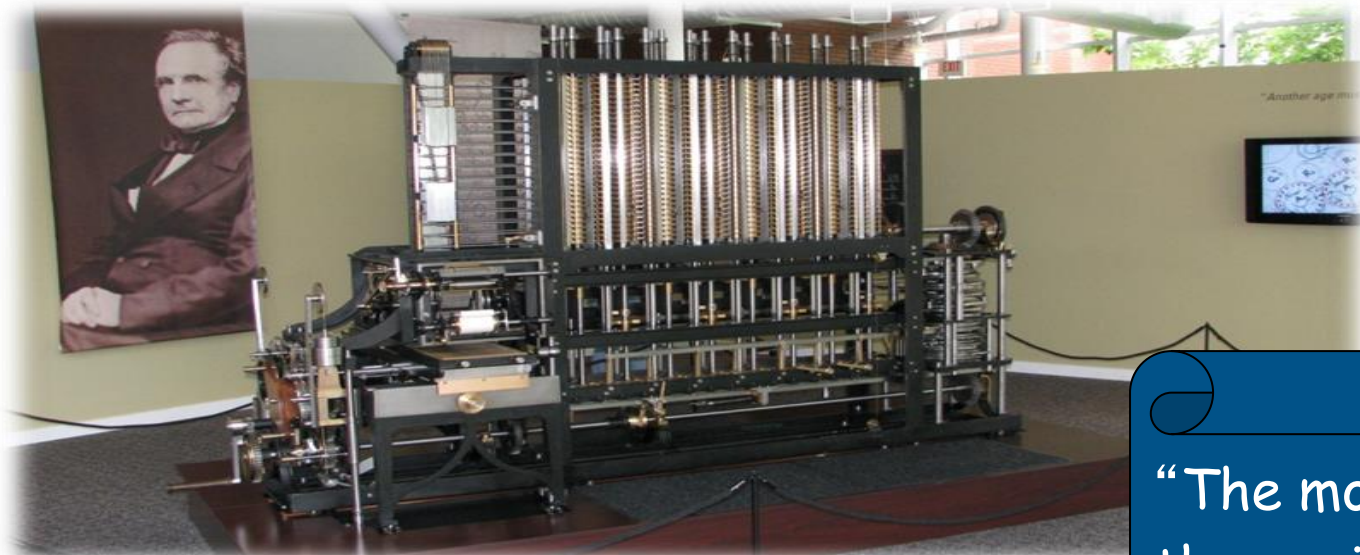


Introduction to Parallel Performance Engineering

VI-HPS Team
Christian Feld
Jülich Supercomputing Centre

(with content used with permission from tutorials by
Bernd Mohr/JSC, Brian Wylie/JSC, Markus Geimer/JSC,
Luiz DeRose/Cray, David Böhme/LLNL,
Andreas Knüpfer/TUD, Jens Doleschal/TUD)

Performance: an old problem



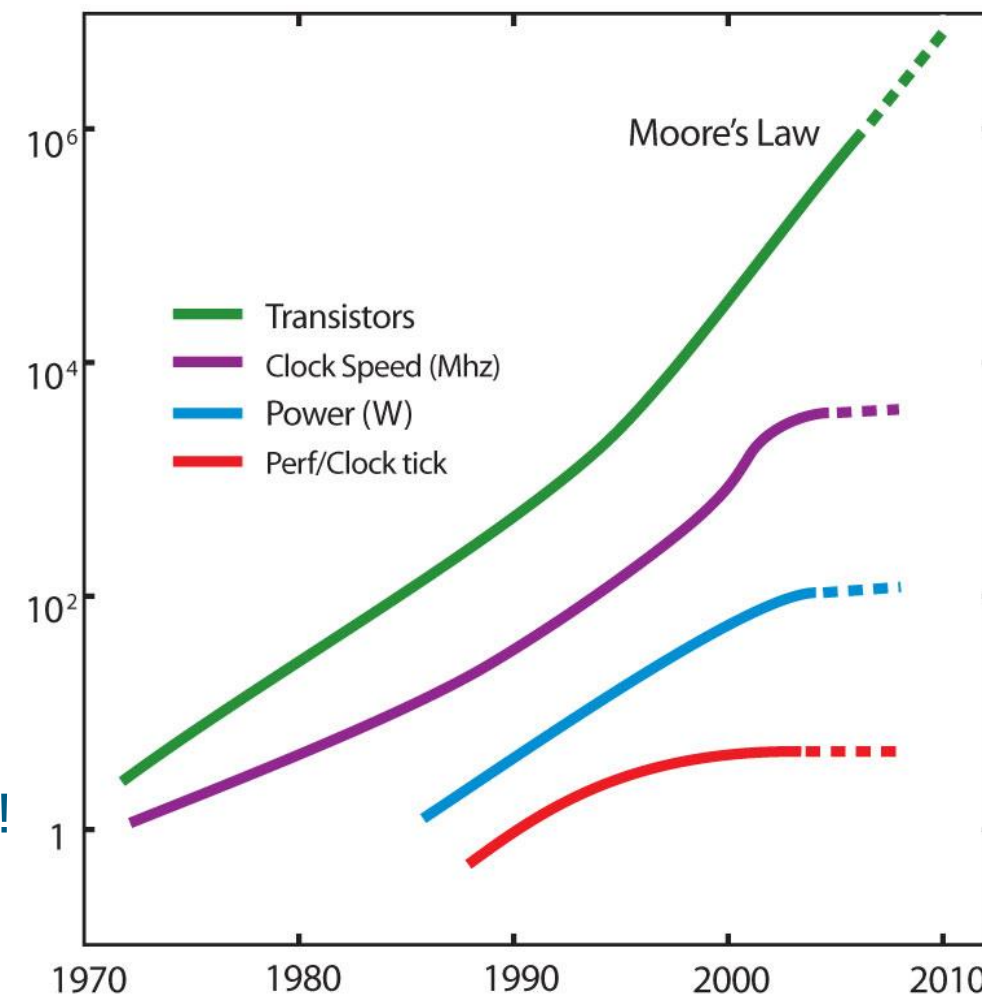
Difference Engine

“The most constant difficulty in contriving the engine has arisen from the desire to reduce the time in which the calculations were executed to the shortest which is possible.”

Charles Babbage
1791 – 1871

Today: the “free lunch” is over

- Moore's law is still in charge, but
 - Clock rates no longer increase
 - Performance gains only through increased parallelism
 - Optimizations of applications more difficult
 - Increasing application complexity
 - Multi-physics
 - Multi-scale
 - Increasing machine complexity
 - Hierarchical networks / memory
 - More CPUs / multi-core
- 👉 Every doubling of scale reveals a new bottleneck!



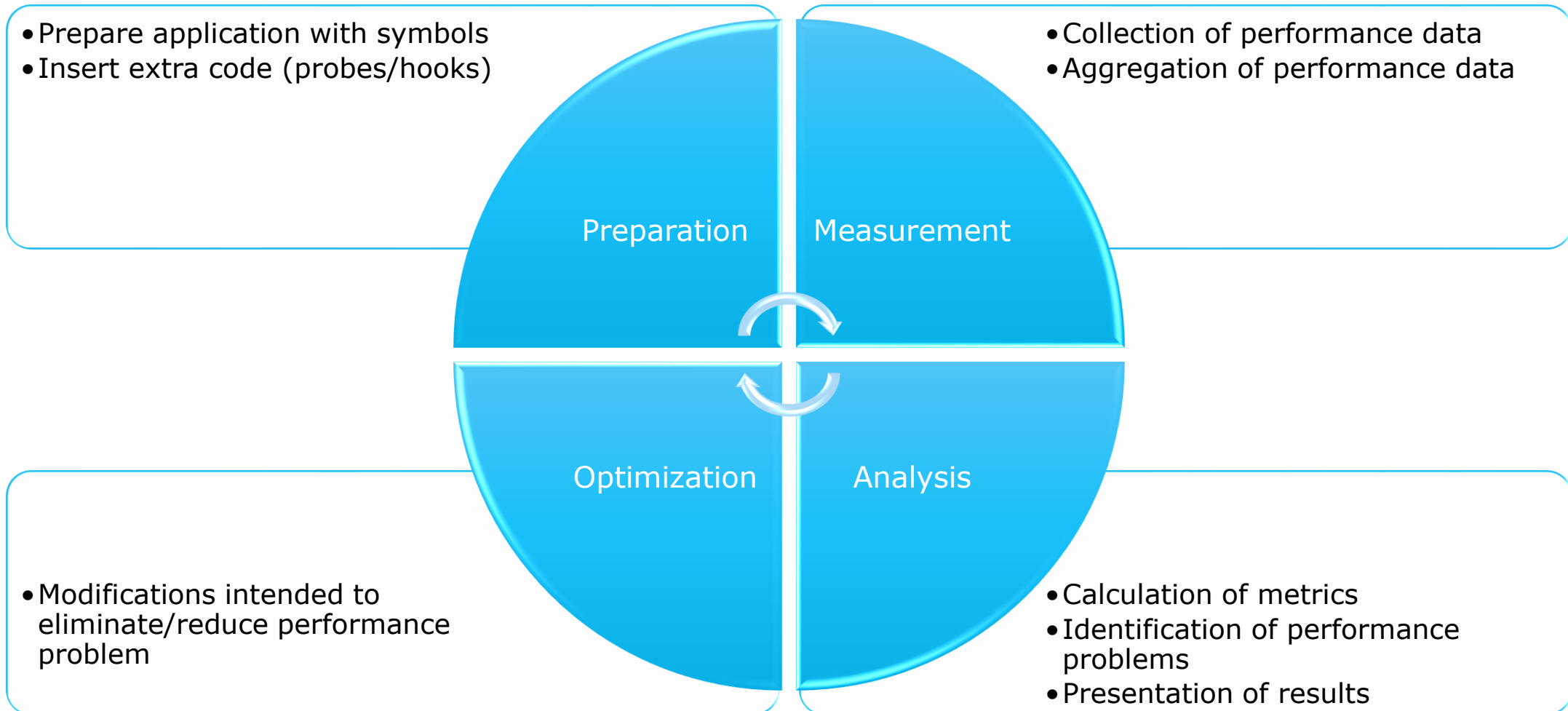
Performance factors of parallel applications

