



MPI Tuning with Intel® Trace Analyzer and Intel® Trace Collector

Intel Software College



Agenda

Introduction

Basic Usage

Focusing on the relevant

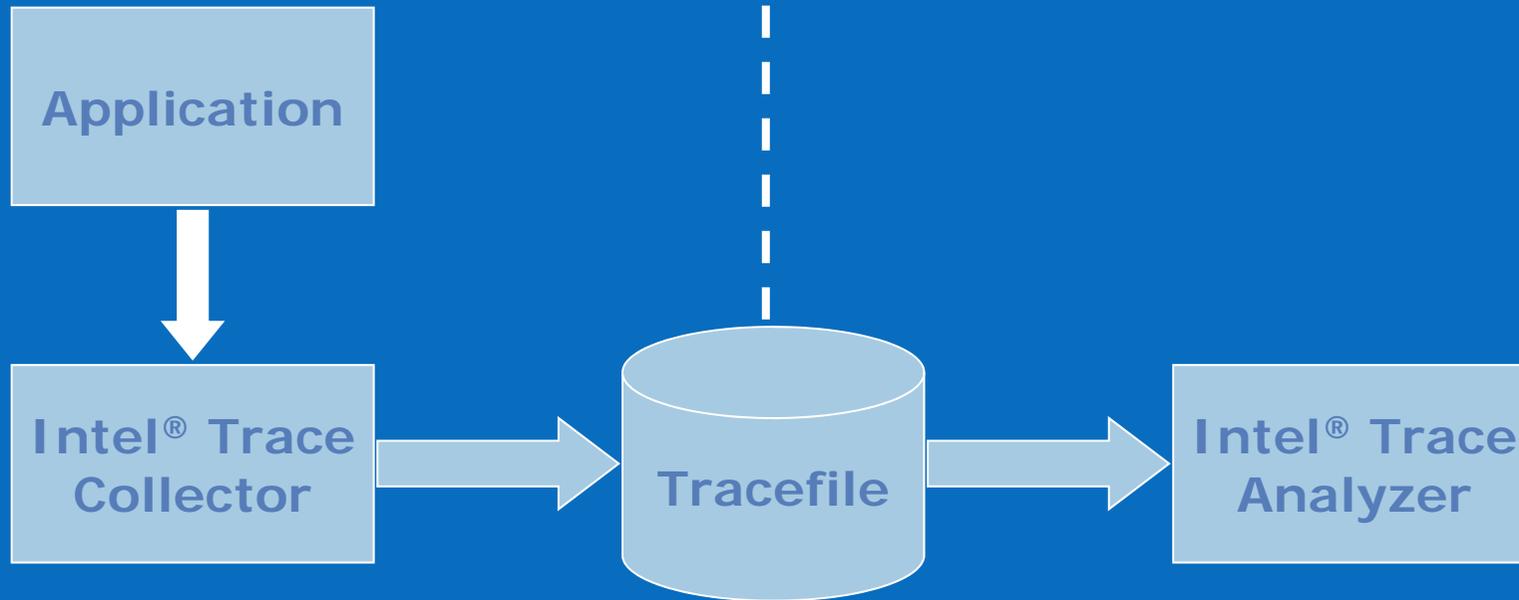
Detecting problems



Introduction – What is Tracing?

- Record program execution
 - Program events such as function enter/exit, communication
- 1:1 protocol of the actual program execution
 - Sampling gathers statistical information
- Accurate data
- Easily get loads of data

