



# INTEL® PARALLEL STUDIO XE 2017 CLUSTER EDITION OVERVIEW

For Distributed Performance

# Intel® Parallel Studio XE 2017 development suite

## Empowering Faster Code Faster

Delivering HPC Development Solutions

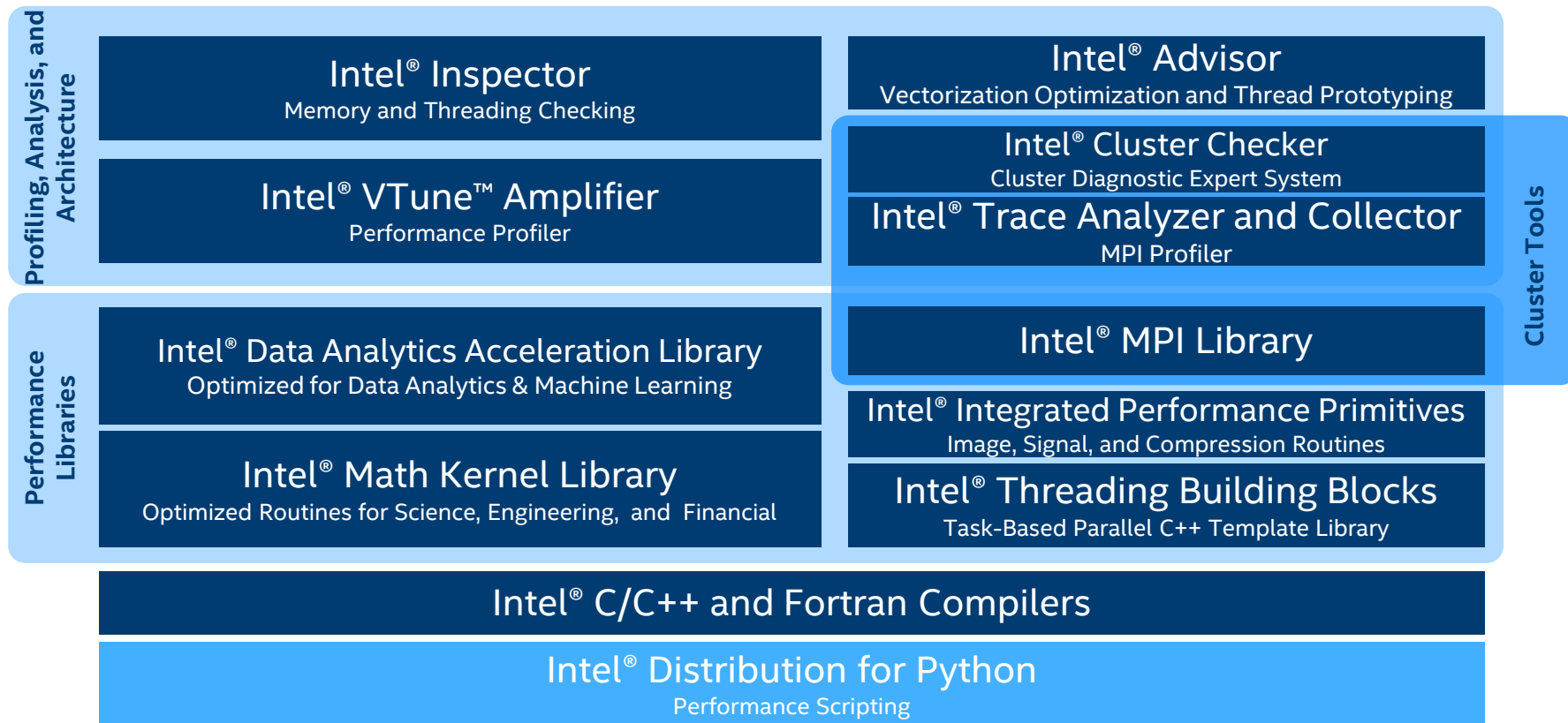
- Over 20 years
- Industry Collaboration on Standards
- Developed with Performance & Scaling with Intel hardware

### Meeting the Challenges

- Boosting Performance
- Increasing Scalability
- Increasing Productivity



# Intel® Parallel Studio XE



## Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.

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



**INTEL<sup>®</sup> MPI LIBRARY**

# Intel® MPI Library Overview

## Optimized MPI application performance

- Application-specific tuning
- Automatic tuning
- New! - Support for Intel® Xeon Phi™ processor (code-named Knights Landing)
- New: Support for Intel® Omni-Path Architecture Fabric

## Lower-latency and multi-vendor interoperability

- Industry leading latency
- Performance optimized support for the fabric capabilities through OpenFabrics\*(OFI)

## Faster MPI communication

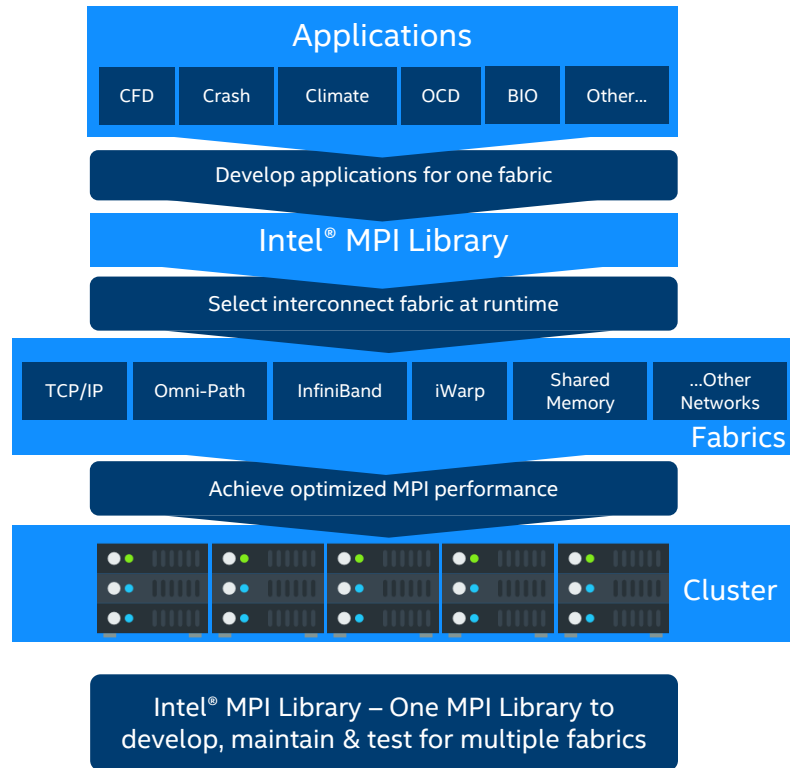
- Optimized collectives

## Sustainable scalability up to 340K cores

- Native InfiniBand\* interface support allows for lower latencies, higher bandwidth, and reduced memory requirements

## More robust MPI applications

- Seamless interoperability with Intel® Trace Analyzer and Collector



# 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

# 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
```



# TUNING MPI APPLICATION PERFORMANCE

# Tuning Methods

Library Tuning (algorithms, fabric parameters)

- mpitune

Application Tuning (load balance, MPI/threaded/serial performance)

- Intel® Trace Analyzer and Collector
- Application Performance Snapshot
- Intel® VTune™ Amplifier XE

# Library Tuning: mpitune

Use the automatic tuning facility to tune the Intel® MPI Library for your cluster or application (done once, may take a long time)

Modes (see `mpitune -h` for options)

- Cluster-wide tuning
- Application-specific tuning

```
mpitune ...
```

```
mpitune -application \"mpirun -n 32 ./exe\" ...
```

Creates options settings which are used with the `-tune` flag