- **“Sequential”** performance factors
 - Computation
 - ☞ Choose right algorithm, use optimizing compiler
 - Cache and memory
 - ☞ Tough! Only limited tool support, hope compiler gets it right
 - Input / output
 - ☞ Often not given enough attention
- **“Parallel”** performance factors
 - Partitioning / decomposition
 - Communication (i.e., message passing)
 - Multithreading
 - Synchronization / locking
 - ☞ More or less understood, good tool support

Tuning basics

- Successful engineering is a combination of
 - Careful setting of various tuning parameters
 - The right algorithms and libraries
 - Compiler flags and directives
 - ...
 - Thinking !!!
- Measurement is better than guessing
 - To determine performance bottlenecks
 - To compare alternatives
 - To validate tuning decisions and optimizations
 - ☞ After each step!

Performance engineering workflow



The 80/20 rule

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

“If you optimize everything,
you will always be unhappy.”

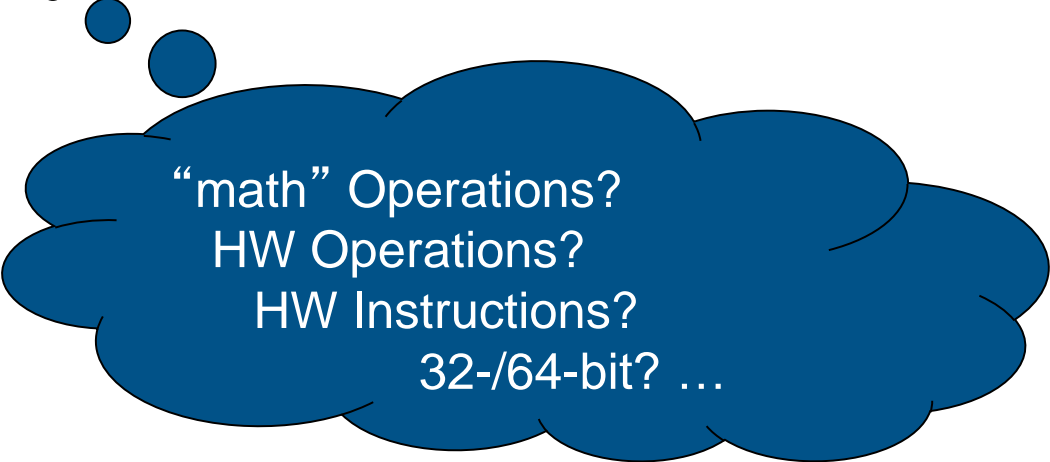
Donald E. Knuth

Metrics of performance

- What can be measured?
 - A **count** of how often an event occurs
 - E.g., the number of MPI point-to-point messages sent
 - The **duration** of some interval
 - E.g., the time spent these send calls
 - The **size** of some parameter
 - E.g., the number of bytes transmitted by these calls
- Derived metrics
 - E.g., rates / throughput
 - Needed for normalization

Example metrics

- Execution time
- Number of function calls
- CPI
 - CPU cycles per instruction
- FLOPS
 - Floating-point operations executed per second.



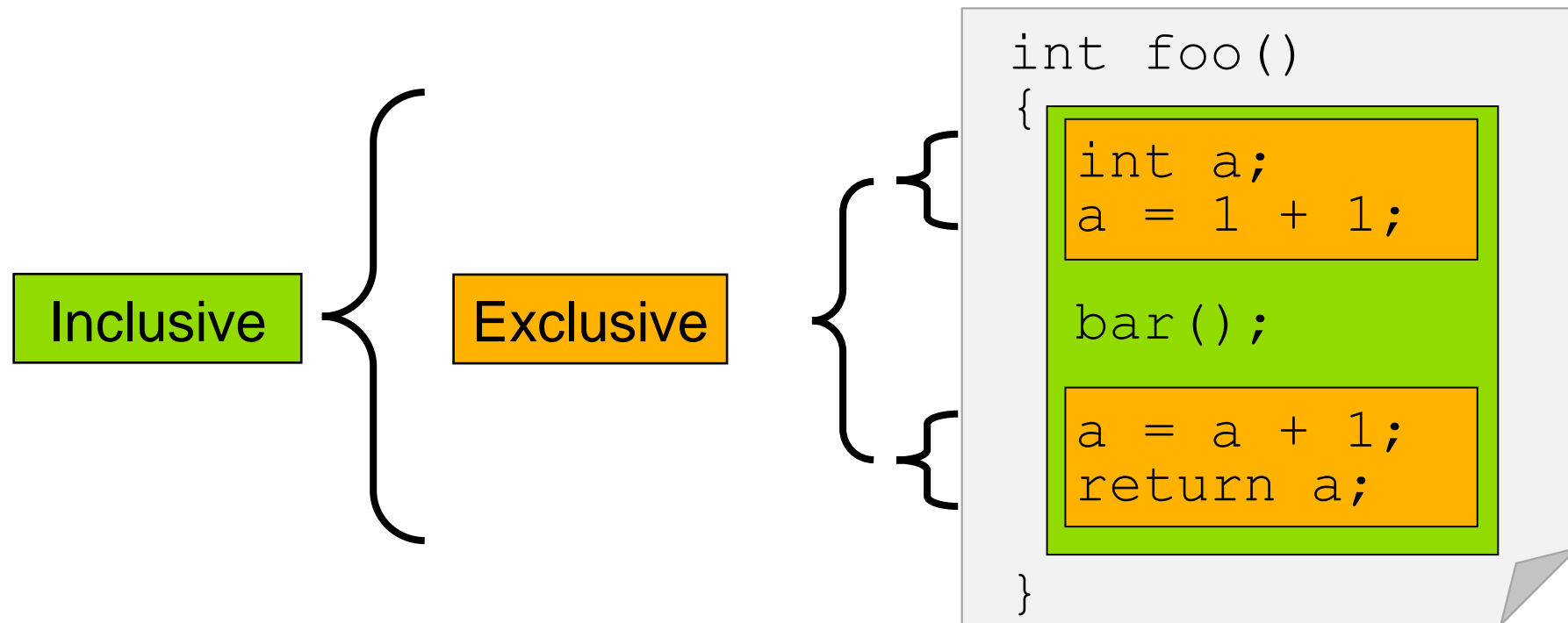
“math” Operations?
HW Operations?
HW Instructions?
32-/64-bit? ...