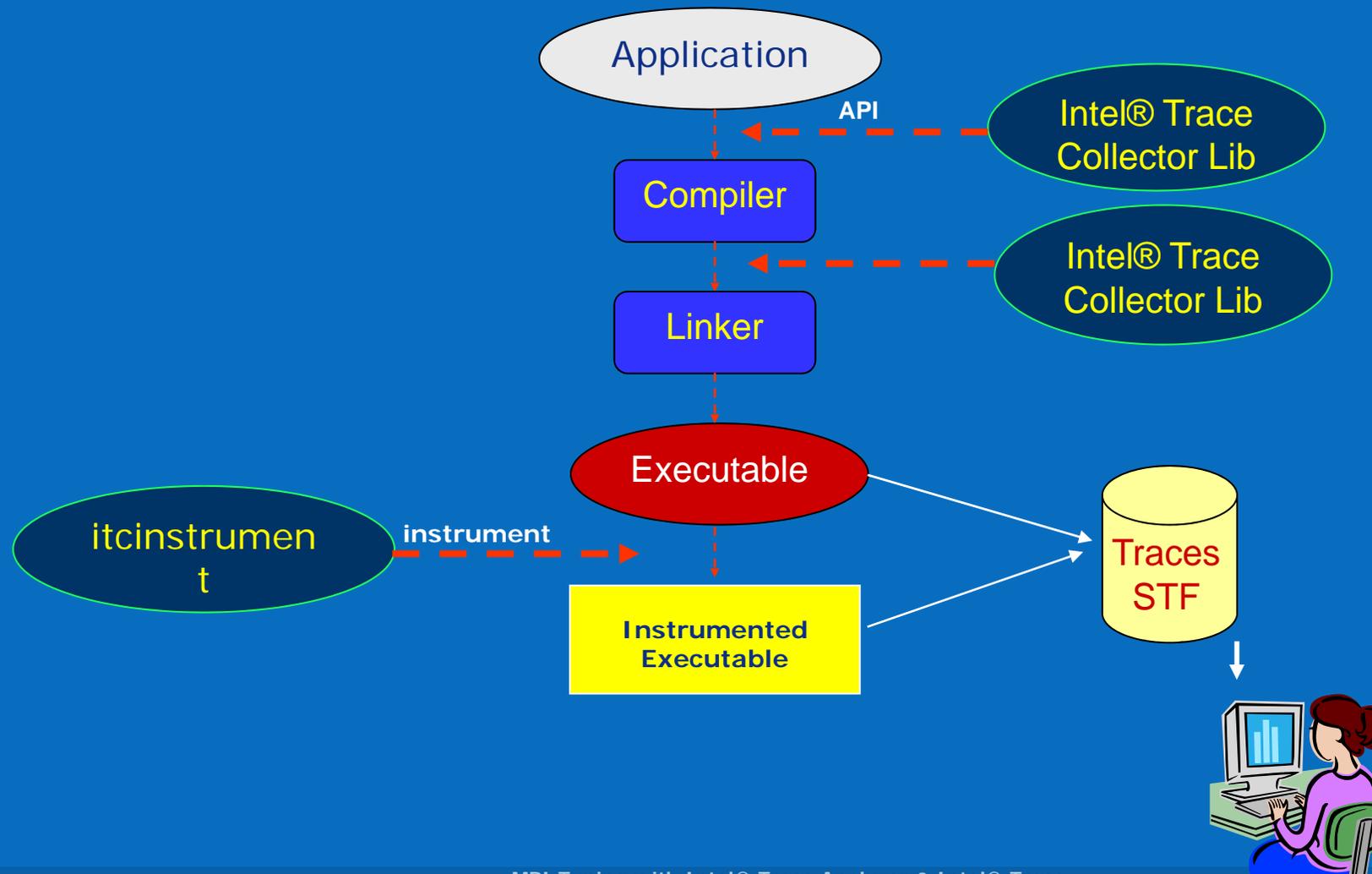
Trace Universe



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Components and Interaction



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Supported Programming Models

MPI: Message Passing

Languages: C/C++, Java*, Fortran

Multithreading: explicit threads - implicit threads (OpenMP)

Hybrid Models: MPI + OpenMP

Supported Environments

Intel® Xeon® Processor

Intel 64 Bit Xeon Processor

Intel Itanium® Processor

Intel MPI, MPICH, LAM-MPI

Linux (RHEL, RHAS, SLES)

SGI* Altix*

Microsoft* Windows* (Analyzer Only)



