

Intel® MPI Library

Intel® MPI Library

Value Proposition

What	<ul style="list-style-type: none">• Intel's High Performance MPI Library
Why	<ul style="list-style-type: none">• Scale Performance – Tuned for Latest Intel Architectures• Scale Forward – Multicore and Manycore Ready• Scale Efficiently – Flexible Fabric Selection & Compatibility
How	<ul style="list-style-type: none">• Standards Based – Built on Open Source MPICH Implementation• Sustained Scalability – Tuning for Low Latencies, High Bandwidth & Increased Processes• Multi Fabric Support – Supports Popular High Performance Networking Fabrics

Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.
*Other names and brands may be claimed as the property of others.



Intel® MPI Library Overview

Streamlined product setup

- Install as root, or as standard user
- Environment variable script mpivars.(c)sh sets paths

Compilation scripts to handle details

- One set to use Intel compilers, one set for user-specified compilers

Environment variables for runtime control

- I_MPI_* variables control many factors at runtime
 - Process pinning, collective algorithms, device protocols, and more

Compiling MPI Programs

Compilation Scripts

- Automatically adds necessary links to MPI libraries and passes options to underlying compiler
- Use *mpiifort*, *mpiicpc*, or *mpiicc* to force usage of the associated Intel compiler
- Use *mpif77*, *mpicxx*, *mpicc*, or others to allow user to specify compiler (I_MPI_F77, ... or `-f77=`, `-cxx=`, ...)
 - Useful for makefiles portable between MPI implementations
- All compilers are found via PATH

MPI Launcher

Robust launch command

```
mpirun <mpi args> executable <program args>
```

Options available for:

- Rank distribution and pinning
- Fabric selection and control
- Environment propagation
- And more

Understanding MPI and Launcher Behavior

`I_MPI_DEBUG=<level>`

Debug Levels (cumulative):

- 0 – *Default*, no debug information
- 1 – Verbose error diagnostics
- 2 – Fabric selection process
- 3 – Rank, PID, node mapping
- 4 – Process pinning
- 5 – Display Intel® MPI Library environment variables
- 6 – Collective operation algorithm controls

`I_MPI_HYDRA_DEBUG=1` turns on Hydra debug output

- Keep in mind that this gives a LOT of output. Only turn on if needed

Process Placement

Default placement puts one rank per core on each node

Use `-ppn` to control processes per node

Use a machinefile to define ranks on each node individually

Use arguments sets or configuration files for precise control for complex jobs

Fabric Selection

`I_MPI_FABRICS=<intranode fabric>:<internode fabric> or <fabric>`

Fabric options

- shm – Shared Memory (only valid for intranode)
- dapl – Direct Access Provider Library*
- ofa – Open Fabric Alliance (OFED* verbs)
- tmi – Tag Matching Interface
- tcp – Ethernet/sockets
- ofi – OpenFabrics Interfaces*

Default behavior goes through a list to find first working fabric combination

If you specify a fabric, fallback is disabled, `I_MPI_FALLBACK=1` to re-enable

Environment Propagation

Use `-[g]env[*]` to control environment propagation

- Adding `g` propagates to all ranks, otherwise only to ranks in current argument set

-env <variable> <value> Set <variable> to <value>

-envuser All user environment variables, with a few exceptions (Default)

-envall All environment variables

-envnone No environment variables

-envlist <variable list> Only the listed variables

What's New: Intel® MPI Library 2018

- Up to 11x faster job start-up performance.
- Up to 25% reduction in job finalization time.
- Added support for the latest Intel® Xeon® Scalable processor.

Handling Heterogeneous Jobs

Global Options vs. Local Options

Global Options are applied to all ranks

- -ppn, -genv, ...

Local Options are applied to a subset of ranks

- -n, -host, -env, ...

WARNING: Some options can be set as local options via environment variable, but must be consistent across job

- Collective algorithms
- Fabric selection and parameters

Configuration Files and Argument Sets

Arguments Sets are used on the command line

Configuration Files are pulled from the file specified by *-configfile <configfile>*

Global arguments appear first (first line, or at beginning of first argument set)

Local arguments for each argument set next

Separated by : on command line (don't separate globals), new line in configfile

Can be used to run heterogeneous binaries, different arguments for each binary, different environment variables, etc.

All ranks combined in order specified into one job

Examples

Configuration File

```
$ cat theconfigfile
-genv OMP_NUM_THREADS 4
-n 6 -host node1 ./exe1
-n 4 -host node2 ./exe2
# -n 4 -host dead_node3 ./exe3
-n 6 -host node4 ./exe4
$ mpirun -configfile theconfigfile
```

Argument Set

```
$ mpirun -genv OMP_NUM_THREADS 4 -n 6 -host node1 ./exe1
: -n 4 -host node2 ./exe2 : -n 6 -host node4 ./exe4
```

Intel® Trace Analyzer and Collector (ITAC)

Intel® Trace Analyzer and Collector

Value Proposition

What	<ul style="list-style-type: none">• Intel's High Performance MPI Communications Profiler & Analyzer for Scalable HPC Development
Why	<ul style="list-style-type: none">• Scale Performance – Perform on More Nodes• Scale Forward – Multicore and Manycore Ready• Scale Efficiently – Tune & Debug on More Nodes
How	<ul style="list-style-type: none">• Visualize - Understand parallel application behavior• Evaluate - Profiling statistics and load balancing• Analyze – Automated analysis of common MPI issues• Identify – Communication hotspots

Intel® Trace Analyzer and Collector Overview

Intel® Trace Analyzer and Collector helps the developer:

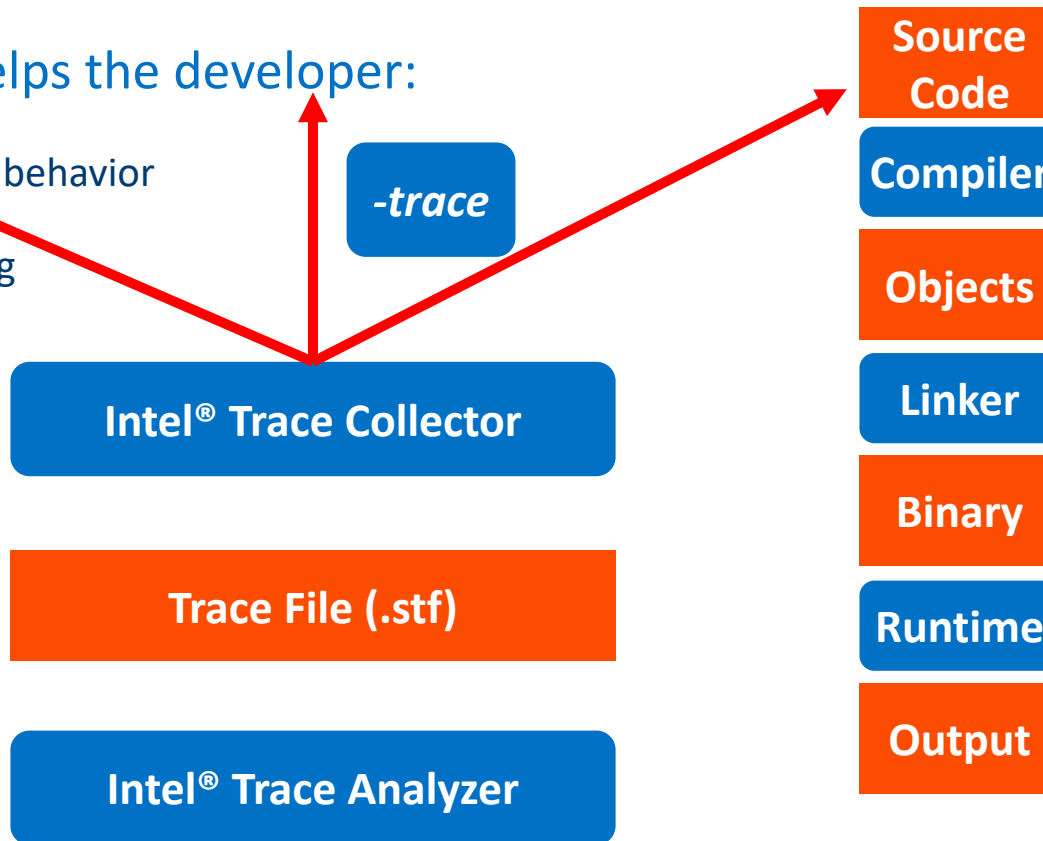
- Visualize and understand parallel application behavior
- Evaluate profiling scenarios
- Identify communication hotspots

API and *-tcollect*

-trace

Features

- Event-based approach
- Low overhead
- Excellent scalability
- Powerful aggregation and filtering functions
- Performance Assistance and Imbalance Tuning



Strengths of Event-based Tracing

Predict

Detailed MPI program behavior

Record

Exact sequence of program states – keep timing consistent

Collect

Collect information about exchange of messages: at what times and in which order

An event-based approach is able to detect temporal dependencies!

Multiple Methods for Data Collection

Collection Mechanism	Advantages	Disadvantages
Run with <code>-trace</code> or preload trace collector library.	Automatically collects all MPI calls, requires no modification to source, compile, or link.	No user code collection.
Link with <code>-trace</code> .	Automatically collects all MPI calls.	No user code collection. Must be done at link time.
Compile with <code>-tcollect</code> .	Automatically instruments all function entries/exits.	Requires recompile of code.
Add API calls to source code.	Can selectively instrument desired code sections.	Requires code modification.

Intel® Trace Analyzer Summary Page



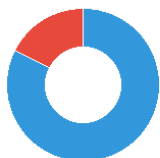
Summary: poisson.sendrecv.single.stf

Total time: 0.675 sec. Resources: 16 processes, 4 nodes.

Continue >

Ratio

This section represents a ratio of all MPI calls to the rest of your code in the application.



Serial Code	- 0.556 sec	82.4 %
OpenMP	- 0 sec	0 %
MPI calls	- 0.118 sec	17.5 %

Top MPI functions

This section lists the most active MPI functions from all MPI calls in the application.

MPI_Sendrecv	0.0643 sec (8.59 %)
MPI_Allreduce	0.0415 sec (5.54 %)
MPI_Finalize	0.00785 sec (1.05 %)
MPI_Bcast	0.00369 sec (0.494 %)
MPI_Errhandler_create	0.000239 sec (0.032 %)

Where to start with analysis

For deep analysis of the MPI-bound application click "Continue >" to open the tracefile View and leverage the Intel® Trace Analyzer functionality:

- *Performance Assistant* - to identify possible performance problems
- *Imbalance Diagram* - for detailed imbalance overview
- *Tagging/Filtering* - for thorough customizable analysis

To optimize node-level performance use:

Intel® VTune™ Amplifier for:

- algorithmic level tuning with hpc-performance and threading efficiency analysis;
- microarchitecture level tuning with general exploration and bandwidth analysis;

Intel® Advisor for:

- vectorization optimization and thread prototyping.

For more information, see documentation for the respective tool:

[Analyzing MPI applications with Intel® VTune™ Amplifier](#)

[Analyzing MPI applications with Intel® Advisor](#)

Show Summary Page when opening a tracefile

Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.

*Other names and brands may be claimed as the property of others.