Execution time

- Wall-clock time
 - Includes waiting time: I/O, memory, other system activities
 - In time-sharing environments also the time consumed by other applications
- CPU time
 - Time spent by the CPU to execute the application
 - Does not include time the program was context-switched out
 - Problem: Does not include inherent waiting time (e.g., I/O)
 - Problem: Portability? What is user, what is system time?
- Problem: Execution time is non-deterministic
 - Use mean or minimum of several runs

Inclusive vs. Exclusive values

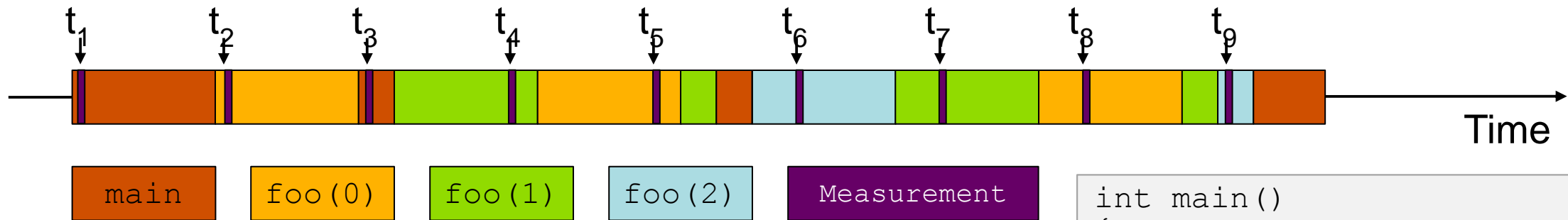
- Inclusive
 - Information of all sub-elements aggregated into single value
- Exclusive
 - Information cannot be subdivided further



Classification of measurement techniques

- **How are performance measurements triggered?**
 - **Sampling**
 - **Code instrumentation**
- How is performance data recorded?
 - Profiling / Runtime summarization
 - Tracing
- How is performance data analyzed?
 - Online
 - Post mortem

Sampling



- Running program is periodically interrupted to take measurement
 - Timer interrupt, OS signal, or HWC overflow
 - Service routine examines return-address stack
 - Addresses are mapped to routines using symbol table information
- Statistical inference of program behavior
 - Not very detailed information on highly volatile metrics
 - Requires long-running applications
- Works with unmodified executables

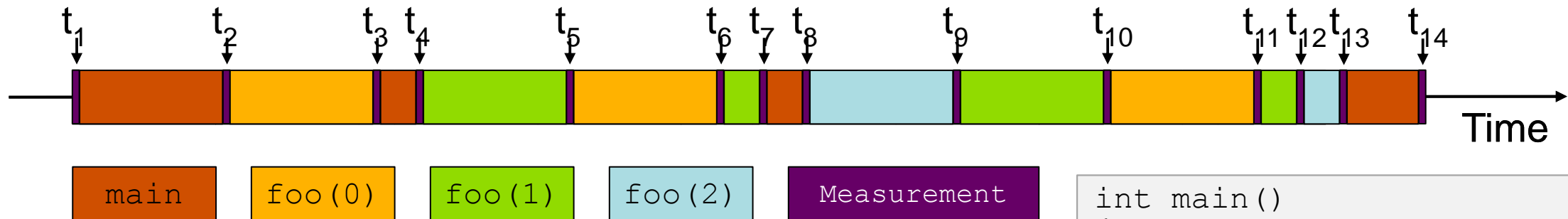
```
int main()
{
    int i;

    for (i=0; i < 3; i++)
        foo(i);

    return 0;
}

void foo(int i)
{
    if (i > 0)
        foo(i - 1);
}
```


Instrumentation



- Measurement code is inserted such that every event of interest is captured directly
 - Can be done in various ways
- Advantage:
 - Much more detailed information
- Disadvantage:
 - Processing of source-code / executable necessary
 - Large relative overheads for small functions

```

int main()
{
    int i;
    Enter("main");
    for (i=0; i < 3; i++)
        foo(i);
    Leave("main");
    return 0;
}

void foo(int i)
{
    Enter("foo");
    if (i > 0)
        foo(i - 1);
    Leave("foo");
}

```

Instrumentation techniques

- **Static** instrumentation
 - Program is instrumented prior to execution
- **Dynamic** instrumentation
 - Program is instrumented at runtime
- Code is inserted
 - Manually
 - Automatically
 - By a preprocessor / source-to-source translation tool
 - By a compiler
 - By linking against a pre-instrumented library / runtime system
 - By binary-rewrite / dynamic instrumentation tool

Critical issues

- Accuracy
 - Intrusion overhead
 - Measurement itself needs time and thus lowers performance
 - Perturbation
 - Measurement alters program behaviour
 - E.g., memory access pattern
 - Accuracy of timers & counters
- Granularity
 - How many measurements?
 - How much information / processing during each measurement?

☞ *Tradeoff: Accuracy vs. Expressiveness of data*

Classification of measurement techniques

- How are performance measurements triggered?
 - Sampling
 - Code instrumentation
- **How is performance data recorded?**
 - **Profiling / Runtime summarization**
 - **Tracing**
- How is performance data analyzed?
 - Online
 - Post mortem

Profiling / Runtime summarization

- Recording of aggregated information
 - Total, maximum, minimum, ...
- For measurements
 - Time
 - Counts
 - Function calls
 - Bytes transferred
 - Hardware counters
- Over program and system entities
 - Functions, call sites, basic blocks, loops, ...
 - Processes, threads

☞ *Profile = summarization of events over execution interval*

Types of profiles

- Flat profile
 - Shows distribution of metrics per routine / instrumented region
 - Calling context is not taken into account
- Call-path profile
 - Shows distribution of metrics per executed call path
 - Sometimes only distinguished by partial calling context (e.g., two levels)
- Special-purpose profiles
 - Focus on specific aspects, e.g., MPI calls or OpenMP constructs
 - Comparing processes/threads

Tracing

- Recording detailed information about significant points (events) during execution of the program
 - Enter / leave of a region (function, loop, ...)
 - Send / receive a message, ...
- Save information in event record
 - Timestamp, location, event type
 - Plus event-specific information (e.g., communicator, sender / receiver, ...)
- Abstract execution model on level of defined events

☞ *Event trace = Chronologically ordered sequence of event records*

Event tracing

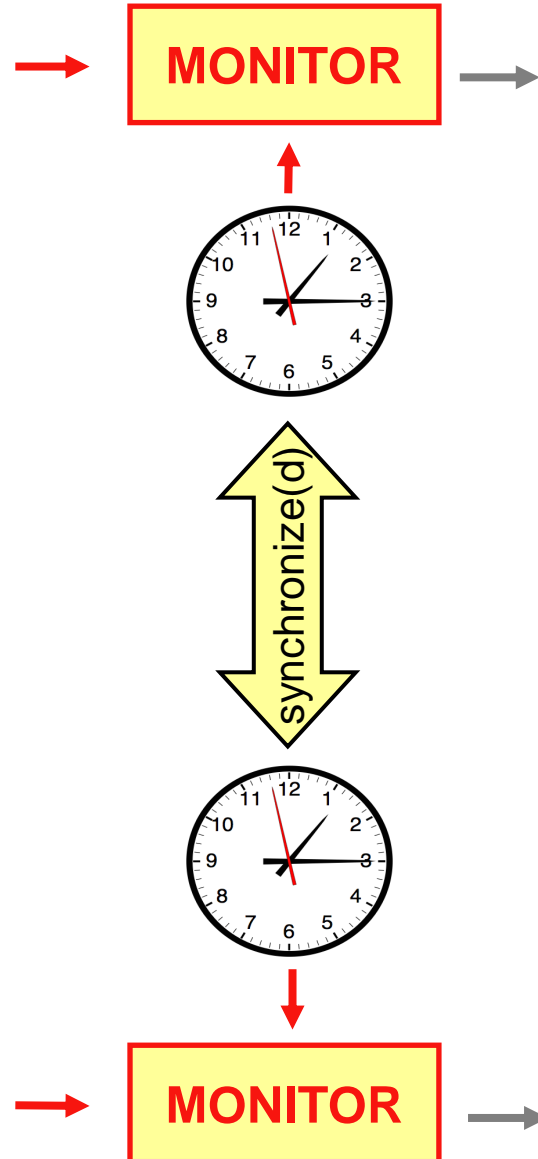
Process A

```
void foo() {
  trc_enter("foo");
  ...
  trc_send(B);
  send(B, tag, buf);
  ...
  trc_exit("foo");
}
```

instrument

Process B

```
void bar() {
  trc_enter("bar");
  ...
  recv(A, tag, buf);
  trc_recv(A);
  ...
  trc_exit("bar");
}
```



Local trace A

...	
58	ENTER foo
62	SEND to B
64	EXIT foo
...	

Local trace B

...	
60	ENTER bar
68	RECV from A
69	EXIT bar
...	

Global trace view

...		
58	A	ENTER foo
60	B	ENTER bar
62	A	SEND to B
64	A	EXIT foo
68	B	RECV from A
69	B	EXIT bar
...		

