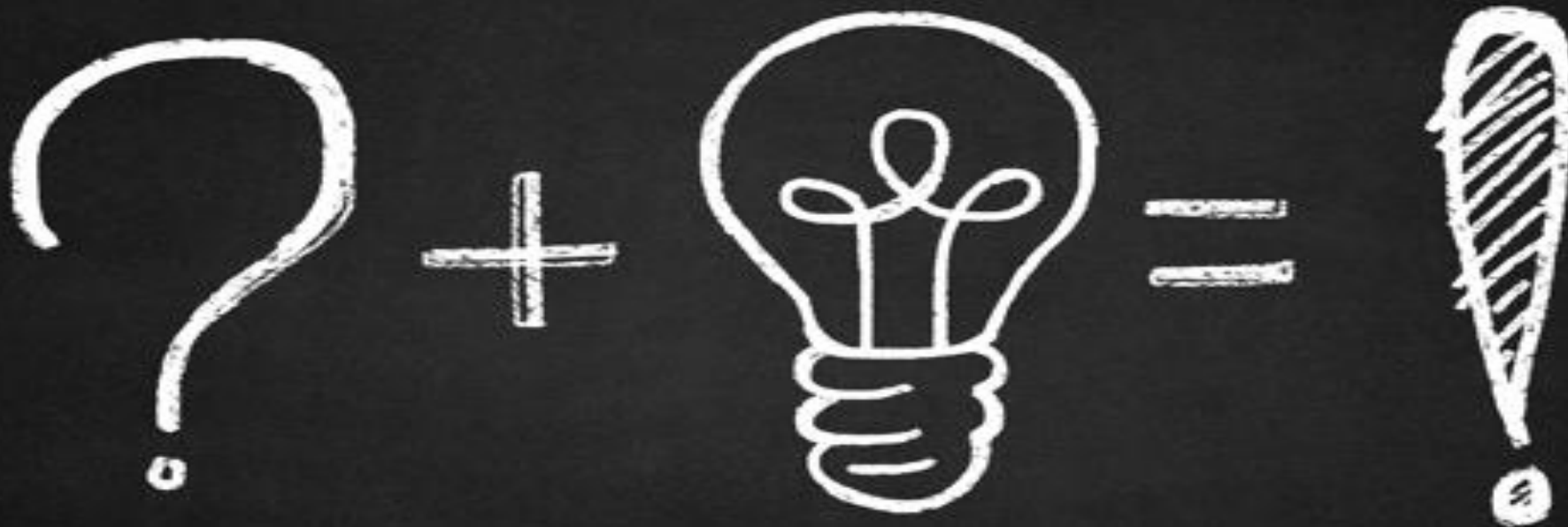
JÜLICH Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# TRACE GENERATION & ANALYSIS W/ SCALASCA

- Enable trace collection & analysis using "`-t`" option of "`scalasca -analyze`":

```
##########################
##  In the job script:  ##
##########################

module load ENV Score-P Scalasca
export SCOREP_TOTAL_MEMORY=120MB    # Consult score report
scalasca -analyze -f filter.txt -t \
    srun -n n [...] ./myprog
```

- **ATTENTION:**
  - Traces can quickly become extremely large!
  - Remember to use proper filtering, selective instrumentation, and Score-P memory specification
  - Before flooding the file system, **ask us for assistance!**

JÜLICH
Forschungszentrum | JÜLICH SUPERCOMPUTING CENTRE

# SCALASCA TRACE ANALYSIS EXAMPLE



Mitglied der Helmholtz-Gemeinschaft

# QUESTIONS

Mitglied der Helmholtz-Gemeinschaft