Views and Charts

Helps navigating through the trace data and keep orientation

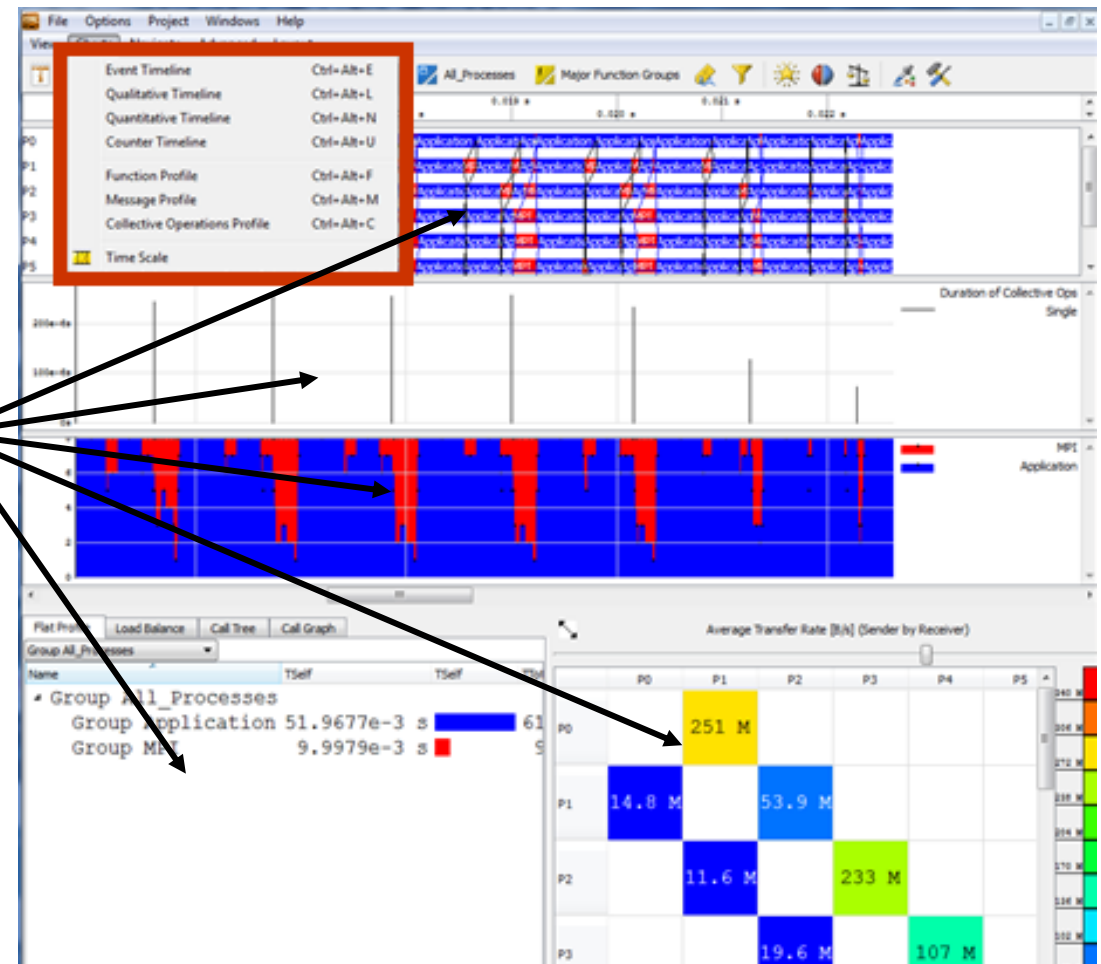
Every View can contain several Charts

All Charts in a View are linked to a single:

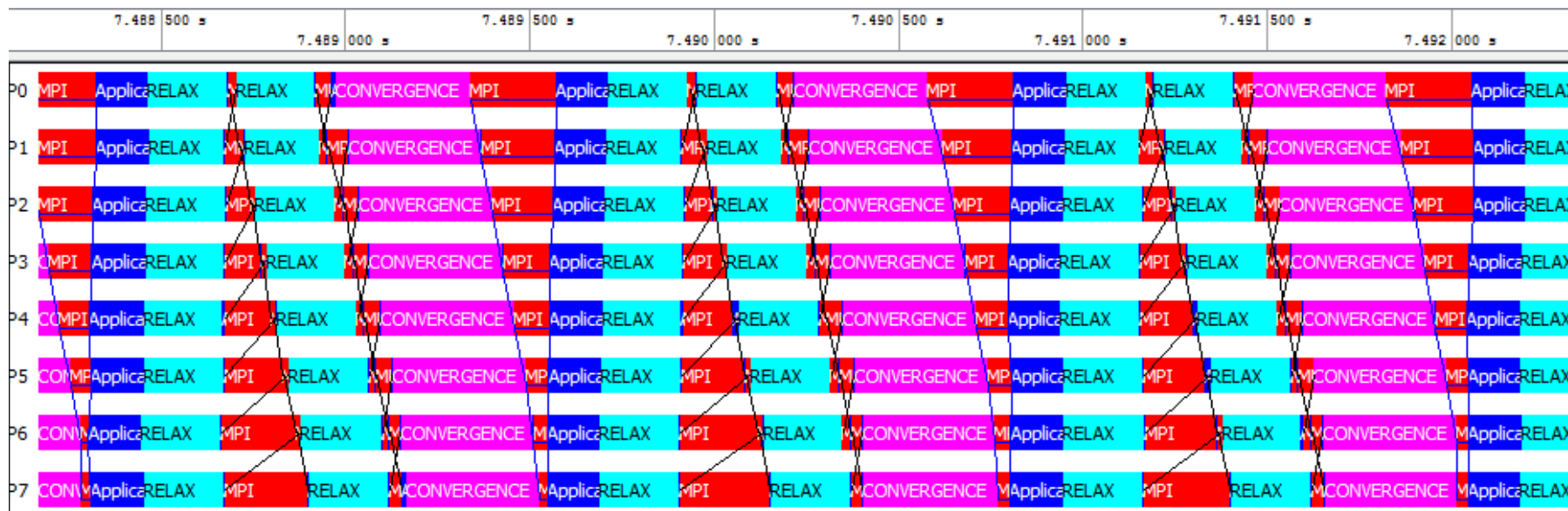
- time-span
- set of threads
- set of functions

All Charts follow changes to View (e.g. zooming)

Chart



Event Timeline



Get detailed impression of program structure

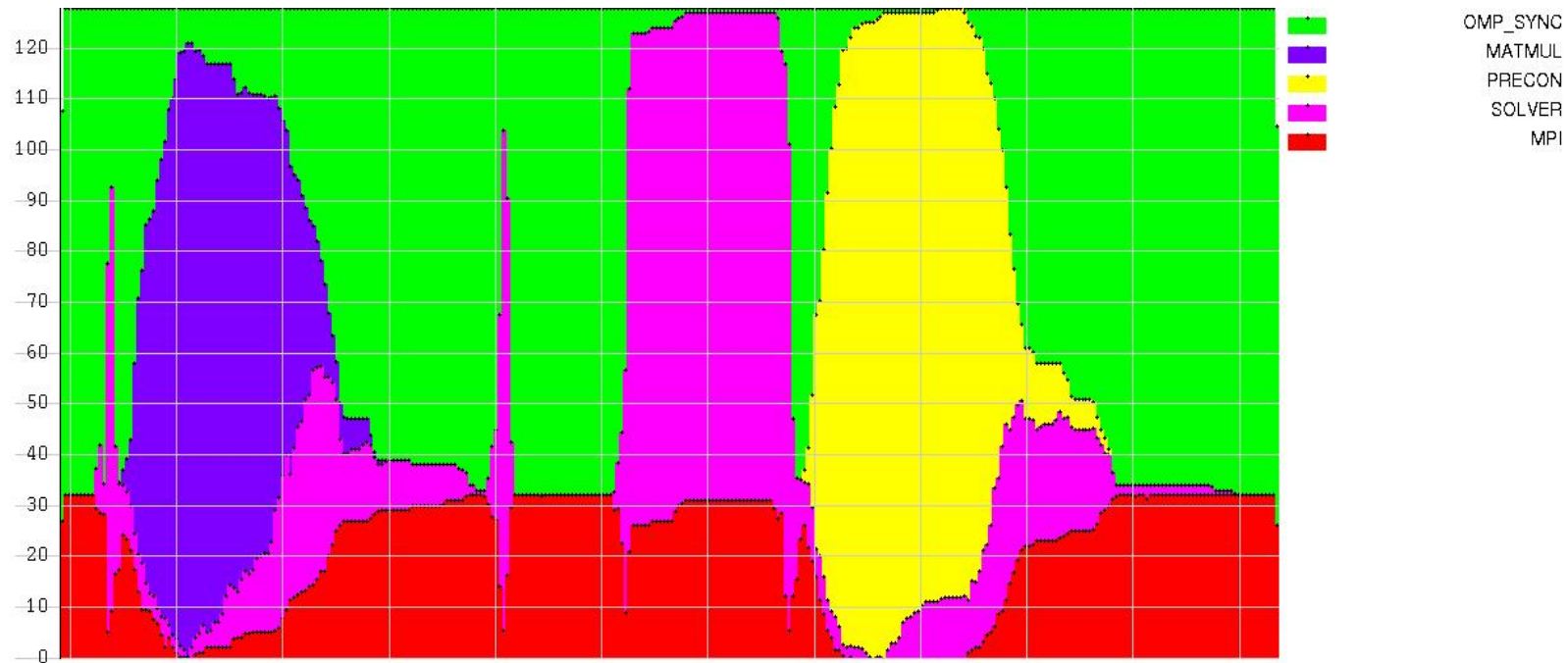
Display functions, messages, and collective operations for each rank/thread along time-axis