(Virtual merge)

Tracing Pros & Cons

- Tracing advantages
 - Event traces preserve the **temporal** and **spatial** relationships among individual events (👉 context)
 - Allows reconstruction of **dynamic** application behaviour on any required level of abstraction
 - Most general measurement technique
 - Profile data can be reconstructed from event traces
- Disadvantages
 - Traces can very quickly become extremely large
 - Writing events to file at runtime may causes perturbation

Classification of measurement techniques

- How are performance measurements triggered?
 - Sampling
 - Code instrumentation
- How is performance data recorded?
 - Profiling / Runtime summarization
 - Tracing
- **How is performance data analyzed?**
 - **Online**
 - **Post mortem**

Online analysis

- Performance data is processed during measurement run
 - Process-local profile aggregation
 - Requires formalized knowledge about performance bottlenecks
 - More sophisticated inter-process analysis using
 - “Piggyback” messages
 - Hierarchical network of analysis agents
- Online analysis often involves application steering to interrupt and re-configure the measurement

Post-mortem analysis

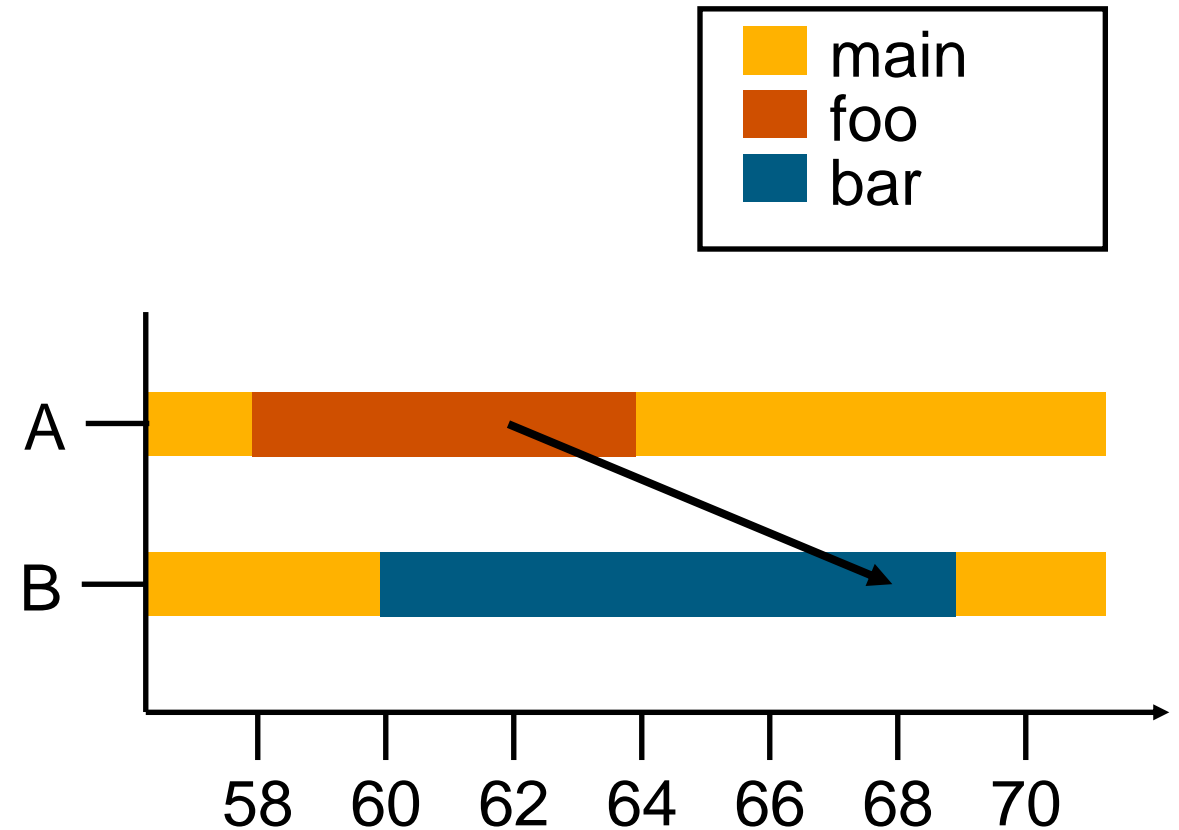
- Performance data is stored at end of measurement run
- Data analysis is performed afterwards
 - Automatic search for bottlenecks
 - Visual trace analysis
 - Calculation of statistics

Example: Time-line visualization

Global trace view

...		
58	A	ENTER foo
60	B	ENTER bar
62	A	SEND to B
64	A	EXIT foo
68	B	RECV from A
69	B	EXIT bar
...		

Post-Mortem
Analysis

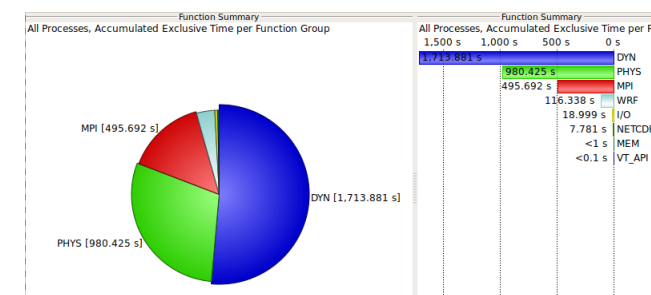
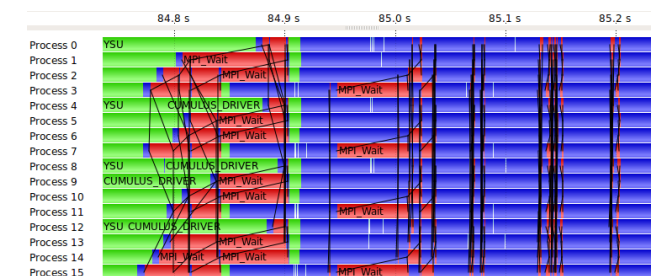


Trace visualizers

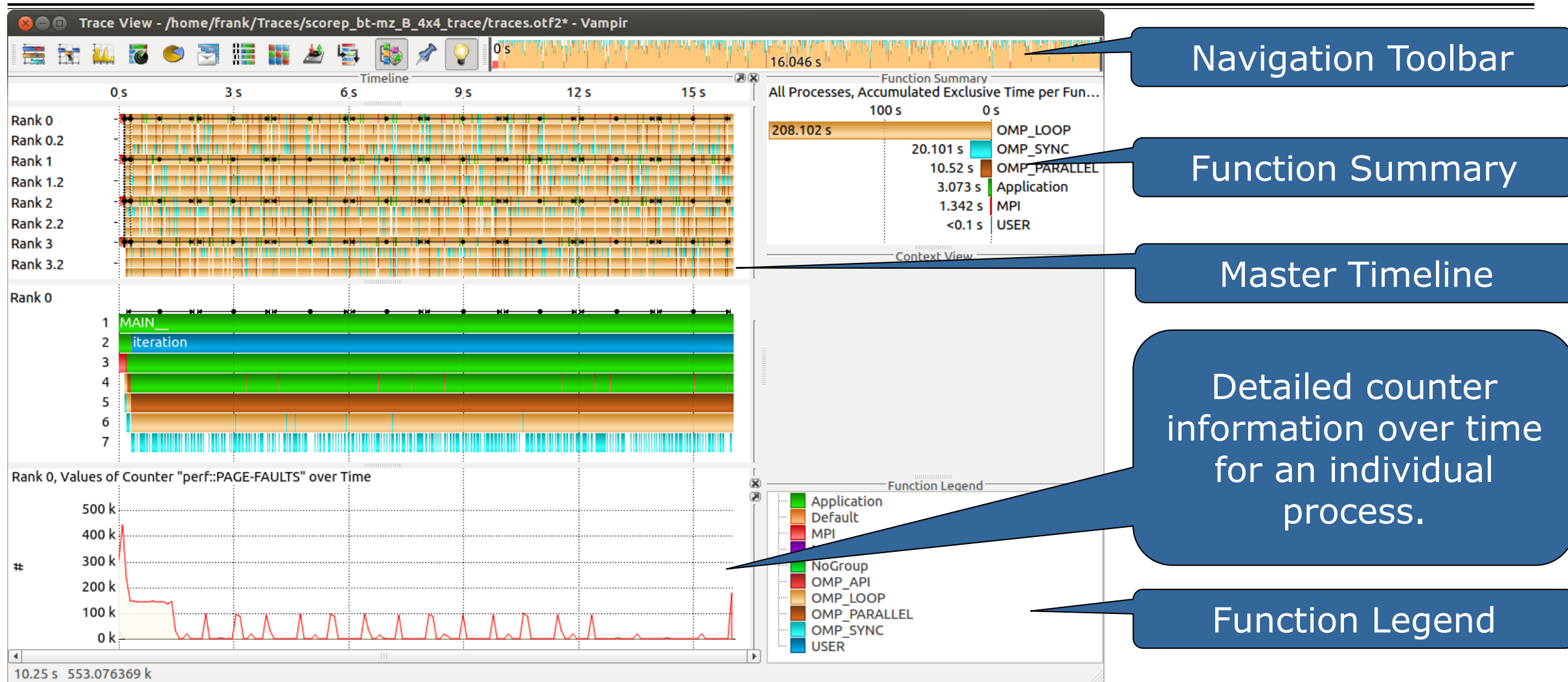
- Jumpshot (ANL) Process-local profile aggregation
 - Free, basic MPI visualizer (routines, messages)
 - SLOG-2 format
 - MPE tracing + converters from TAU, (EPILOG)
- Paraver (BSC)
 - Free, extremely flexible and programmable visualizer
 - PRV format
 - Extrae tracing + converters from TAU, EPILOG, (OTF)
- Vampir (TUD)
 - Commercial portable trace visualizer
 - OTF2, OTF, EPILOG format
- Intel trace collector and analyzer
 - Commercial, Intel-only trace collection and visualizer