```
mpirun -tune ...
```

# **INTEL<sup>®</sup> TRACE ANALYZER AND COLLECTOR (ITAC)**

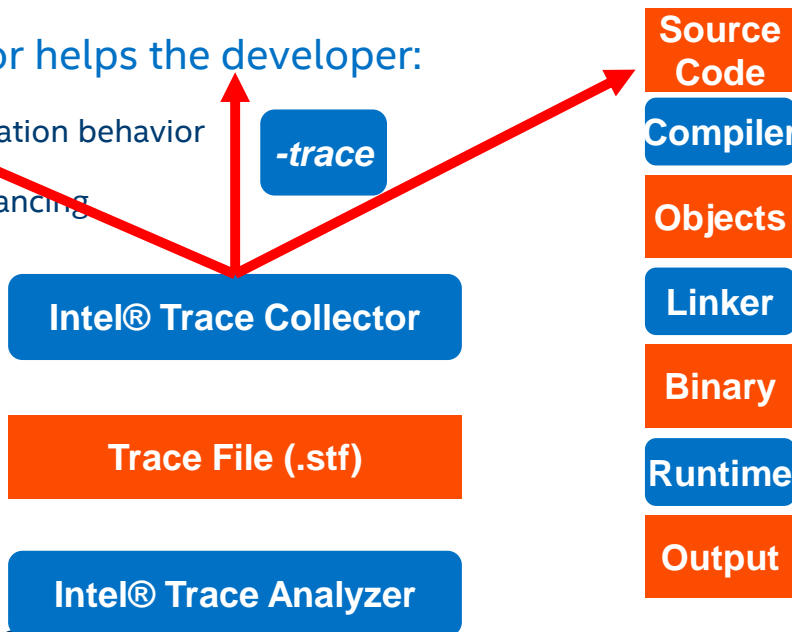
# Intel® Trace Analyzer and Collector Overview

Intel® Trace Analyzer and Collector helps the developer:

- Visualize and understand parallel application behavior
- Evaluate profiling **API and *-tcollect*** balancing
- Identify communication hotspots

## 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

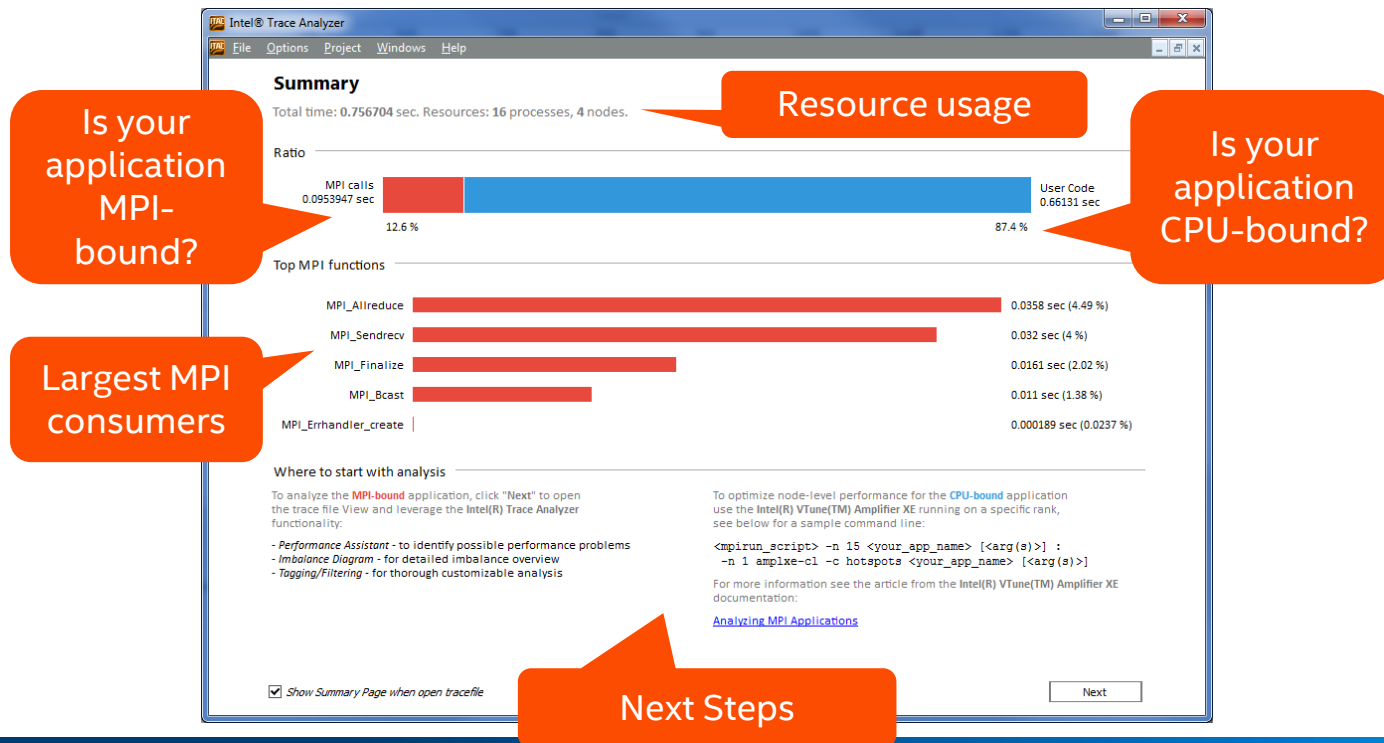
Collect information about exchange of messages: at what times and in which order

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

# Multiple Methods for Data Collection

Collection Mechanism	Advantages	Disadvantages
<b>Run with <code>-trace</code> or preload trace collector library.</b>	<b>Automatically collects all MPI calls, requires no modification to source, compile, or link.</b>	<b>No user code collection.</b>
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.

# Summary page shows computation vs. communication breakdown





# Views and Charts

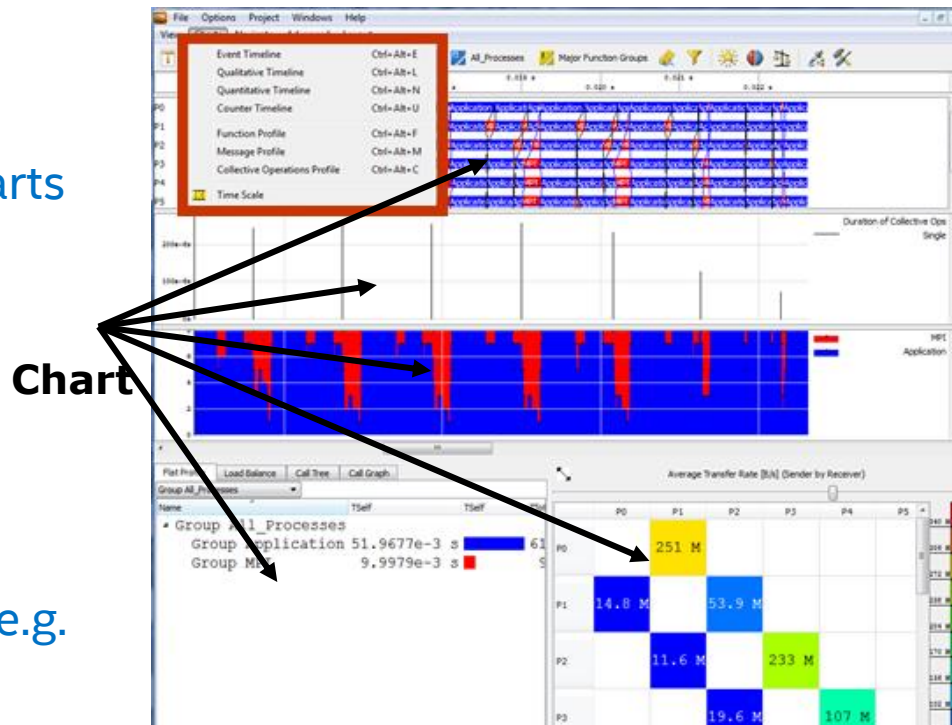
Helps navigating through the trace data and keep orientation

Every View can contain several Charts

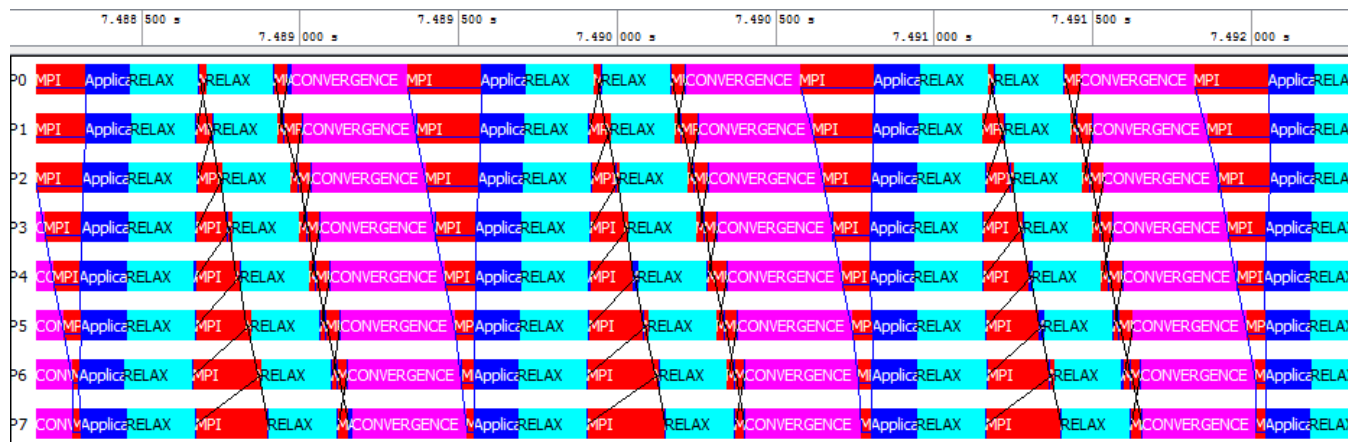
All Charts in a View are linked to a single:

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

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



# Event Timeline



Get detailed impression of program structure

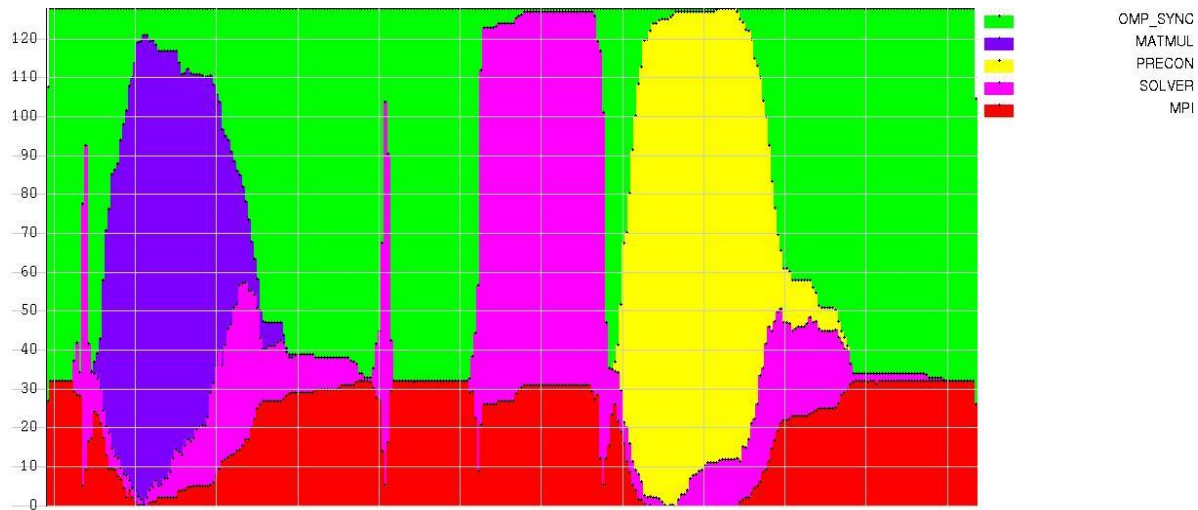
Display functions, messages, and collective operations for each rank/thread along time-axis

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

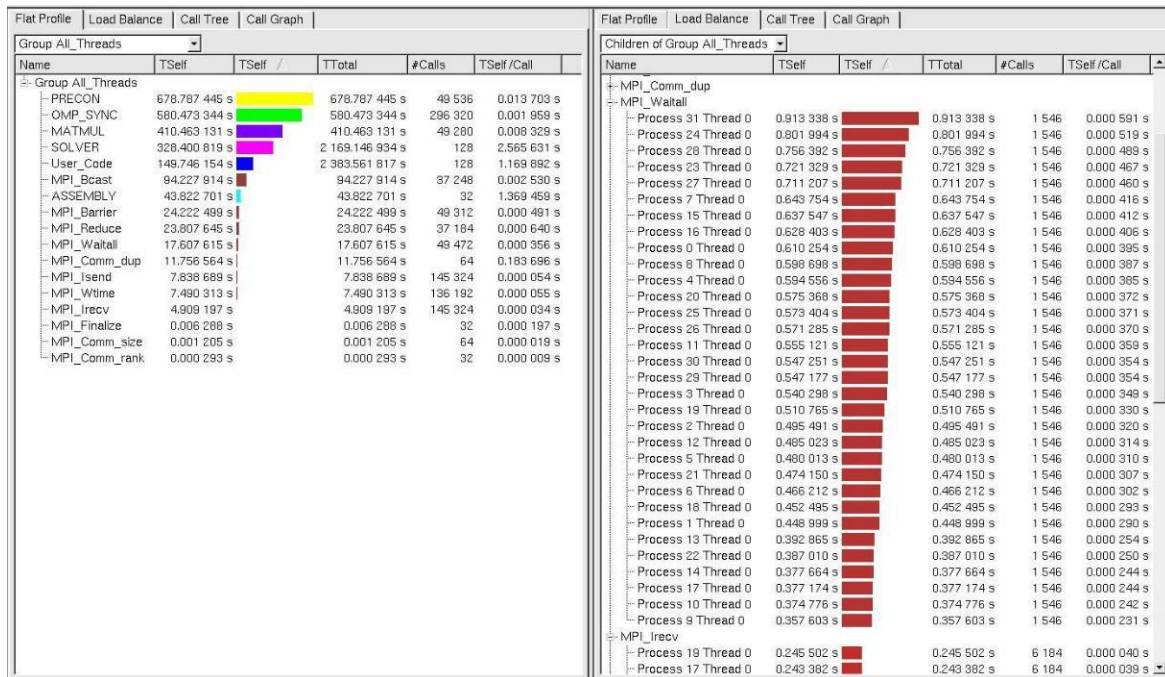


## Optimization Notice

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

# Flat Function Profile

## Statistics about functions



### Optimization Notice

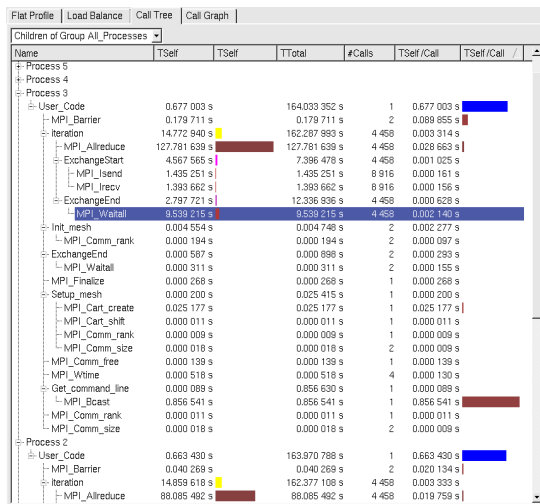
Copyright © 2017, Intel Corporation. All rights reserved.  
\*Other names and brands may be claimed as the property of others.



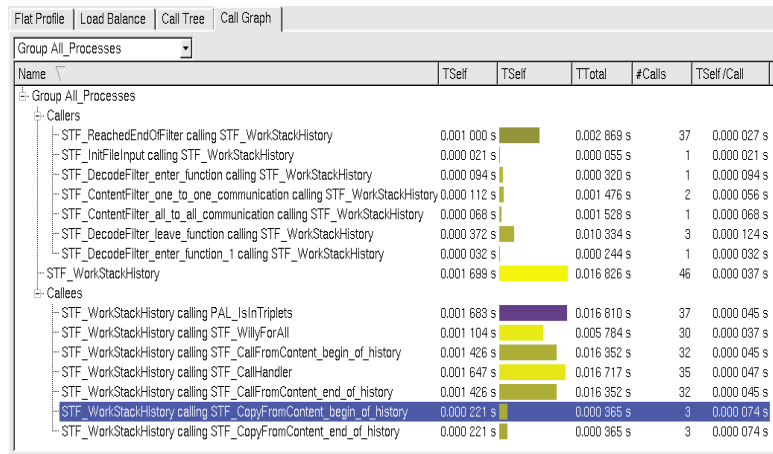
# Call Tree and Call Graph

## Function statistics including calling hierarchy

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



Name	TSelf	TSelf	TTotol	#Calls	TSelf/Call	TSelf/Call
Process 5						
Process 4						
Process 3						
User_Code	0.677 003 s		164.033 352 s	1	0.677 003 s	
MPI_Barrier	0.179 711 s		0.179 711 s	2	0.089 855 s	
Iteration	14.772 940 s		162.287 893 s	4 458	0.003 314 s	
MPI_Allreduce	12.781 639 s		12.781 639 s	4 458	0.028 663 s	
ExchangeStart	4.567 565 s		7.396 478 s	4 458	0.001 025 s	
MPI_Isend	1.435 251 s		1.435 251 s	8 916	0.000 161 s	
MPI_Irecv	1.393 662 s		1.393 662 s	8 916	0.000 156 s	
ExchangeEnd	2.797 721 s		12.336 938 s	4 458	0.000 628 s	
MPI_Waitall	0.000 000 s		0.000 000 s	4 458	0.000 000 s	
Init_mesh	0.004 554 s		0.004 749 s	2	0.002 277 s	
MPI_Comm_rank	0.000 194 s		0.000 194 s	2	0.000 097 s	
ExchangeEnd	0.000 587 s		0.000 898 s	2	0.000 233 s	
MPI_Waitall	0.000 311 s		0.000 311 s	2	0.000 155 s	
MPI_Finalize	0.000 268 s		0.000 268 s	1	0.000 268 s	
Setup_mesh	0.000 200 s		0.025 415 s	1	0.000 200 s	
MPI_Cart_create	0.025 177 s		0.025 177 s	1	0.025 177 s	
MPI_Cart_shift	0.000 011 s		0.000 011 s	1	0.000 011 s	