Retrieval of detailed event information

Quantitative Timeline

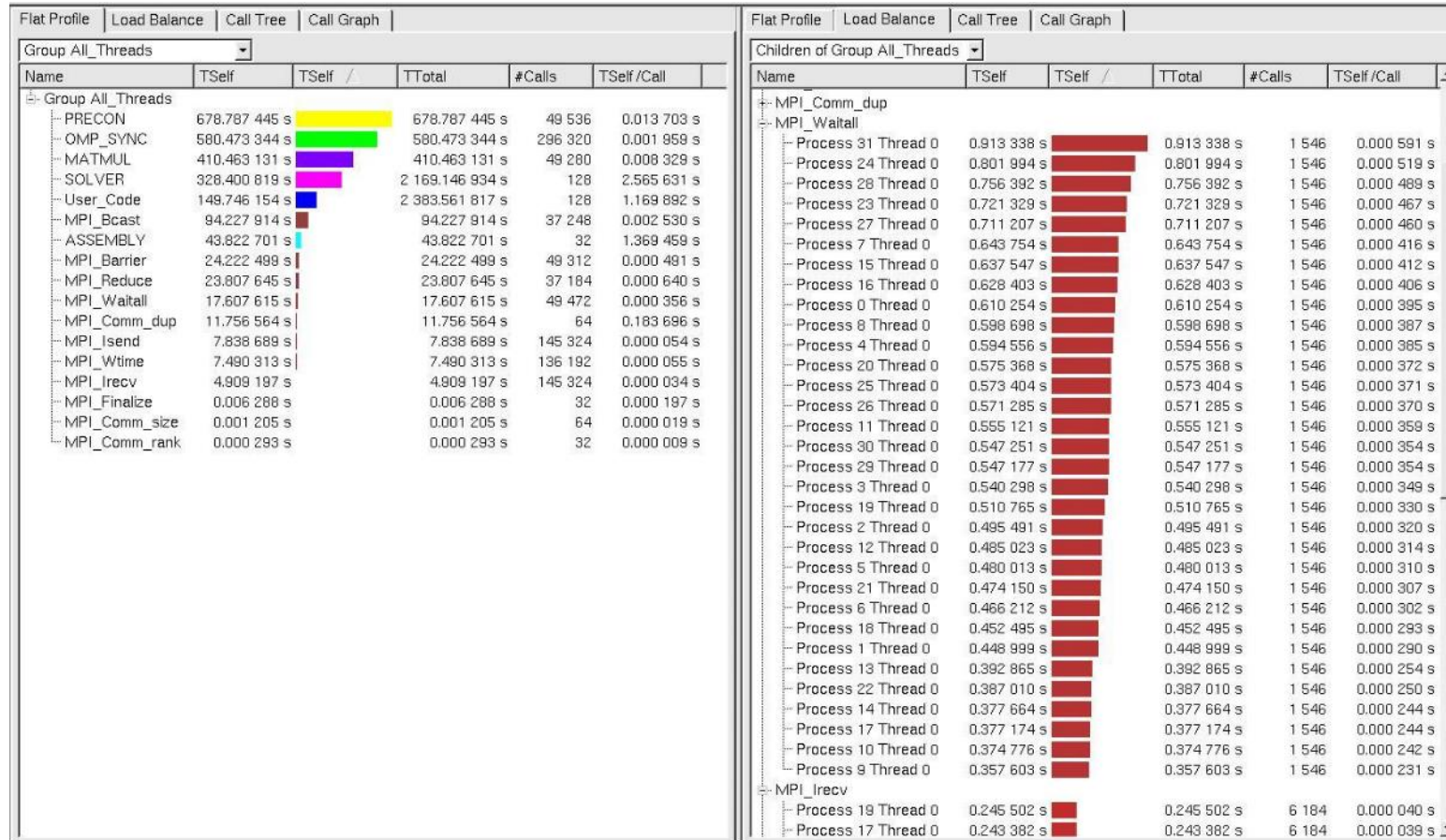
Get impression on parallelism and load balance