Introduction – Intel® Trace Collector

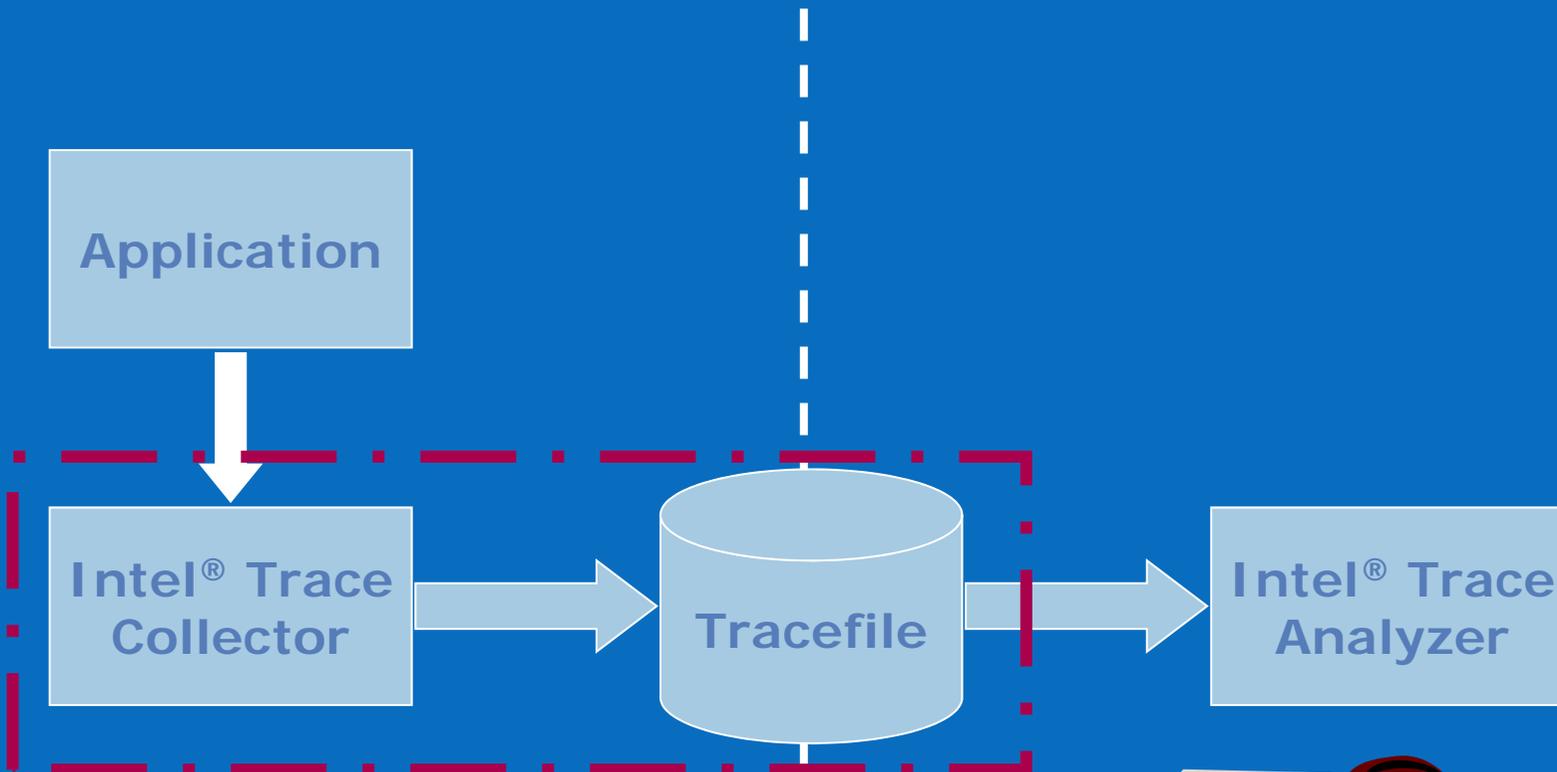
Key features

Usage

MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Trace Universe – Intel® Trace Collector



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Intel® Trace Collector - Overview

Event based approach

Low impact on application performance

Provides API to instrument user code

Trace optimized program runs

Analyzes communication layer (default)



Event based approach

Event = time stamp + thread ID + description

Function entry/exit

Messages

Collective operations

Counter samples

Strengths of Event-based Tracing

Predict detailed program behavior

Record **exact sequence** of program states – keep timing consistent

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

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

Key Features

Catch all MPI events

Strong configuration mechanism

- Filters, settings, features

Automatic source-code references

Instrumentation

- Rich API
- Binary instrumentation (itcinstrument)
- Compiler based (beta)

Fail-safe version



How to use Intel® Trace Collector

Step 1: Generate instrumented binary re-link your application:

```
mpicc -shlib tst.o $(LFLAGS) -lVT -lmpi -ldwarf -lelf -lnsl  
-lm -o tst
```

Step 2: Produce tracefile

run the instrumented binary for a representative amount of time
(reduce initialization influences) on representative data (no corner
cases)

Step 3: (optional): Instrument binary

```
itcinstrument -input lm -output lm_instr
```

Lab 1 - Set up Environment



Get hands-on package

Edit make_base to reflect local configuration

SETUP ENV

Get hands-on package

Package name:

```
/home/public/itac-bcor2005.tgz
```

Contents:

- make_base
- Makefile
 - task01
 - **README**
 - **Description**
 - **source files**
 - task02
 - task03
 - ...