MPI_Comm_rank	0.000 009 s		0.000 009 s	1	0.000 009 s	
MPI_Comm_size	0.000 018 s		0.000 018 s	2	0.000 009 s	
MPI_Comm_free	0.000 139 s		0.000 139 s	1	0.000 139 s	
MPI_Wtime	0.000 518 s		0.000 518 s	4	0.000 130 s	
Get_command_line	0.000 089 s		0.856 630 s	1	0.000 089 s	
MPI_Bcast	0.856 541 s		0.856 541 s	1	0.856 541 s	
MPI_Comm_rank	0.000 011 s		0.000 011 s	1	0.000 011 s	
MPI_Comm_size	0.000 018 s		0.000 018 s	2	0.000 009 s	
Process 2						
User_Code	0.663 430 s		163.970 788 s	1	0.663 430 s	
MPI_Barrier	0.040 269 s		0.040 269 s	2	0.020 134 s	
Iteration	14.059 618 s		162.377 108 s	4 458	0.003 333 s	
MPI_Allreduce	88.085 492 s		88.085 492 s	4 458	0.019 759 s	



Name	TSelf	TSelf	TTotol	#Calls	TSelf/Call	TSelf/Call
Group All_Processes						
Callers						
STF_ReachedEndOfFilter calling STF_WorkStackHistory	0.001 000 s		0.002 868 s	37	0.000 027 s	
STF_InitFileInput calling STF_WorkStackHistory	0.000 021 s		0.000 055 s	1	0.000 021 s	
STF_DecodeFilter_enter_function calling STF_WorkStackHistory	0.000 094 s		0.000 320 s	1	0.000 094 s	
STF_ContentFilter_one_to_one_communication calling STF_WorkStackHistory	0.000 112 s		0.001 476 s	2	0.000 056 s	
STF_ContentFilter_all_to_all_communication calling STF_WorkStackHistory	0.000 068 s		0.001 528 s	1	0.000 068 s	
STF_DecodeFilter_leave_function calling STF_WorkStackHistory	0.000 372 s		0.010 334 s	3	0.000 124 s	
STF_DecodeFilter_enter_function_1 calling STF_WorkStackHistory	0.000 032 s		0.000 244 s	1	0.000 032 s	
STF_WorkStackHistory	0.001 699 s		0.016 826 s	46	0.000 037 s	
Callees						
STF_WorkStackHistory calling PAL_IsInTriplets	0.001 683 s		0.016 810 s	37	0.000 045 s	
STF_WorkStackHistory calling STF_WillyForAll	0.001 104 s		0.005 784 s	30	0.000 037 s	
STF_WorkStackHistory calling STF_CallFromContent_begin_of_history	0.001 426 s		0.016 352 s	32	0.000 045 s	
STF_WorkStackHistory calling STF_CallHandler	0.001 647 s		0.016 717 s	35	0.000 047 s	
STF_WorkStackHistory calling STF_CallFromContent_end_of_history	0.001 426 s		0.016 352 s	32	0.000 045 s	
STF_WorkStackHistory calling STF_CopyFromContent_begin_of_history	0.000 221 s		0.000 365 s	3	0.000 074 s	
STF_WorkStackHistory calling STF_CopyFromContent_end_of_history	0.000 221 s		0.000 365 s	3	0.000 074 s	

### Optimization Notice

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



# 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

Total Time [s] (Collective Operation by Process)											
	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	0.853	0.855	0.860	0.861	6.010	0.751	0.284
MPI_Allreduce	87.299	120.579	88.085	127.782	89.071	124.261	109.330	137.064	883.576	110.447	18.704
Sum	87.362	121.590	88.990	128.818	90.182	125.187	110.248	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	58.263	51.318	64.359			52.973

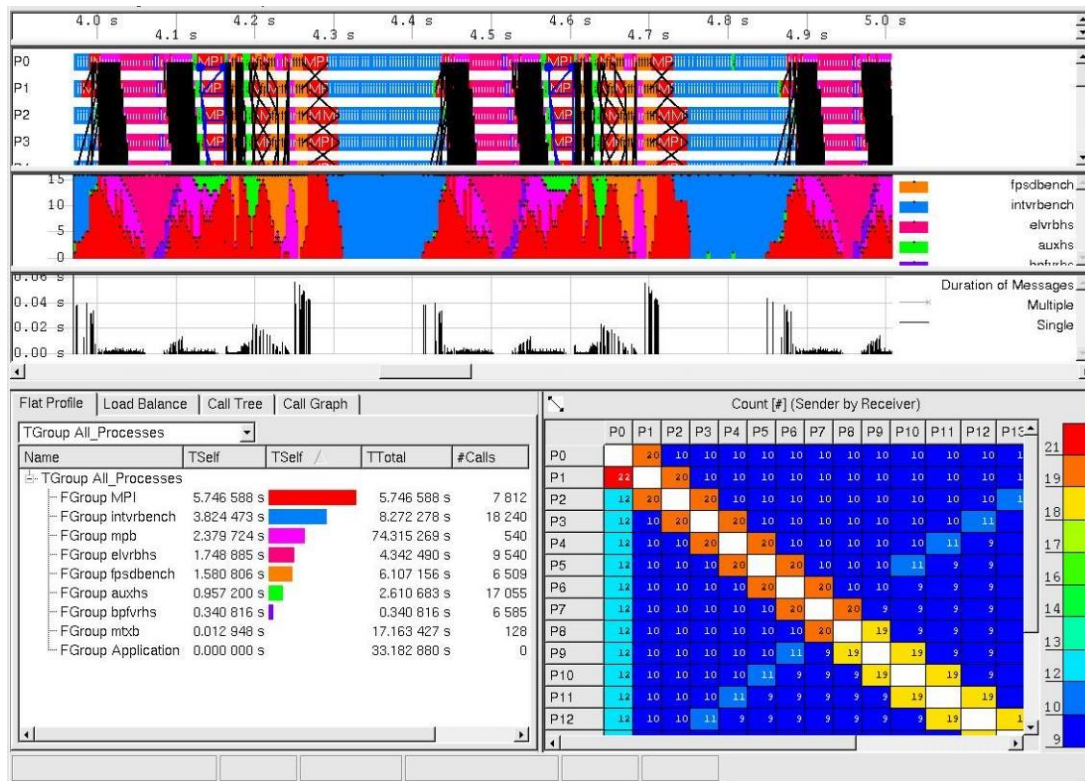
Total Time [s] (Sender by Receiver)											
	P0	P1	P2	P3	P4	P5	P6	P7	Sum	Mean	StdDev
P0		78.641							74.641	74.641	0.000
P1	21.901		45.149						69.152	34.576	10.673
P2		31.539		47.964					99.551	49.776	1.814
P3			41.605		36.994				78.509	39.254	2.351
P4				51.558		54.114			105.672	52.836	1.278
P5					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	124.231	86.854	99.519	74.788	91.793	58.646	35.861	597.535		
Mean	23.903	62.116	43.427	49.759	37.394	45.896	29.323	35.861	42.681		