Show for every function how many threads/ranks are currently executing it



Flat Function Profile

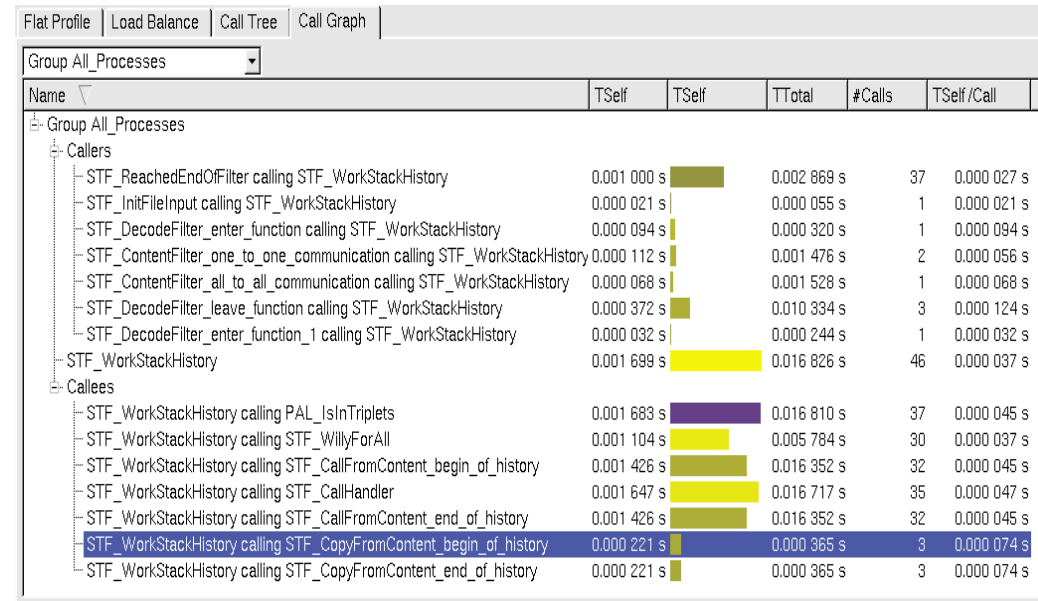
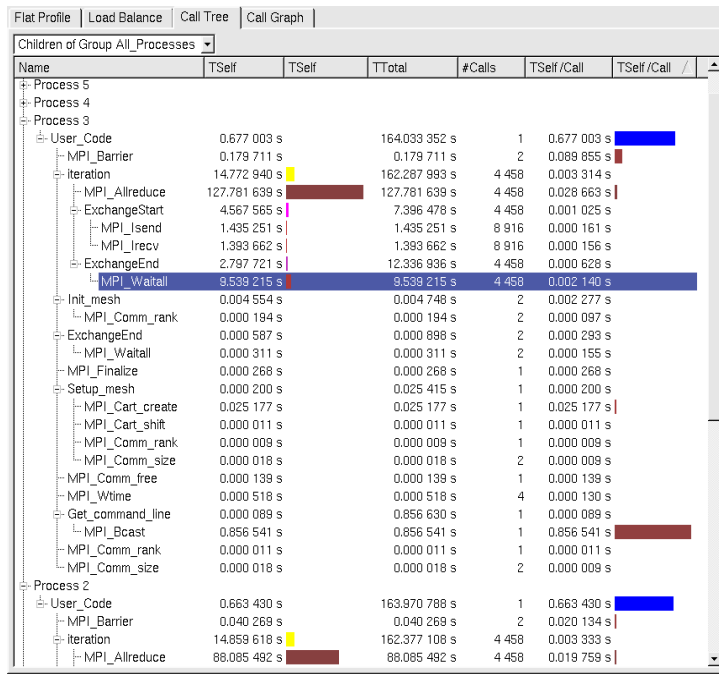
Statistics about functions



Call Tree and Call Graph

Function statistics including calling hierarchy

- Call Tree shows call stack
- Call Graph shows calling dependencies



Communication Profiles

Statistics about point-to-point or collective communication

Matrix supports grouping by attributes in each dimension

- Sender, Receiver, Data volume per msg, Tag, Communicator, Type

Available attributes

- Count, Bytes transferred, Time, Transfer rate

	P0	P1	P2	P3	P4	P5	P6	P7	Sum	Mean	StdDev
MPI_Barrier	0.063	0.052	0.040	0.180	0.258	0.066	0.079	0.215	0.952	0.119	0.080
MPI_Bcast	0.000	0.860	0.865	0.857	0.953	0.855	0.860	0.861	6.010	0.751	0.284
MPI_Allreduce	87.299	120.679	88.085	127.782	89.071	124.266	109.330	137.064	883.576	110.447	18.704
Sum	87.362	121.590	88.990	128.818	90.182	125.187	110.268	138.141	890.538		
Mean	29.121	40.530	29.663	42.939	30.061	41.729	36.756	46.047		37.106	
StdDev	41.139	56.675	41.312	59.993	41.727	58.363	51.318	64.359			52.973

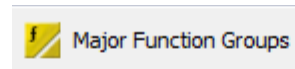
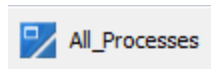
	P0	P1	P2	P3	P4	P5	P6	P7	Sum	Mean	StdDev
P0		74.641							74.641	74.641	0.000
P1	23.903		45.249						69.152	34.576	10.873
P2		51.590		47.361					99.551	49.776	1.014
P3			41.605		36.904				78.509	39.254	2.351
P4				51.558		54.114			105.672	52.836	1.278
P5					37.884		34.262		72.146	36.073	1.811
P6						37.619		35.861	73.480	36.740	0.879
P7							24.384		24.384	24.384	0.000
Sum	23.903	126.231	86.854	99.519	74.788	91.793	58.646	35.861	597.535		
Mean	23.903	63.116	43.427	49.759	37.394	45.866	29.323	35.861		42.681	
StdDev	0.000	11.526	1.822	1.798	0.490	8.248	4.939	0.000			12.629

Grouping and Aggregation

Allow analysis on different levels of detail by aggregating data upon group-definitions