Edit make_base to reflect local configuration

CC = [path to MPI compiler wrapper]
MPI_ROOT = [path to MPI implementation]
VT_ROOT = [path to Intel® Trace Collector]
RUN_CMD = [command to run mpi exec's]

Example:

```
CC = mpiicc  
MPI_ROOT = /opt/intel/mpi/2.0  
VT_ROOT = /opt/intel/itc6  
RUN_CMD = /opt/intel/mpi/2.0/bin64/mpiexec -n
```

SETUP ENV

Intel® Trace Collector installation is in [ITC]: `/opt/intel/itc6`

- `csch:` `source [ITC]/sourceme.csh`
- `sh:` `sh [ITC]/sourceme.sh`

Intel® Trace Analyzer installation is in [ITA]: `/opt/intel/ita6`

- `csch:`
`set path = (. [ITA]/bin [ITC]/bin $path)`
- `bash:`
`export PATH=$PATH:[ITA]/bin:[ITC]/itc6/bin`
- `start:` `traceanalyzer`

Lab 2 – Trace File Generation

Objective: How to generate a tracefile.

This directory contains **cpi.c** which is a simple example program.

Re-link (or compile and link) with Intel® Trace Collector library.

Run again – check for console message on 'writing tracefile'.

Open tracefile with Intel® Trace Analyzer.

```
[traceanalyzer cpi_itc.stf]
```

Lab 3 – Trace File Generation

MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Introduction – Intel® Trace Analyzer

Key features

Usage

MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



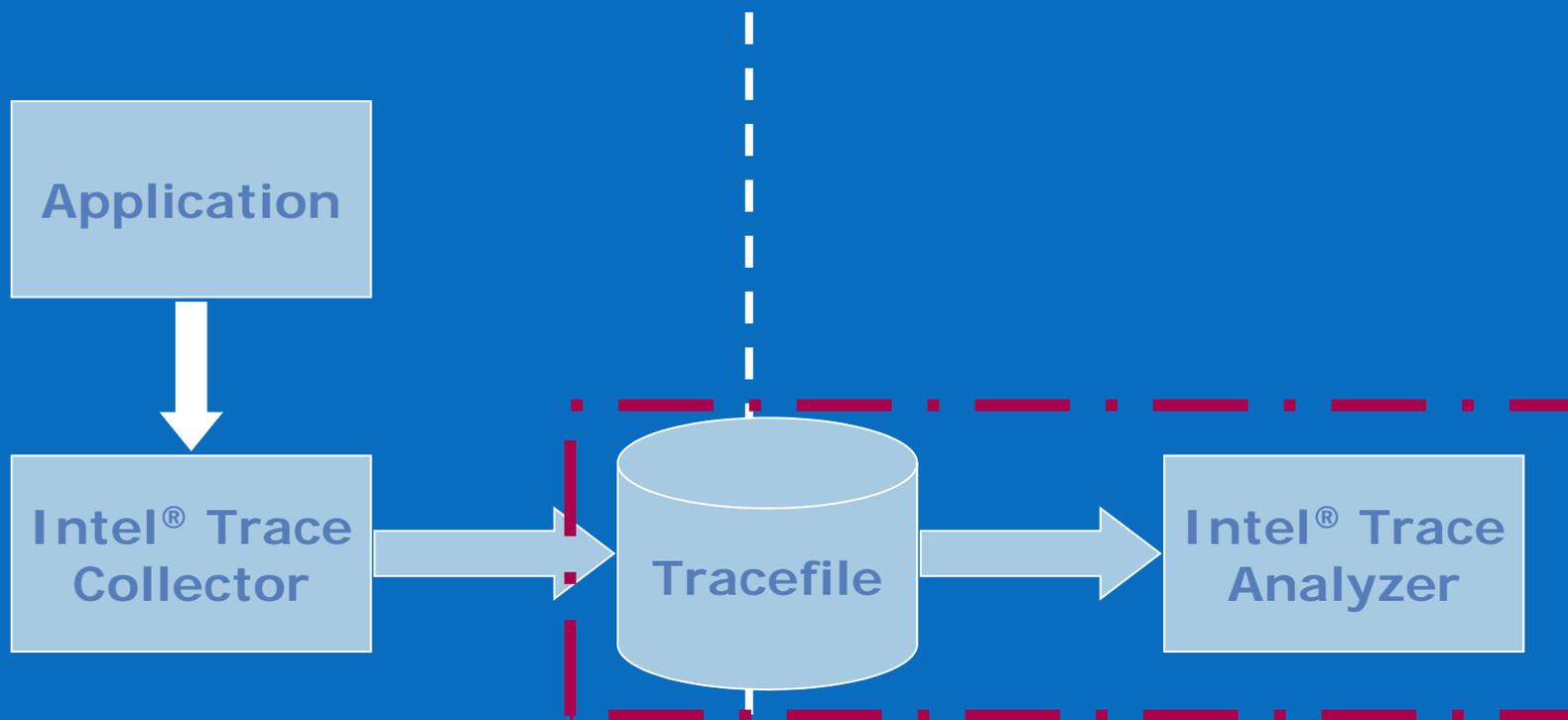
21

Copyright © 2006, Intel Corporation. All rights reserved.