StdDev	0.000	11.526	1.822	1.798	0.430	0.248	4.939	0.000			12.629

## Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.

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

# Zooming

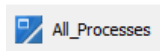


# Grouping and Aggregation

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

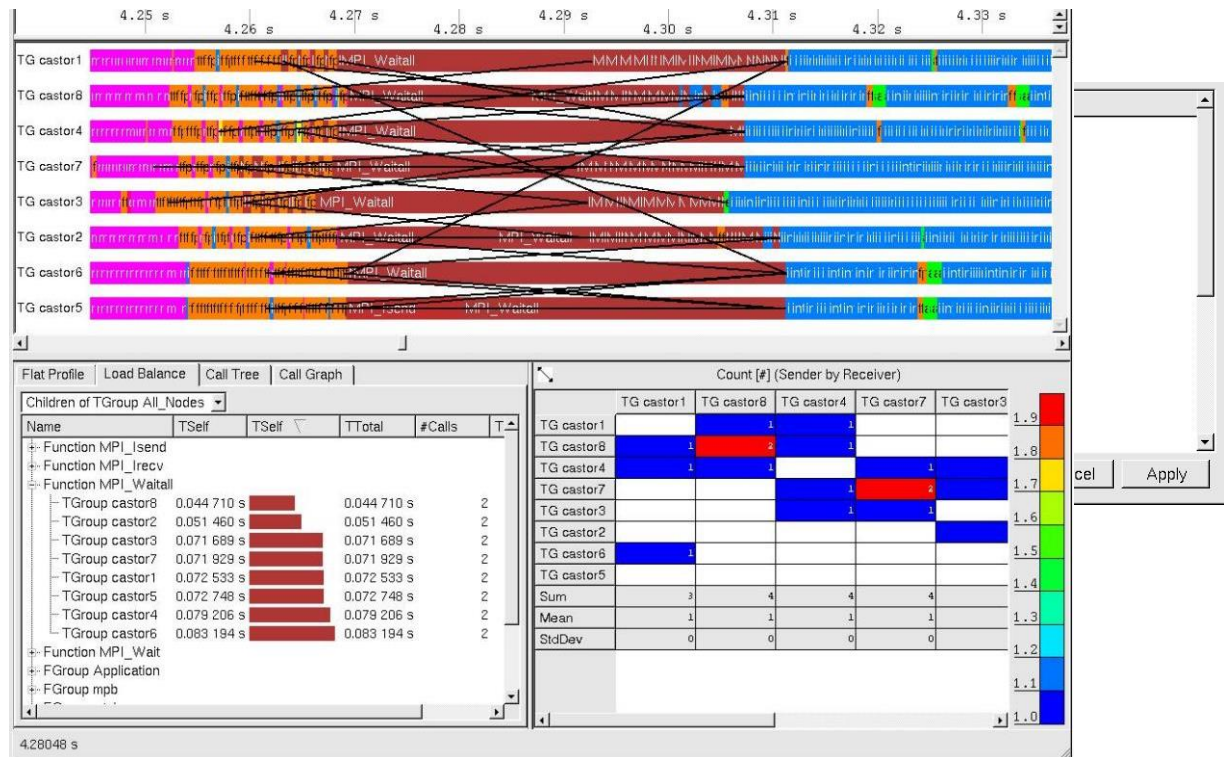
- Functions/threads on the same level as groups
- User can define his/her own groups

Aggregation is part of View-definition

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



# Aggregation Example



## Optimization Notice

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



# Tagging and Filtering



Help concentrating on relevant parts

Avoid getting lost in huge amounts of trace data

Define a set of interesting data

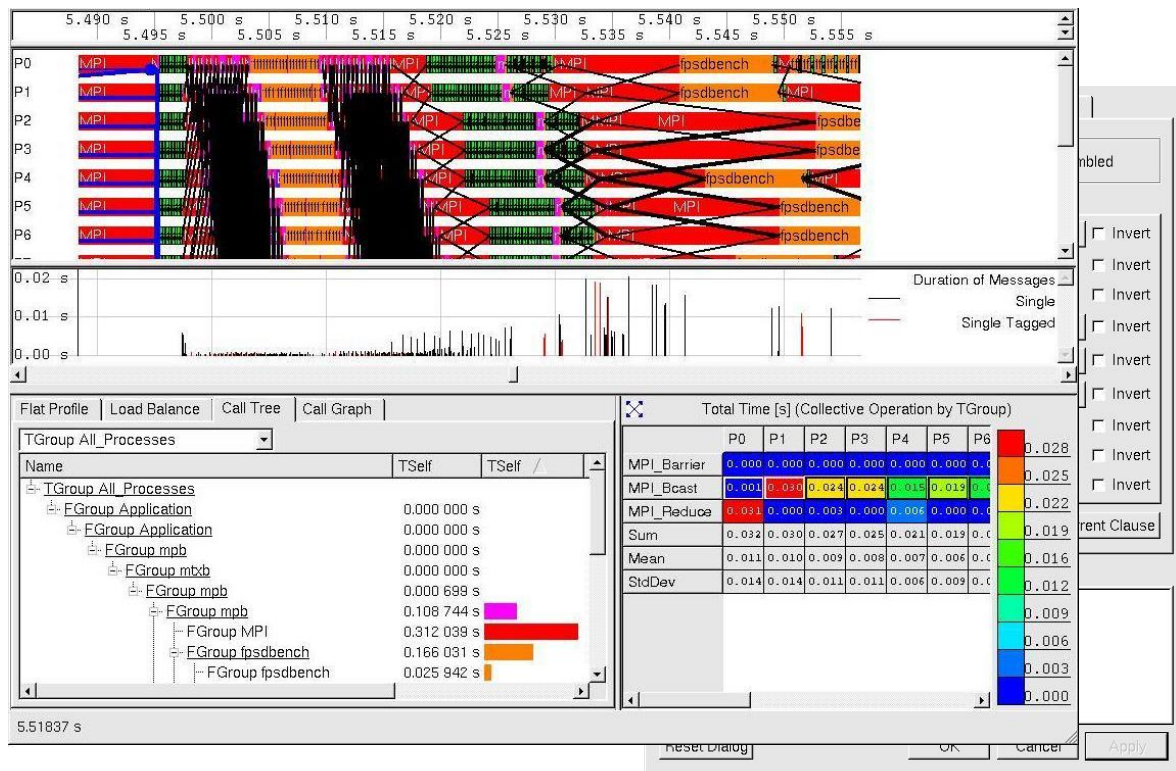
- E.g. all occurrences of function x
- E.g. all messages with tag y on communicator z

Combine several filters:  
Intersection, Union, Complement

Apply it

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

# Tagging Example

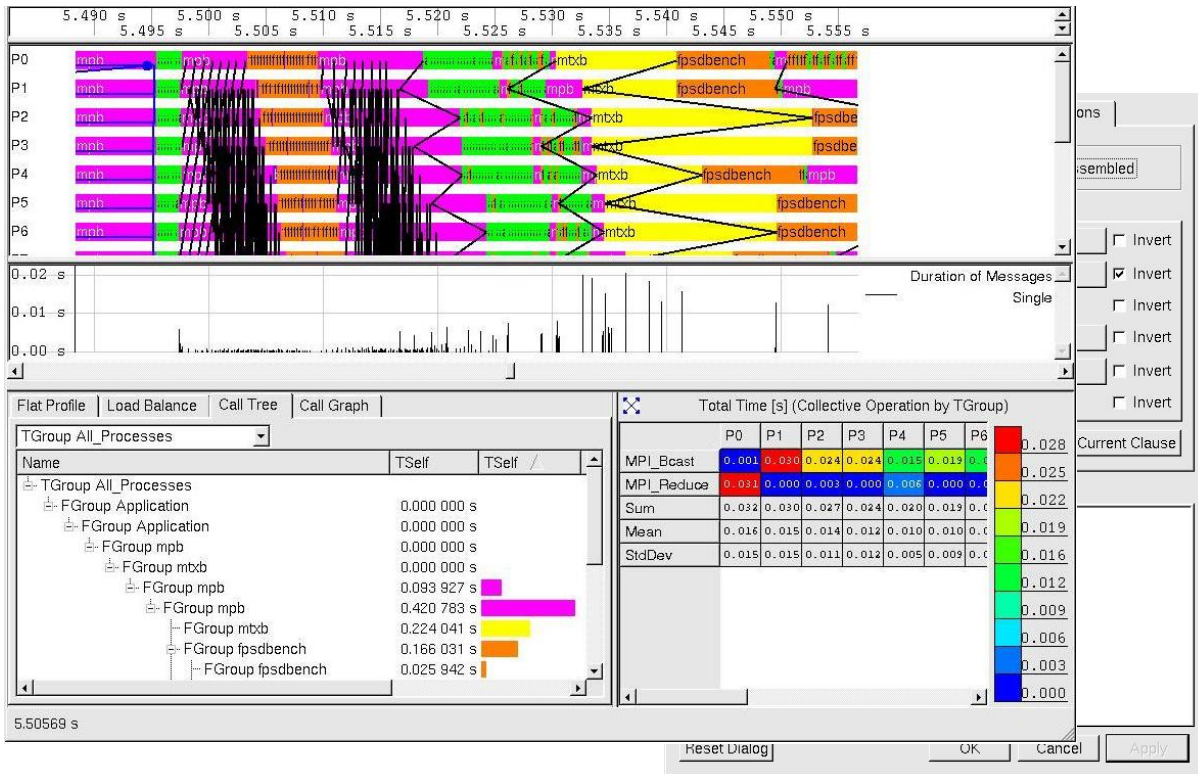


## Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.

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

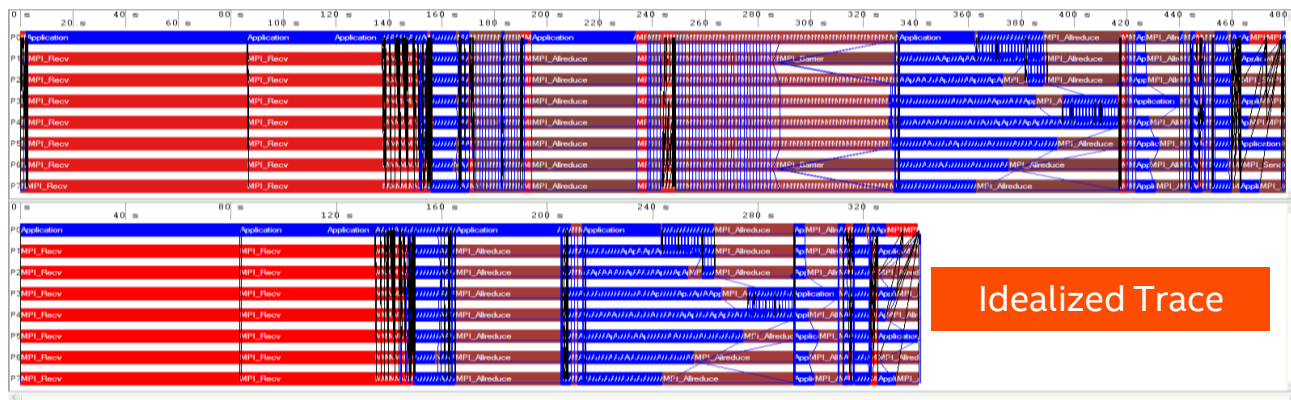
# Filtering Example



# Ideal Interconnect Simulator (Idealizer)

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

Actual trace

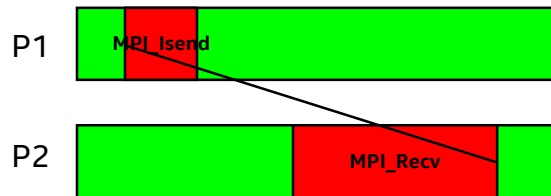


Idealized Trace

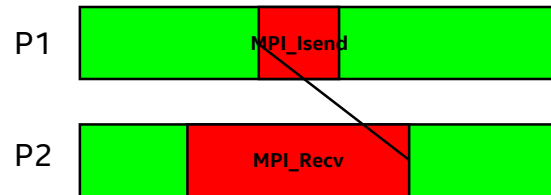
Easy way to identify application bottlenecks

# Building Blocks: Elementary Messages

Early Send /  
Late Receive

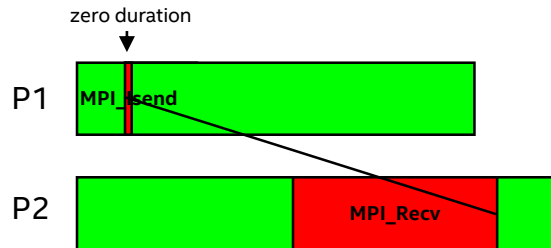


Late Send /  
Early Receive

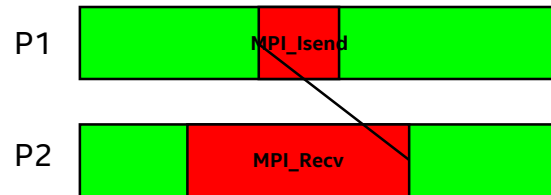


# Building Blocks: Elementary Messages

Early Send /  
Late Receive

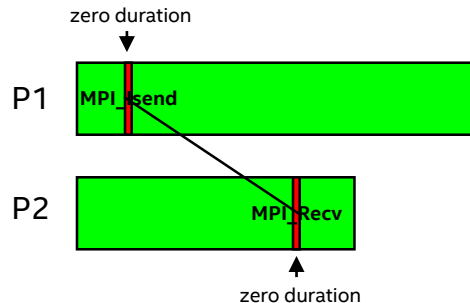


Late Send /  
Early Receive

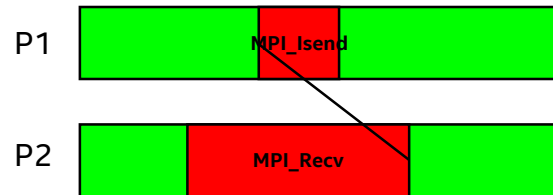


# Building Blocks: Elementary Messages

Early Send /  
Late Receive



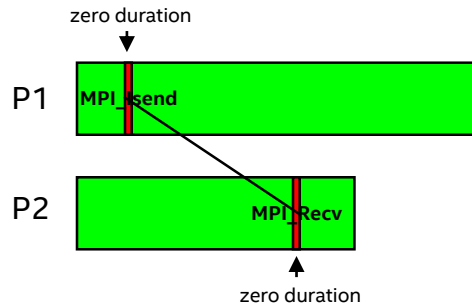
Late Send /  
Early Receive



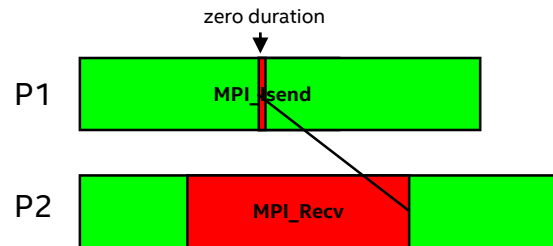


# Building Blocks: Elementary Messages

Early Send /  
Late Receive

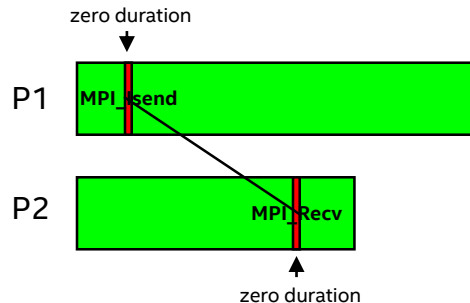


Late Send /  
Early Receive

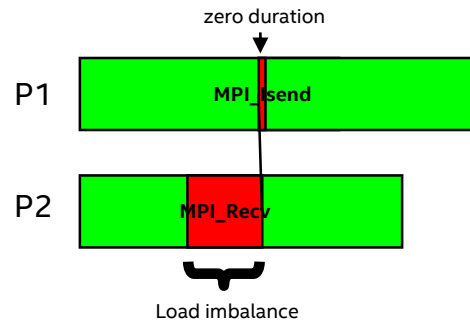


# Building Blocks: Elementary Messages

Early Send /  
Late Receive



Late Send /  
Early Receive



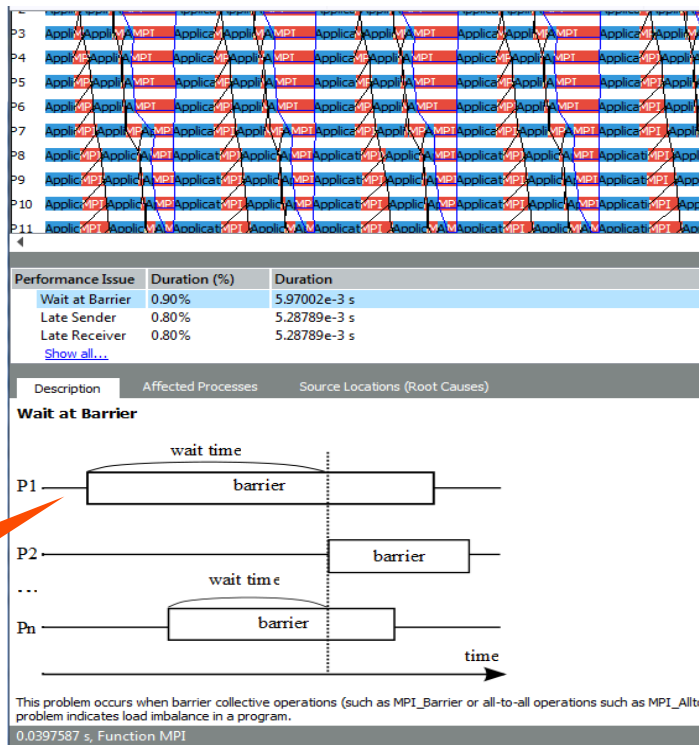
# MPI Performance Assistance

Automatic Performance Assistant

Detect common MPI performance issues

Automated tips on potential solutions

Automatically detect performance issues and their impact on runtime

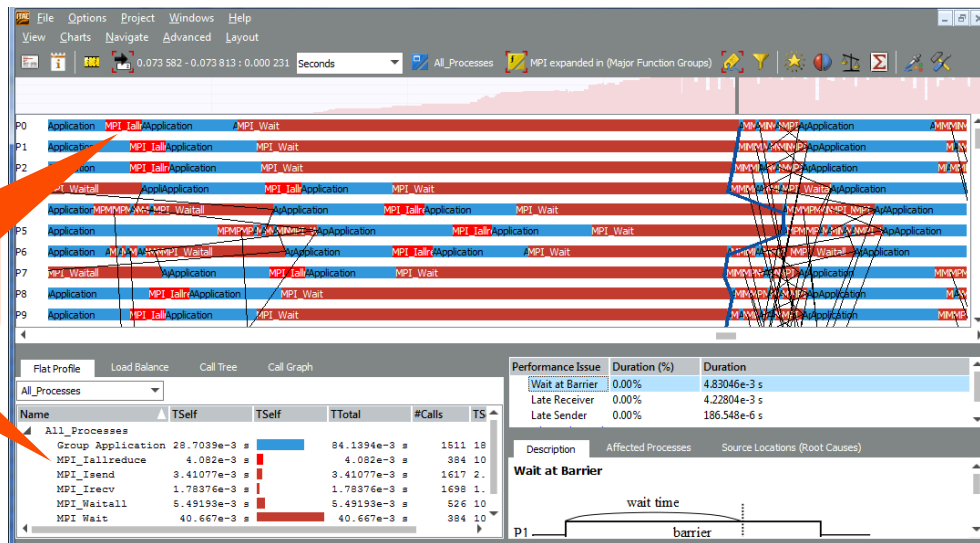


# MPI-3.0 Support

## Support for major MPI-3.0 features

- Non-blocking collectives
- Fast RMA
- Large counts

Non-blocking  
Allreduce  
(MPI\_iallreduce)



### Optimization Notice

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

**INTEL<sup>®</sup> VTUNE<sup>™</sup> AMPLIFIER XE**

# Using Intel® VTune™ Amplifier XE on MPI programs

Run VTune underneath MPI

NEW! – VTune can run multiple instances per node

- Results are grouped into one result per node
  - <result folder>.<node name>
- Within result, ranks indicate rank number