Event Trace Visualization with Vampir

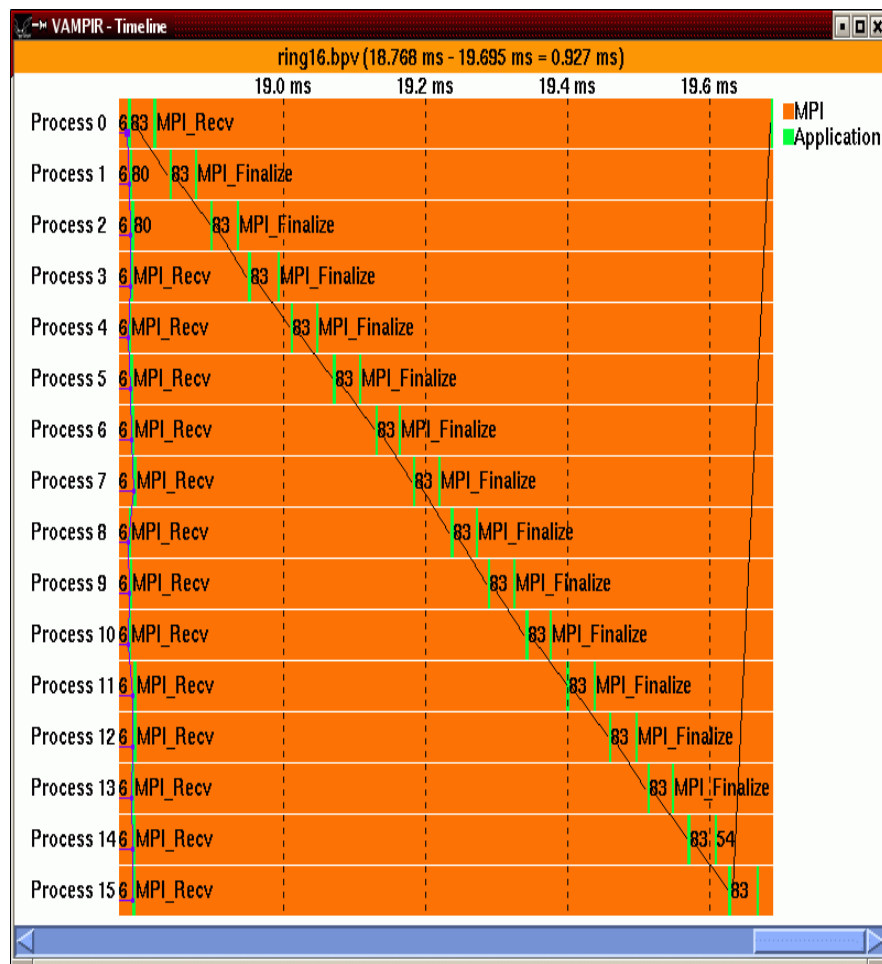
- Alternative and supplement to automatic analysis
 - Show dynamic run-time behavior graphically at any level of detail
 - Provide statistics and performance metrics
- **Timeline charts**
 - Show application activities and communication along a time axis
 - **Summary charts**
 - Provide quantitative results for the currently selected time interval



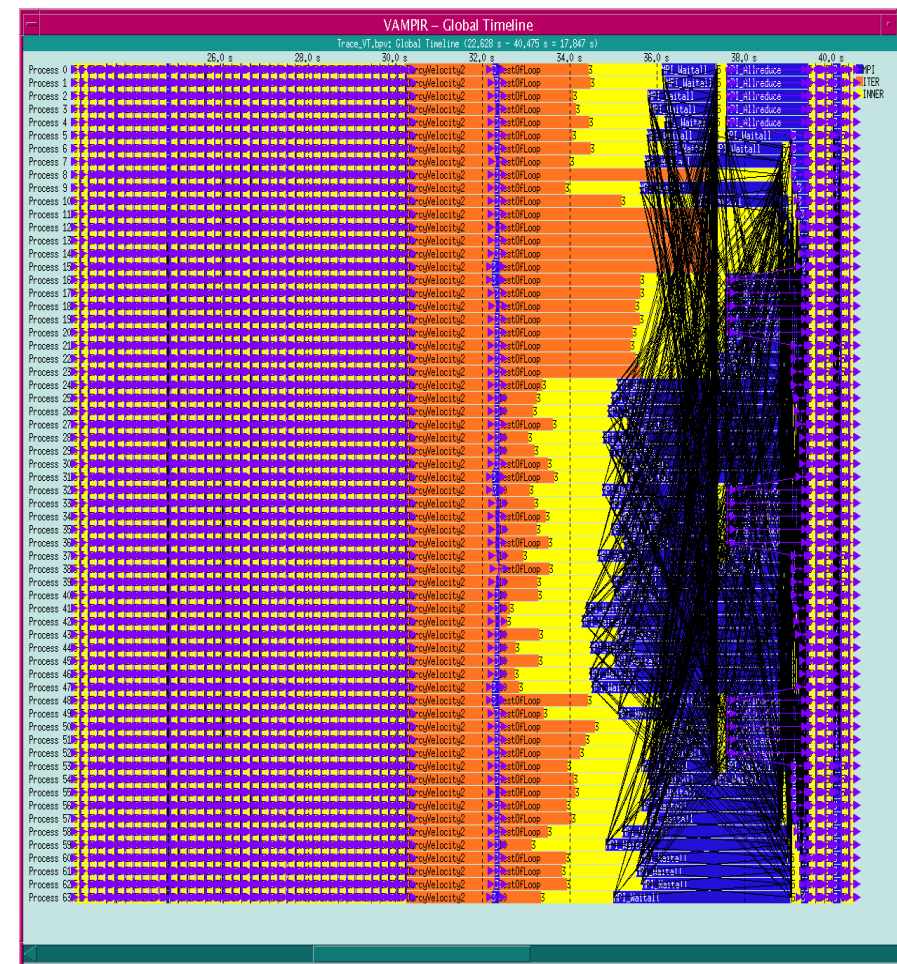
Visualization of the NPB-MZ-MPI / BT trace



“A picture is worth a 1000 words ...”



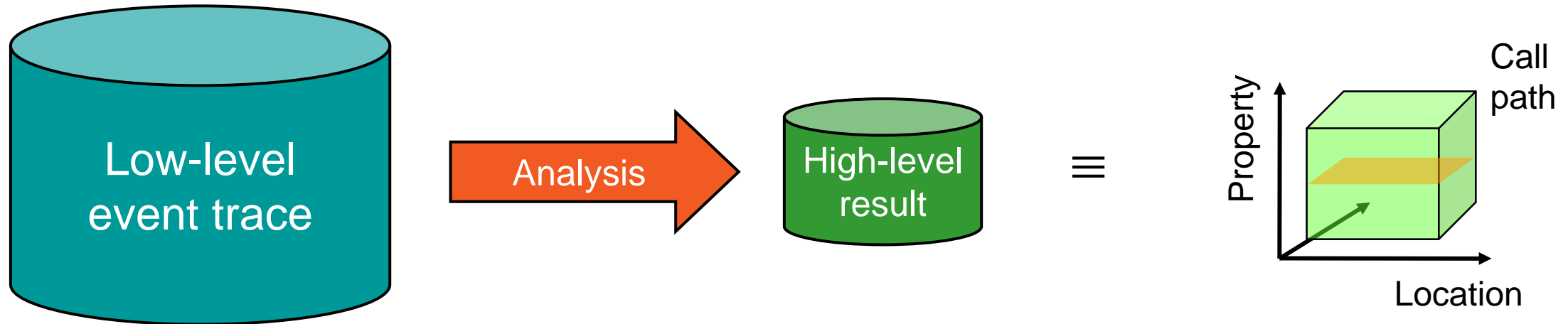
MPI ring example



“Real world” example

Automatic trace analysis

- Idea
 - Automatic search for patterns of inefficient behavior
 - Classification of behavior & quantification of significance



- Guaranteed to cover the entire event trace
- Quicker than manual/visual trace analysis
- Parallel replay analysis exploits available memory & processors to deliver scalability