Functions and threads can be grouped hierarchically

- Process Groups and Function Groups



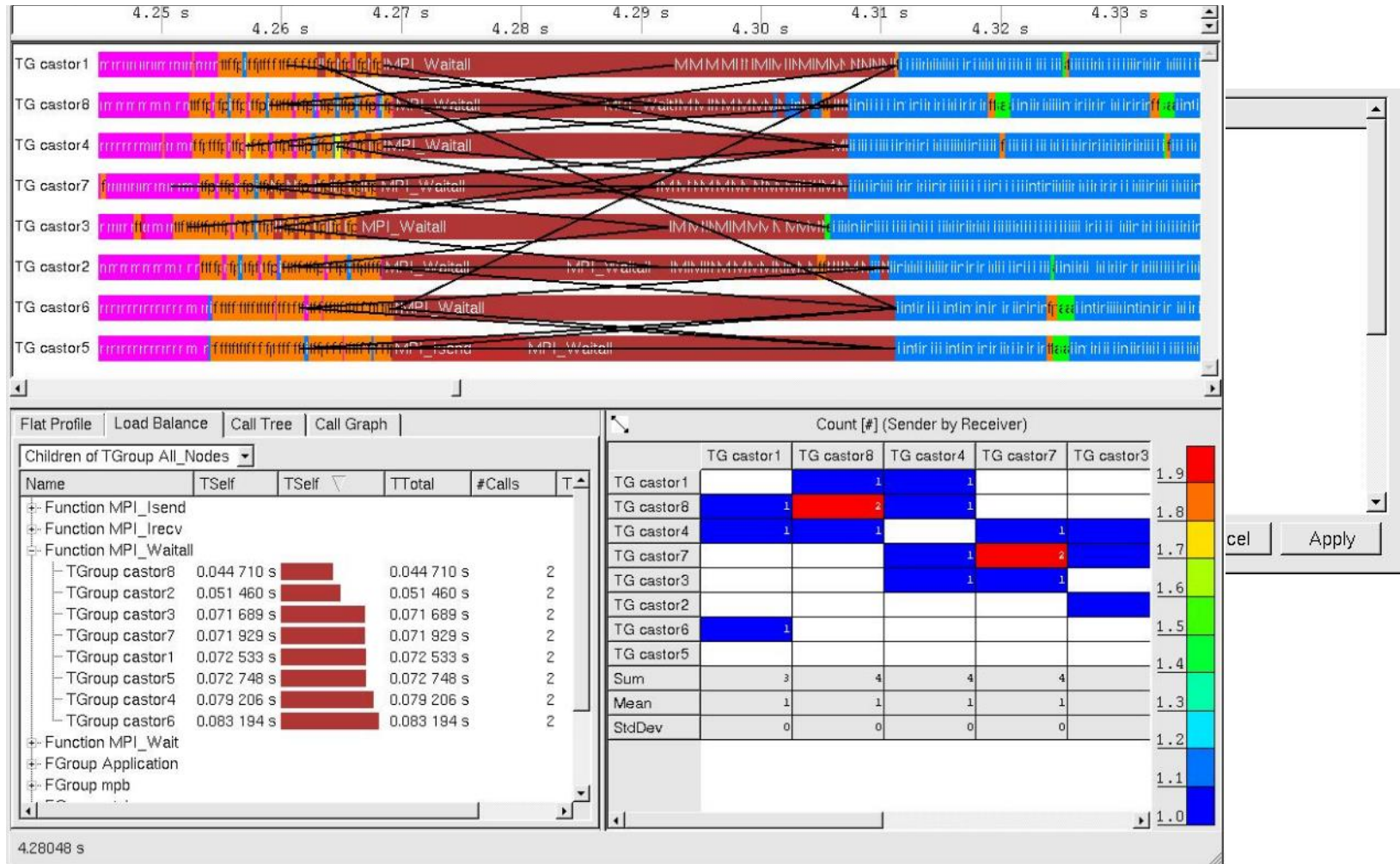
Arbitrary nesting is supported

- Functions/threads on the same level as groups
- User can define his/her own groups

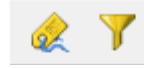
Aggregation is part of View-definition

- All charts in a View adapt to requested grouping
- All charts support aggregation

Aggregation Example



Tagging and Filtering



Help concentrating on relevant parts

Avoid getting lost in huge amounts of trace data

Define a set of interesting data

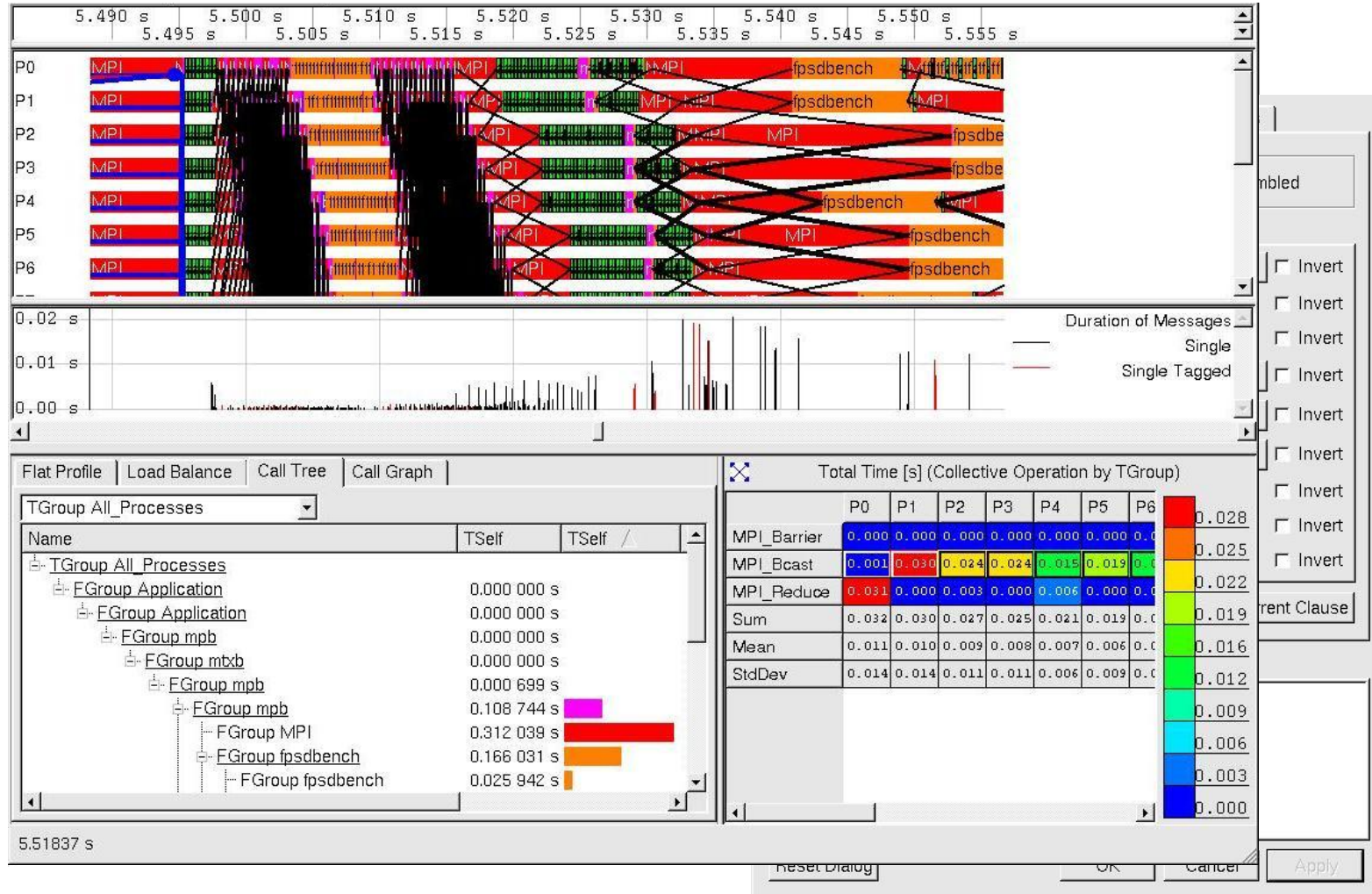
- E.g. all occurrences of function x
- E.g. all messages with tag y on communicator z