```
$ mpirun <mpi args> amplxe-cl <vtune args> -- <application and args>
```

# Easier Multi-Rank Analysis of MPI + OpenMP

Tune hybrid parallelism using ITAC + VTune Amplifier

Tune OpenMP performance of high impact ranks in VTune Amplifier

Ranks sorted by MPI Communication Spins – ranks on the critical path are on the top

Process names link to OpenMP metrics

Detailed OpenMP metrics per MPI ranks

Top OpenMP Processes by MPI Communication Spin Time								
This section lists processes sorted by MPI Communication Spin time. The lower MPI Communication Spin time is a critical path of MPI application execution. Explore OpenMP efficiency metrics by MPI processes laying on the critical path.								
Process	PID	MPI Communication Spinning <sup>①</sup> (%) <sup>②</sup>		OpenMP Potential Gain <sup>①</sup> (%) <sup>②</sup>				
<a href="#">heart_demo (rank 7)</a>	32394	5.122s	8.1%	19.929s	31.3%	2.875s	4.5%	
<a href="#">heart_demo (rank 10)</a>	32397	5.463s	8.6%	19.482s	30.6%	2.867s	4.5%	
<a href="#">heart_demo (rank 11)</a>	32398	5.593s	8.8%	20.183s	31.7%	2.873s	4.5%	
<a href="#">heart_demo (rank 6)</a>	32393	6.264s	9.8%	19.429s	30.5%	2.868s	4.5%	
<a href="#">heart_demo (rank 9)</a>	32396	6.595s	10.4%	19.379s	30.5%	2.864s	4.5%	

Per-rank OpenMP Potential Gain and Serial Time metrics

Advanced Hotspots Hotspots viewpoint (change)													
Collection Log Analysis Target Analysis Type Summary Bottom-up Caller/Callee Top-down Tree Platform													
Process / OpenMP Region / OpenMP Barrier-to-Barrier Segment / Function / Call Stack				OpenMP Potential Gain						MPI Com...	Numb... of Open...	Ins... Cou...	Avg Open...
Elapsed Time	Imbalance	Loc. Co...	Cre...	Sche...	Red...	Ato...	Other			Spin...	threads		Loop Itera...
heart_demo (rank 7)	33.938s	16.780s	0s	0s	0.086s	0s	0s	0.095s	5.122s	2.536s	18	1	
[Serial - outside any region]	2.875s							0s	2.536s				
solve\$omp\$parallel:18@unknown:527:561	31.062s	16.780s	0s	0s	0.086s	0s	0s	0.095s	2.586s	18		1	
make_rk_step\$omp\$loop_barrier_segment@unknown:352	8.887s	1.124s	0s	0s	0.012s	0s	0s	0.013s	0s	18		98	Static 1,762
update_coupling_v2\$omp\$loop_barrier_segment@unknov	16.461s	11.565s	0s	0s	0.025s	0s	0s	0.043s	0.281s	18		98	Static 1,762
make_rk_step\$omp\$barrier_segment@unknown:247	5.715s	4.090s	0s	0s	0.049s	0s	0s	0.039s	2.305s	18			
heart_demo (rank 10)	63.619s	19.108s	0s	0s	0.188s	0s	0s	0.187s	5.463s			1	
heart_demo (rank 11)	63.615s	19.824s	0s	0s	0.176s	0s	0s	0.183s	5.593s			1	
heart_demo (rank 6)	63.617s	19.091s	0s	0s	0.148s	0s	0s	0.190s	6.264s			1	
heart_demo (rank 9)	63.616s	19.002s	0s	0s	0.173s	0s	0s	0.203s	6.595s			1	

## Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.

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



# CHECKING MPI APPLICATION CORRECTNESS



# MPI Correctness Checking

## Solves two problems:

- Finding programming mistakes in the application which need to be fixed by the application developer
- Detecting errors in the execution environment

## Two aspects:

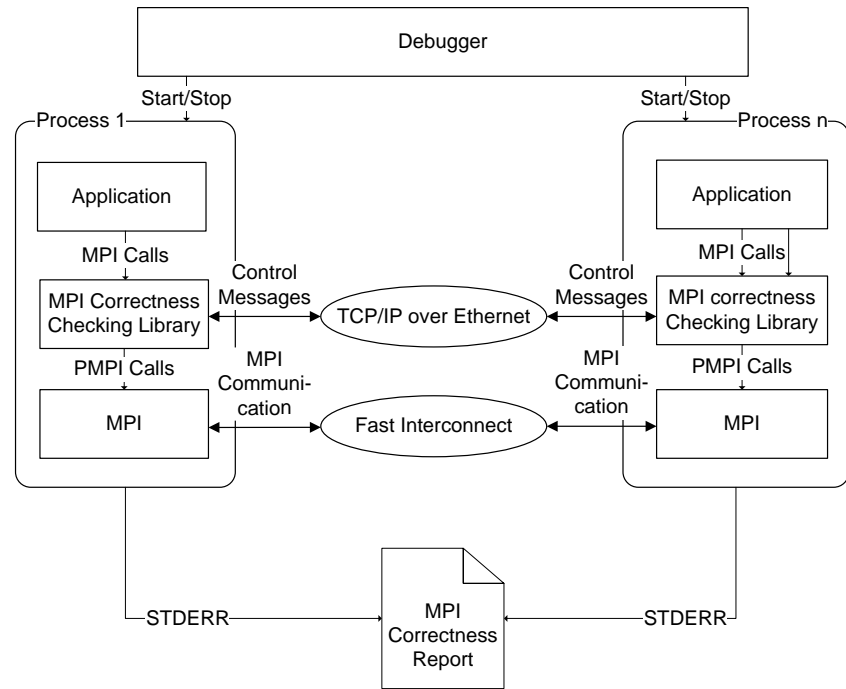
- Error Detection – done automatically by the tool
- Error Analysis – manually by the user based on:
  - Information provided about an error
  - Knowledge of source code, system, ...

# How Correctness Checking Works

All checks are done at runtime in MPI wrappers

Detected problems are reported on stderr immediately in textual format

A debugger can be used to investigate the problem at the moment when it is found



# Categories of Checks

## Local checks: isolated to single process

- Unexpected process termination
- Buffer handling
- Request and data type management
- Parameter errors found by MPI

## Global checks: all processes

- Global checks for collectives and p2p ops
  - Data type mismatches
  - Corrupted data transmission
  - Pending messages
  - Deadlocks (hard & potential)
- Global checks for collectives – one report per operation
  - Operation, size, reduction operation, root mismatch
  - Parameter error
  - Mismatched MPI\_Comm\_free()

# Severity of Checks

## Levels of severity:

- *Warnings*: application can continue
- *Error*: application can continue but almost certainly not as intended
- *Fatal error*: application must be aborted

## Some checks may find both warnings and errors

- Example: CALL\_FAILED check due to invalid parameter
  - Invalid parameter in MPI\_Send() => msg cannot be sent => *error*
  - Invalid parameter in MPI\_Request\_free() => resource leak => *warning*

# Correctness Checking on Command Line

Command line option via `-check_mpi` flag for Intel MPI Library:

```
$ mpirun -check_mpi -n 2 overlap
[...]  
[0] WARNING: LOCAL:MEMORY:OVERLAP: warning  
[0] WARNING: New send buffer overlaps with currently active send buffer at address 0x7fbffec10.  
[0] WARNING: Control over active buffer was transferred to MPI at:  
[0] WARNING: MPI_Isend(*buf=0x7fbffec10, count=4, datatype=MPI_INT, dest=0, tag=103,  
comm=COMM_SELF [0], *request=0x508980)  
[0] WARNING: overlap.c:104  
[0] WARNING: Control over new buffer is about to be transferred to MPI at:  
[0] WARNING: MPI_Isend(*buf=0x7fbffec10, count=4, datatype=MPI_INT, dest=0, tag=104,  
comm=COMM_SELF [0], *request=0x508984)  
[0] WARNING: overlap.c:105
```

# Correctness Checking in GUI

Enable correctness checking info to be added to the trace file:

- Enable VT\_CHECK\_TRACING environment variable:

```
$ mpirun -check_mpi -genv VT_CHECK_TRACING on -n 4 ./a.out
```



Errors



Warnings

## Optimization Notice

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



# Viewing Source Code

Function		Issue			
Process	Show Source	Time [s]	Type	Level	Description
P4		11.909 909	LOCAL:MPI:CALL_FAILED	warning	Null MPI_Request

*Warnings* indicate potential problems that could cause unexpected behavior (e.g., incomplete message requests, overwriting a send/receive buffer, potential deadlock, etc.).

*Errors* indicate problems that violate the MPI standard or definitely cause behavior not intended by the programmer (e.g., incomplete collectives, API errors, corrupting a send/receive buffer, deadlock, etc.).

**Source View: CCR in Process 1**

View: 1: C:/Work/development/ITA/main/Traces/mcerrorhandlingsuppre  
Chart:3: Event Timeline

Process 1