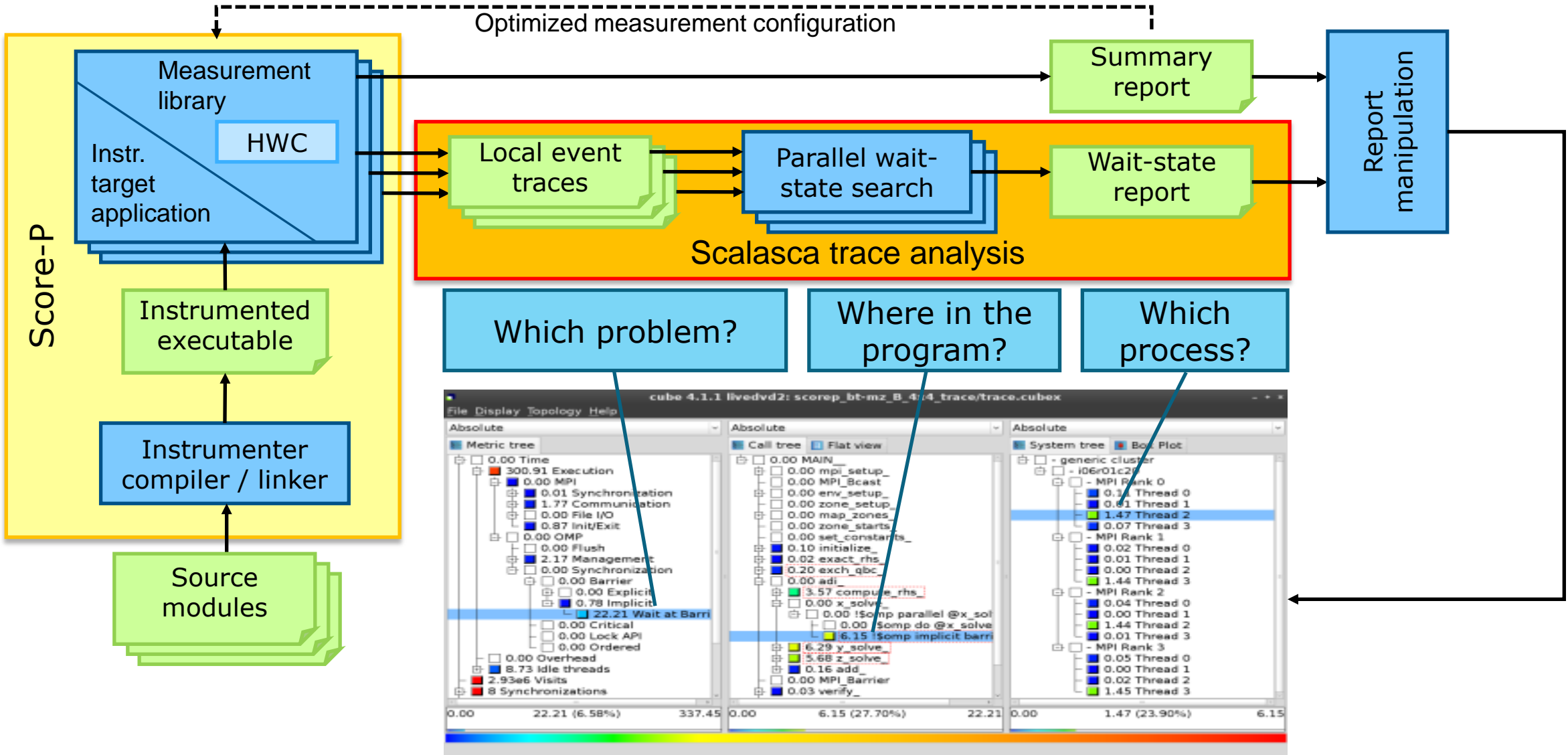
The Scalasca project: Objective

- Development of a **scalable** performance analysis toolset for most popular parallel programming paradigms
- Specifically targeting **large-scale** parallel applications
 - Such as those running on IBM Blue Gene or Cray systems with one million or more processes/threads
- Latest release:
 - Scalasca v2.3.1 coordinated with Score-P v2.0.2 (May 2016)

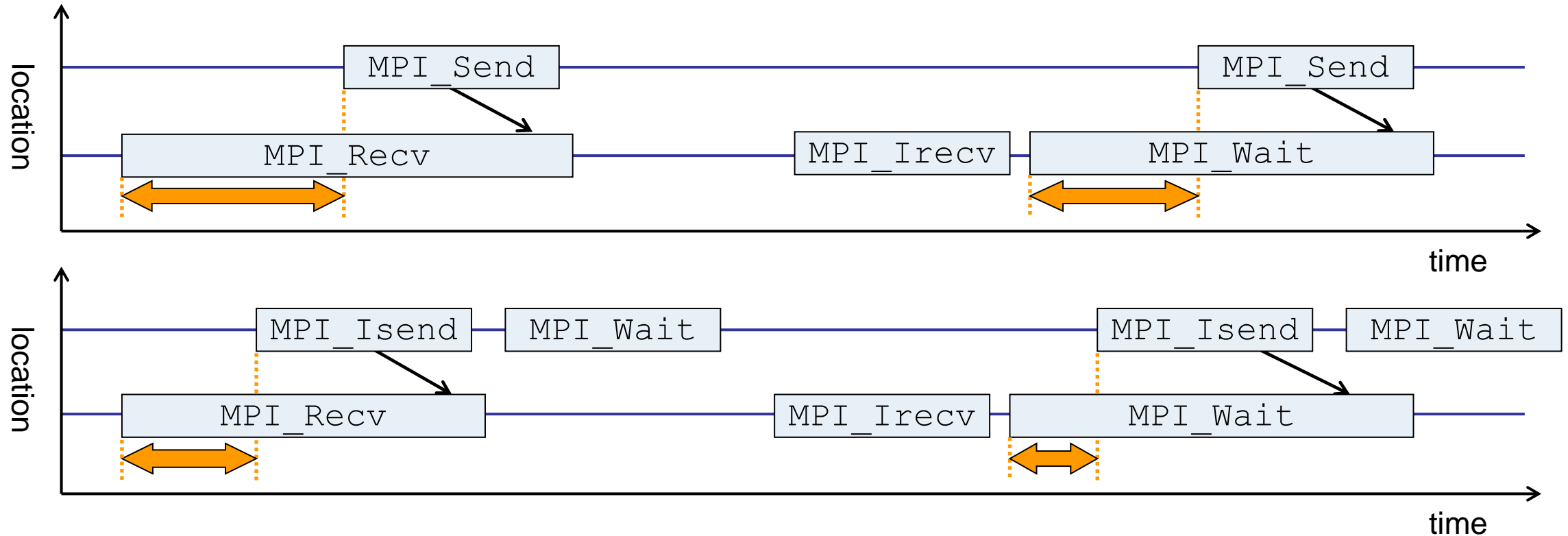
Scalasca features

- Open source, 3-clause BSD license
- Fairly portable
 - IBM Blue Gene, Cray XT/XE/XK/XC, SGI Altix, Fujitsu FX10/100 & K computer, Linux clusters (x86, Power, ARM), Intel Xeon Phi, ...
- Uses Score-P instrumenter & measurement libraries
 - Scalasca 2 core package focuses on trace-based analyses
 - Supports common data formats
 - Reads event traces in OTF2 format
 - Writes analysis reports in CUBE4 format
- Current limitations:
 - Unable to handle traces containing CUDA or SHMEM events, or OpenMP nested parallelism
 - PAPI/rusage metrics for trace events are ignored