Combine several filters:
Intersection, Union, Complement

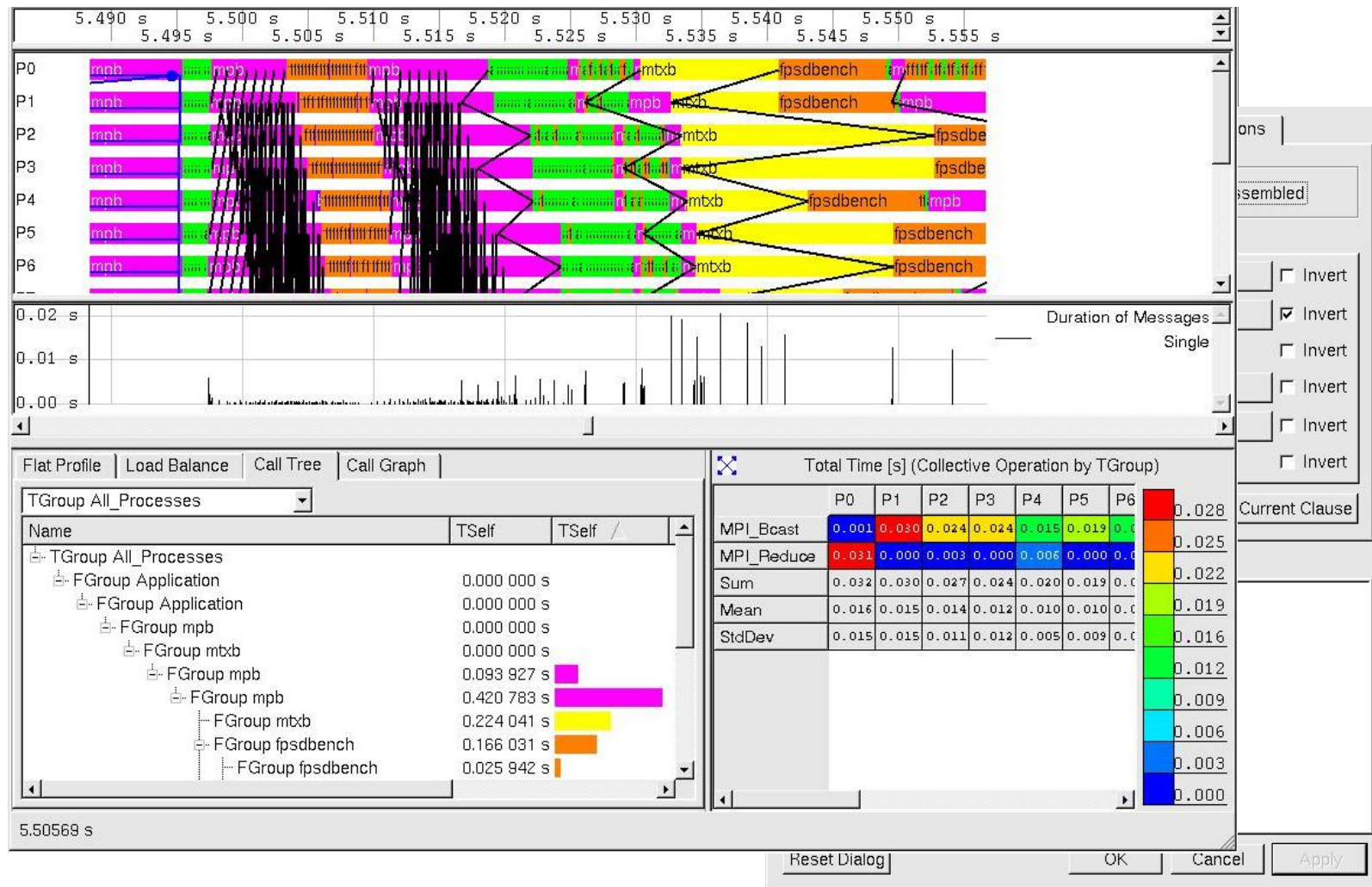
Apply it

- Tagging: Highlight messages
- Filtering: Suppress all non-matching events

Tagging Example



Filtering Example



Optimization Notice

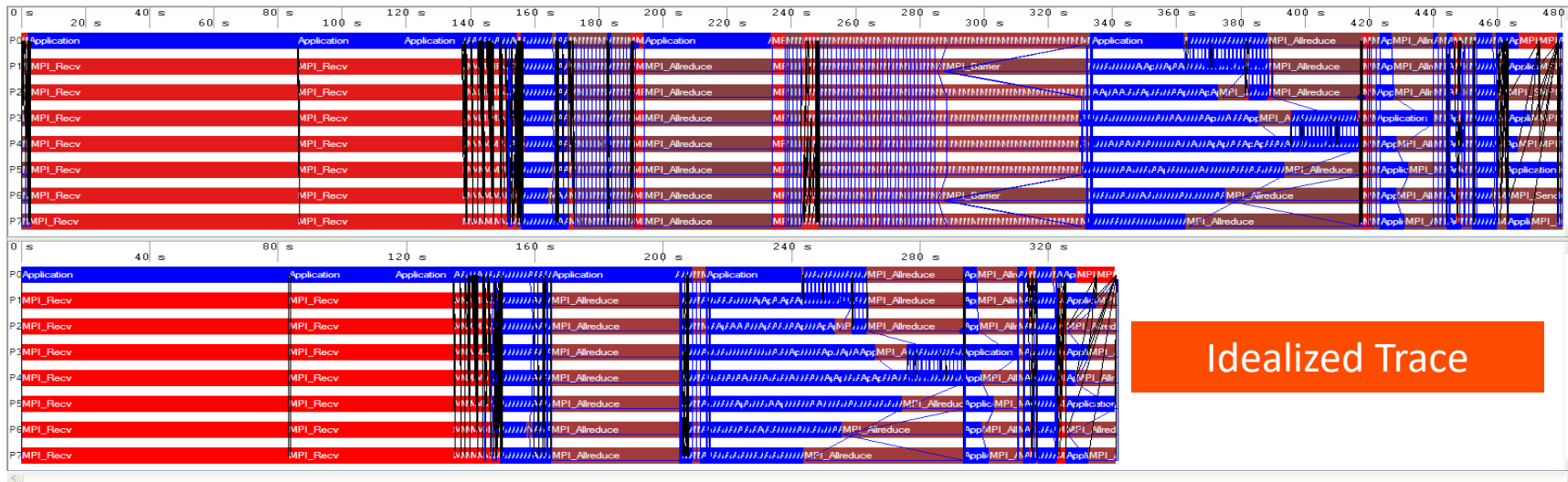
Copyright © 2017, Intel Corporation. All rights reserved.
 *Other names and brands may be claimed as the property of others.