```
058         } else {  
059             MPI_Isend( &send, 1, MPI_CH  
060             MPI_Isend( &send, 1, MPI_CH  
061             MPI_Waitall( 2, reqs, statu  
062         }  
063     }  
064 }  
065  
066 MPI_Barrier( MPI_COMM_WORLD );  
067  
068 /* warning: free an invalid request */  
069 req = MPI_REQUEST_NULL;  
070 MPI_Request_free( &req );  
071  
072 MPI_Barrier( MPI_COMM_WORLD );
```

Function		Issue			
Process	Show Source	Time [s]	Type	Level	Description
P1		13.109 900	GLOBAL:MSG:DATATYPE:MISMATCH	error	Datatype signature mismatch.

# Debugger Integration

Debugger must be in control of application before error is found

A breakpoint must be set in `MessageCheckingBreakpoint()`

Documentation contains instructions for automating this process for TotalView\*, gdb, and idb.



# Intel® Inspector

## Dynamic Analysis

### Launch Intel® Inspector

- Use mpirun
- List your app as a parameter

### Results organized by MPI rank

### Review results

- Graphical user interface
- Command line report

Find errors earlier when they are less expensive to fix

## Static Analysis

Source analyzed for errors (similar to a build)

### Review results

- Graphical user interface

# Using Intel® Inspector with MPI

Use the command-line tool under the MPI run script to gather report data