Scalasca workflow

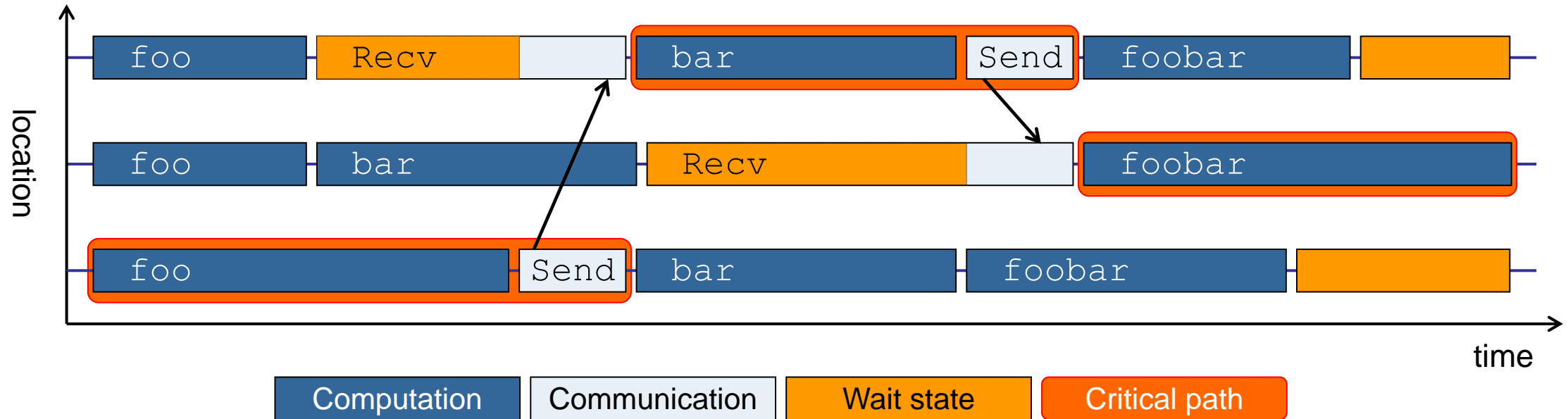


Example: “Late Sender” wait state



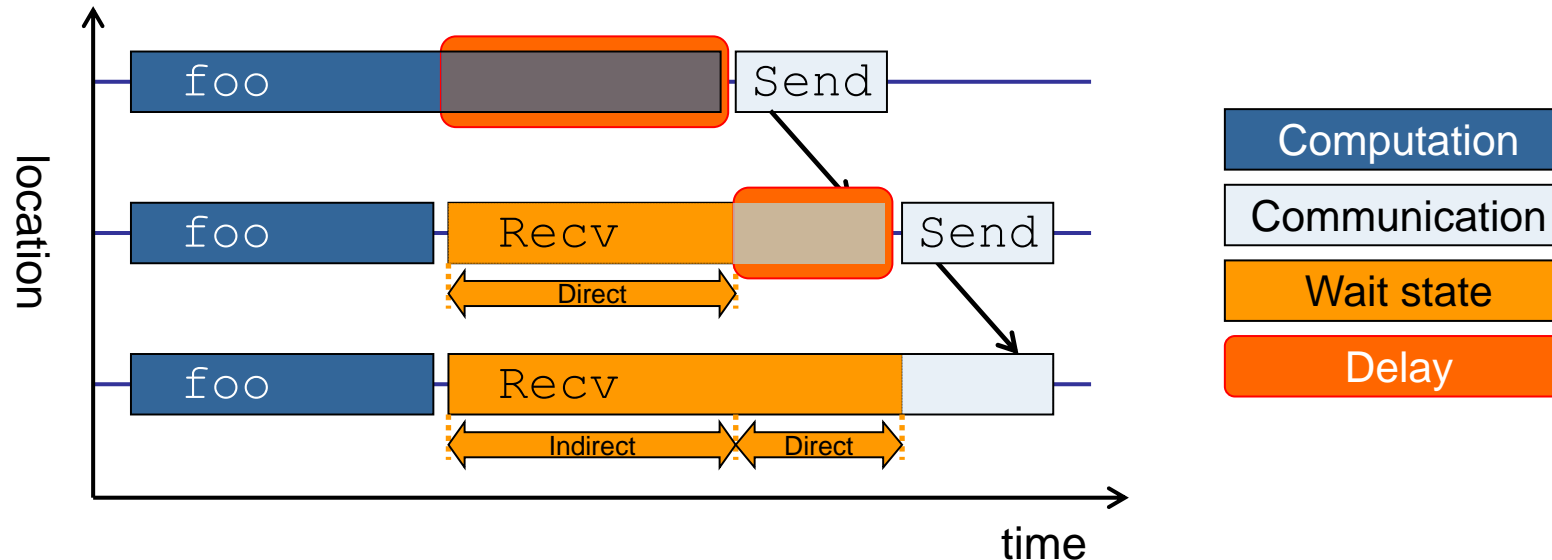
- Waiting time caused by a blocking receive operation posted earlier than the corresponding send
- Applies to blocking as well as non-blocking communication

Example: Critical path



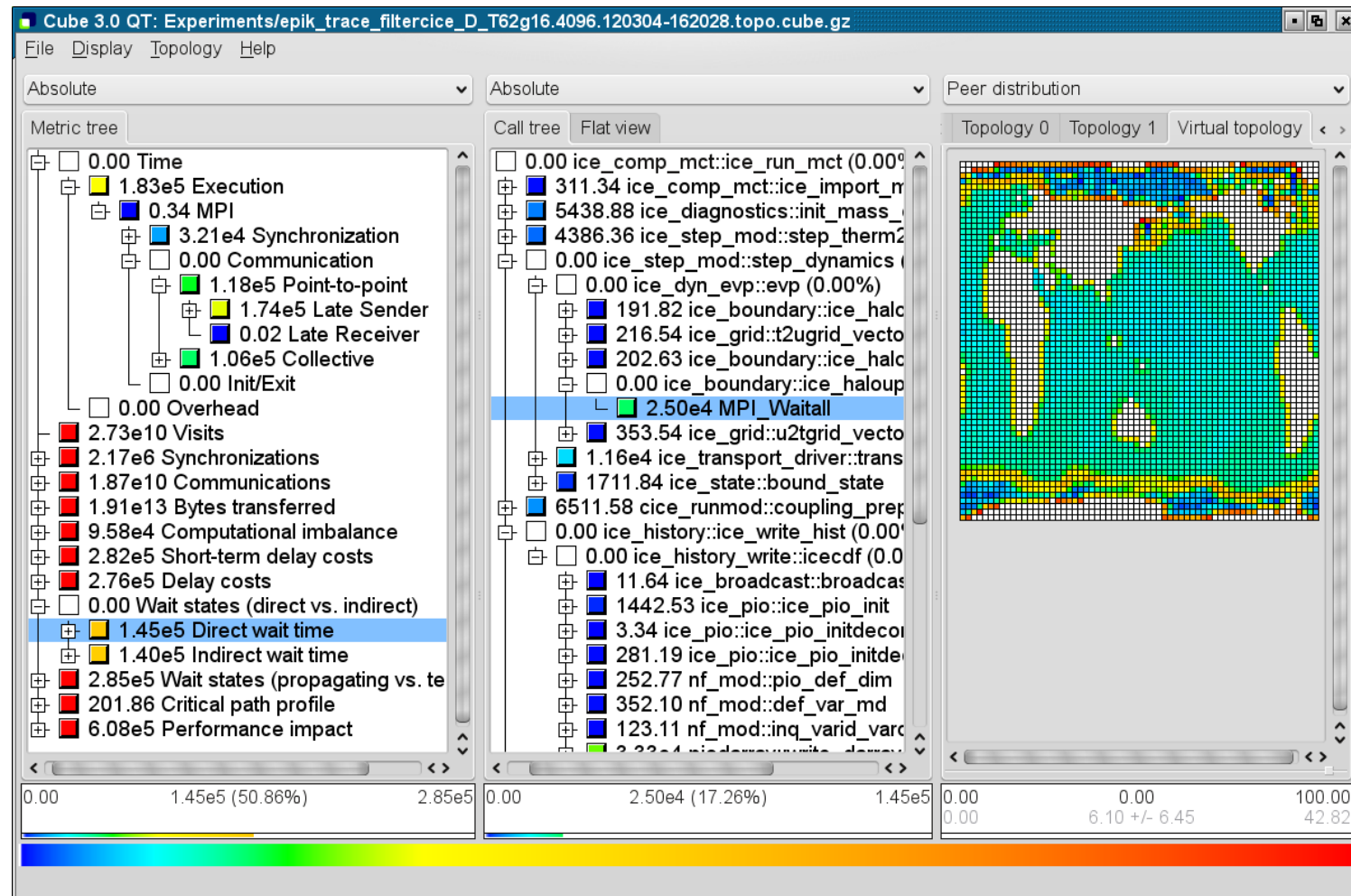
- Shows call paths and processes/threads that are responsible for the program's wall-clock runtime
- Identifies good optimization candidates and parallelization bottlenecks