Ideal Interconnect Simulator (Idealizer)

Helps to figure out application's imbalance simulating its behavior in the "ideal communication environment"

Actual trace



Idealized Trace

Easy way to identify application bottlenecks

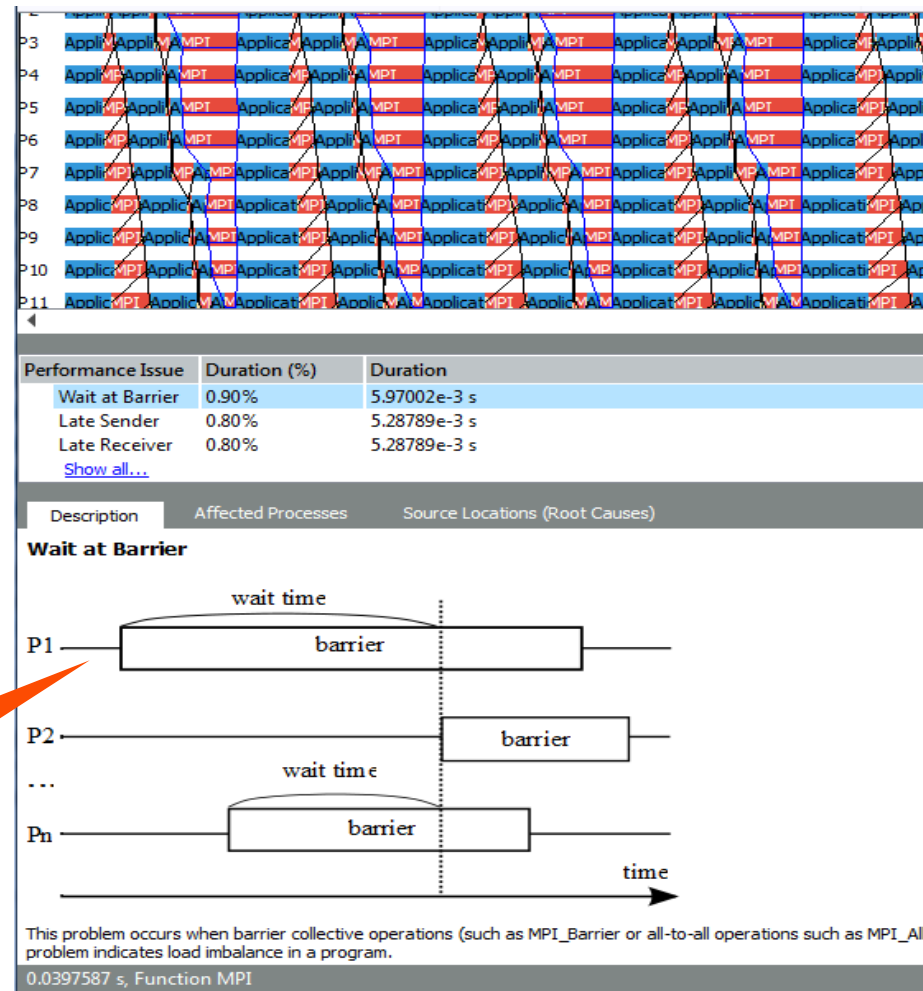
MPI Performance Assistance

Automatic Performance Assistant

Detect common MPI performance issues

Automated tips on potential solutions

Automatically detect performance issues and their impact on runtime

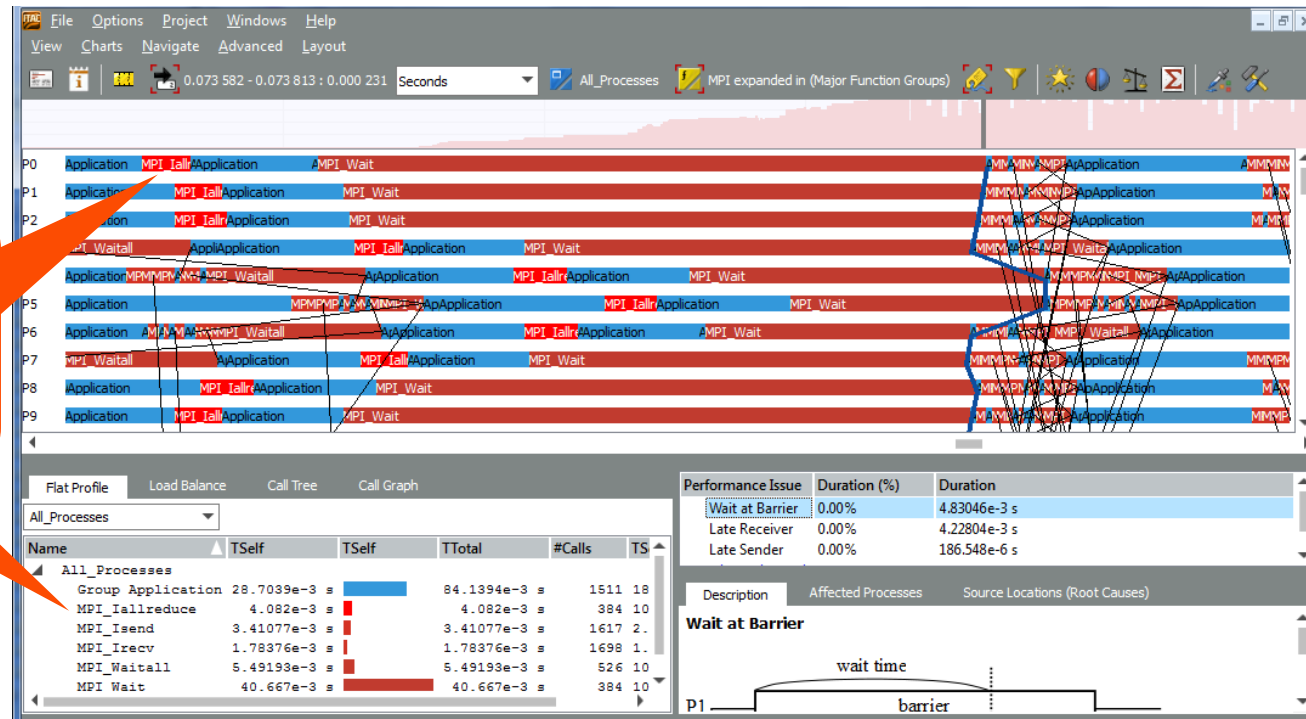


MPI-3.0 Support

Support for major MPI-3.0 features

- Non-blocking collectives
- Fast RMA
- Large counts

Non-blocking
Allreduce
(MPI_iallreduce)



What's New in Intel® Trace Analyzer and Collector 2018

Support for OpenSHMEM* applications.

Support for the latest Intel® Xeon® Scalable and Intel® Xeon Phi™ processors.