```
$ mpirun -n 4 inspxe-cl -r my_result -collect mi1 -- ./test
```

Argument Sets can be used for more control

- Only collect data on certain ranks
- Different collections or options on different ranks

A unique results directory is created for each analyzed MPI rank

Launch the GUI and view the results for each rank

# **BENCHMARKING MPI AND CLUSTER PERFORMANCE**

# Intel® MPI Benchmarks

Standard benchmarks with OSI-compatible CPL license

- Enables testing of interconnects, systems, and MPI implementations
- Comprehensive set of MPI kernels that provide performance measurements for:
  - Point-to-point message-passing
  - Global data movement and computation routines
  - One-sided communications
  - File I/O
  - Supports MPI-1.x, MPI-2.x, and MPI-3.x standards

What's New:

Introduction of new benchmarks

- Measure cumulative bandwidth and message rate values

The Intel® MPI Benchmarks provide a simple and easy way to measure MPI performance on your cluster

# Online Resources

## Intel® MPI Library product page

- [www.intel.com/go/mpi](http://www.intel.com/go/mpi)

## Intel® Trace Analyzer and Collector product page

- [www.intel.com/go/traceanalyzer](http://www.intel.com/go/traceanalyzer)

## Intel® Clusters and HPC Technology forums

- <http://software.intel.com/en-us/forums/intel-clusters-and-hpc-technology>

## Intel® Xeon Phi™ Coprocessor Developer Community

- <http://software.intel.com/en-us/mic-developer>



# Legal Disclaimer & Optimization Notice

INFORMATION IN THIS DOCUMENT IS PROVIDED "AS IS". NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO THIS INFORMATION INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products.

Copyright © 2016, Intel Corporation. All rights reserved. Intel, Pentium, Xeon, Xeon Phi, Core, VTune, Cilk, and the Intel logo are trademarks of Intel Corporation in the U.S. and other countries.

## Optimization Notice

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.

Notice revision #20110804

**BACKUP**