Example: Root-cause analysis

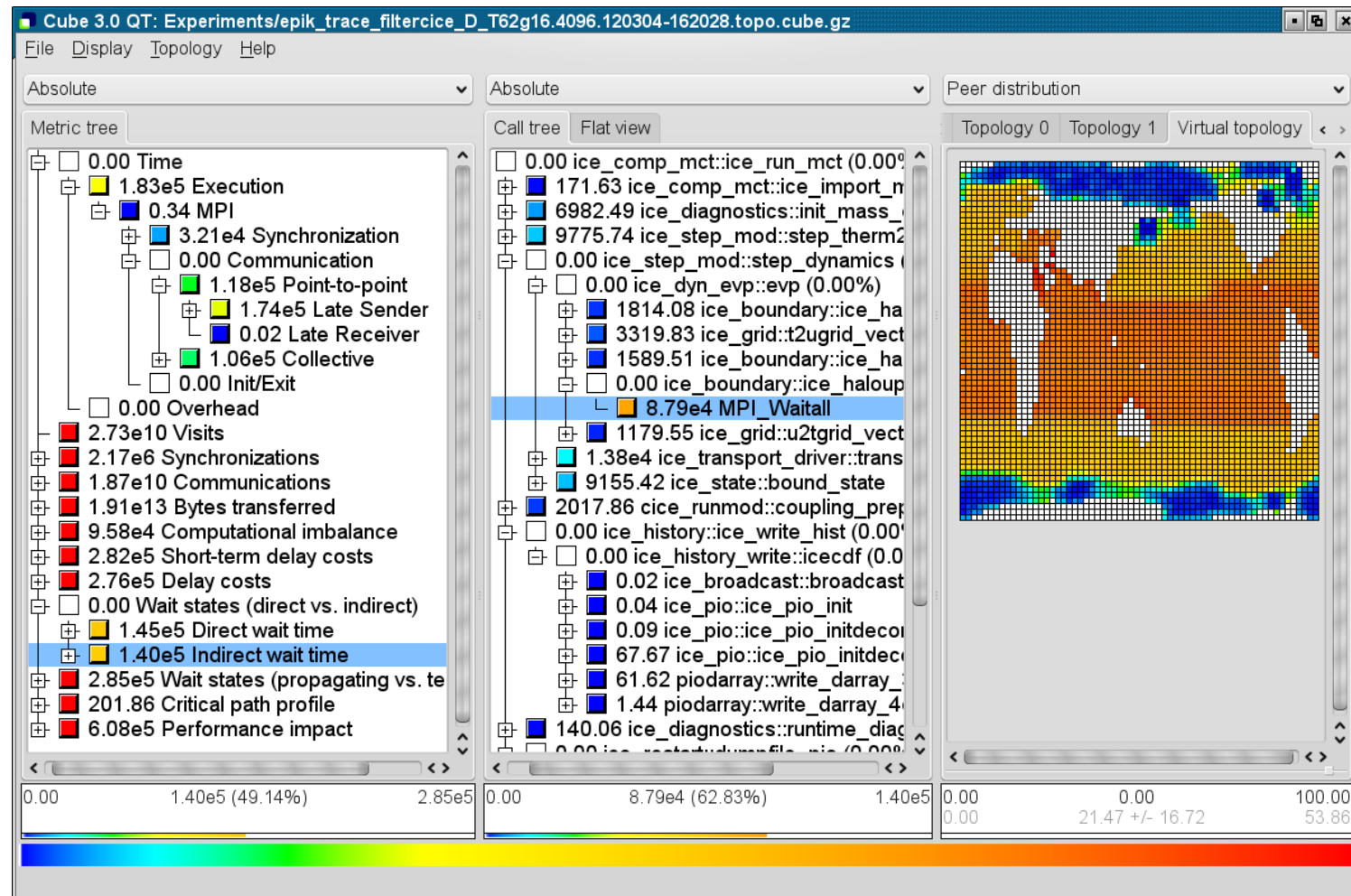


- Classifies wait states into direct and indirect (i.e., caused by other wait states)
- Identifies *delays* (excess computation/communication) as root causes of wait states
- Attributes wait states as *delay costs*

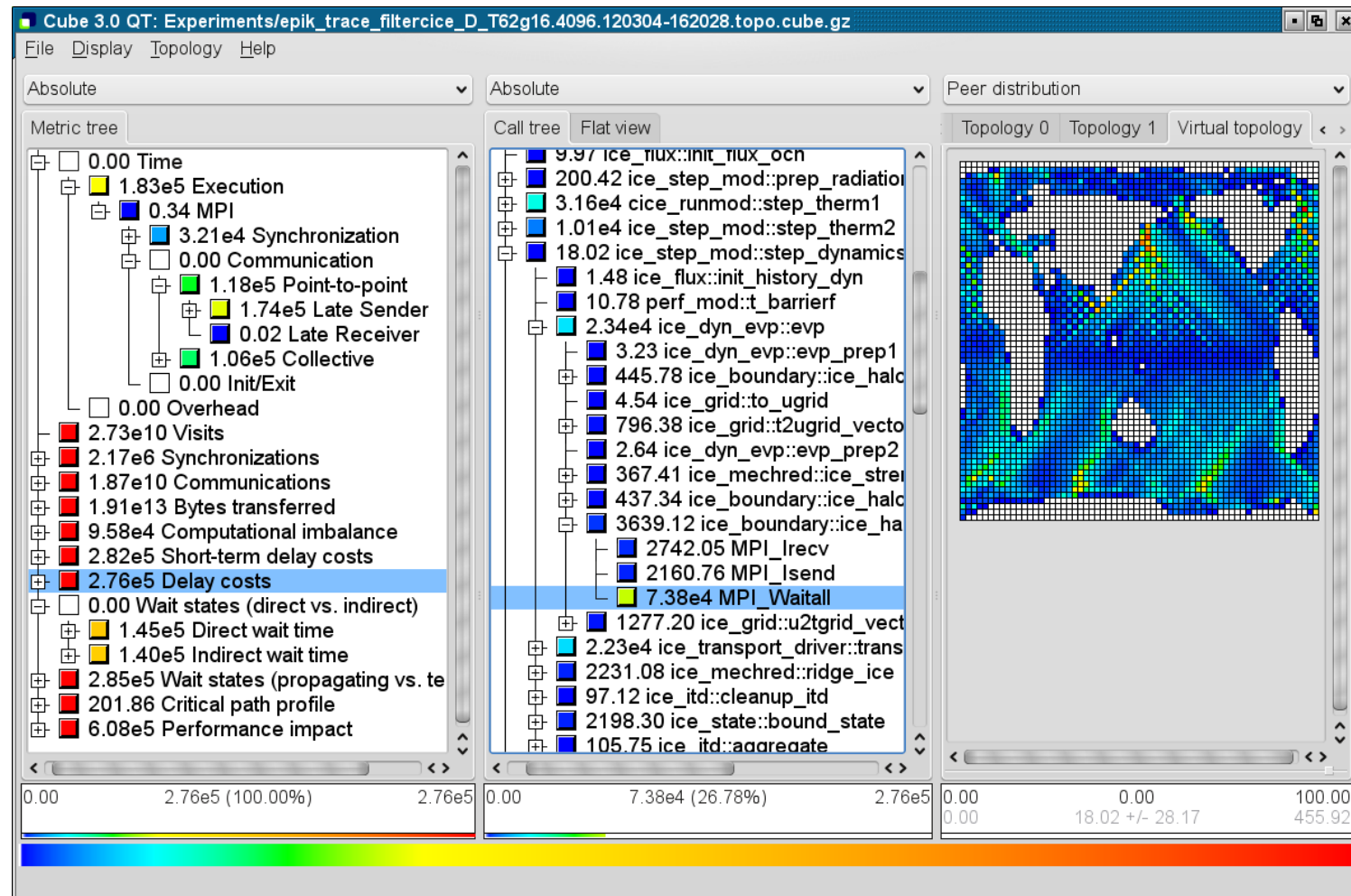
Example: Root-cause analysis - CESM Sea Ice Module



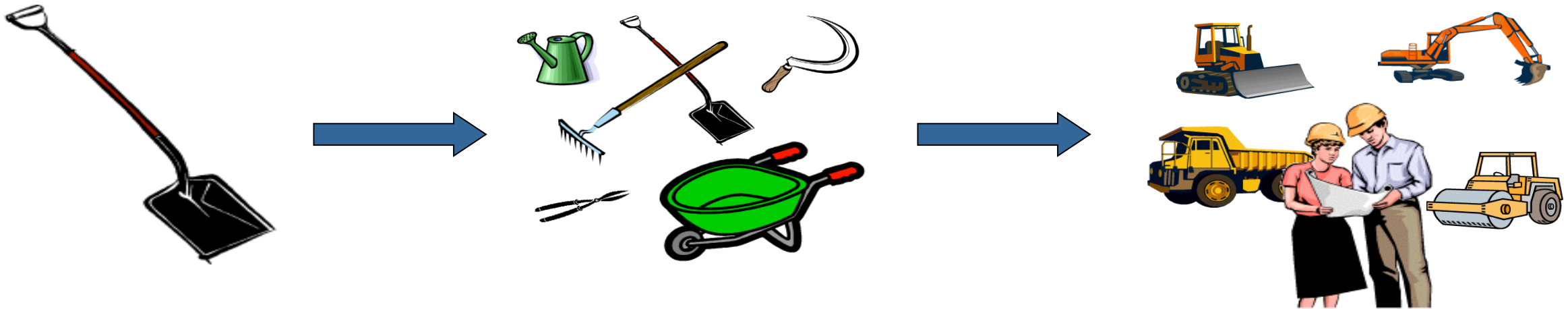
Example: Root-cause analysis - CESM Sea Ice Module



Example: Root-cause analysis - CESM Sea Ice Module



No single solution is sufficient!



A combination of different methods, tools and techniques is typically needed!

- Analysis
 - Statistics, visualization, automatic analysis, data mining, ...
- Measurement
 - Sampling / instrumentation, profiling / tracing, ...
- Instrumentation
 - Source code / binary, manual / automatic, ...

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