Intel and the Intel logo are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States or other countries. *Other brands and names are the property of their respective owners.



Trace Universe - Intel® Trace Analyzer



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Intel® Trace Analyzer - Overview

Enables the user to quickly focus at the appropriate level of detail Offline trace analysis (**postmortem**)

Use of **hierarchical techniques** to address **scalability** in function, time and processor–space

High–performance graphics, excellent zooming and filtering

Views and Charts - I

What's a Chart?

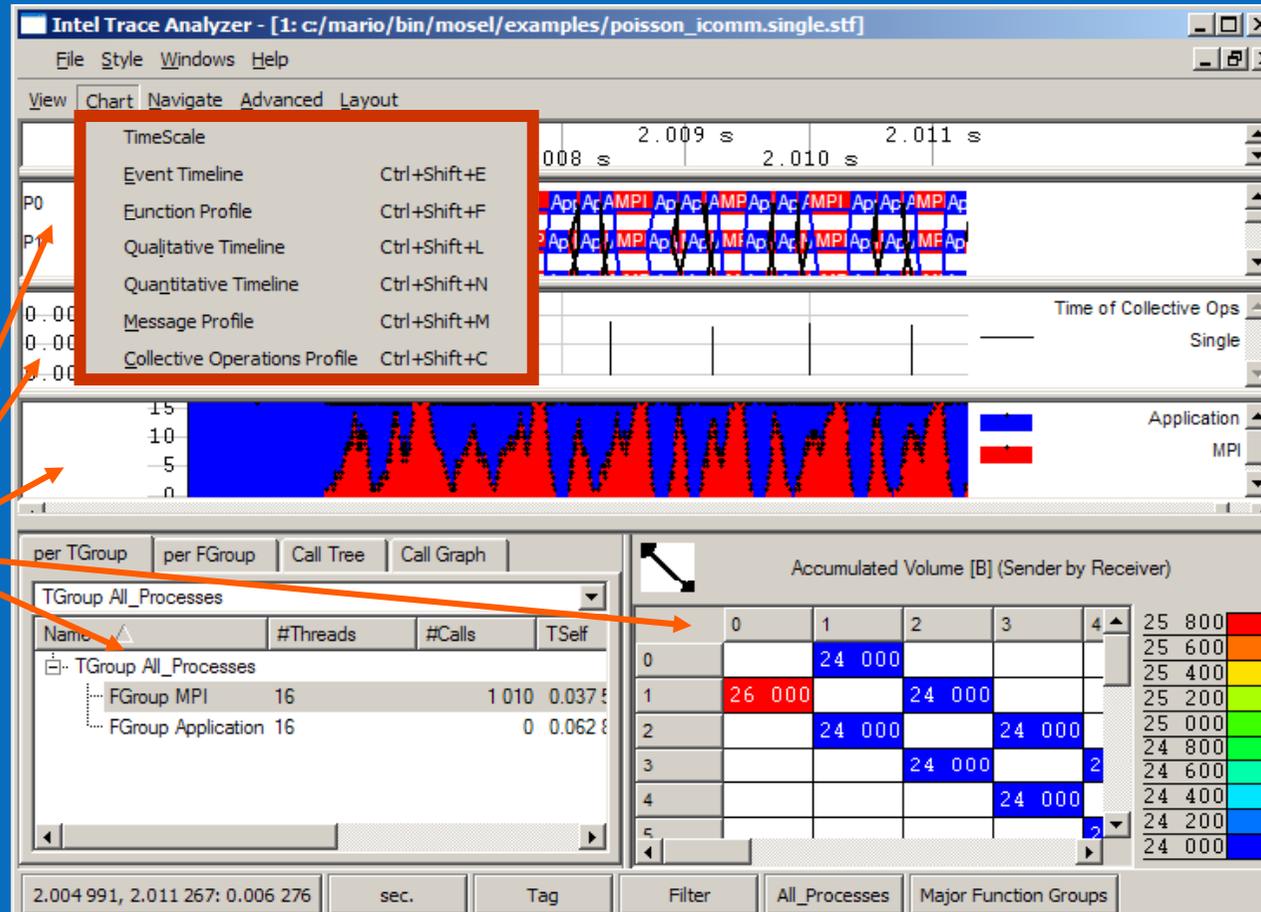
What Charts are provided?

What's a View?

What's the purpose of a View?

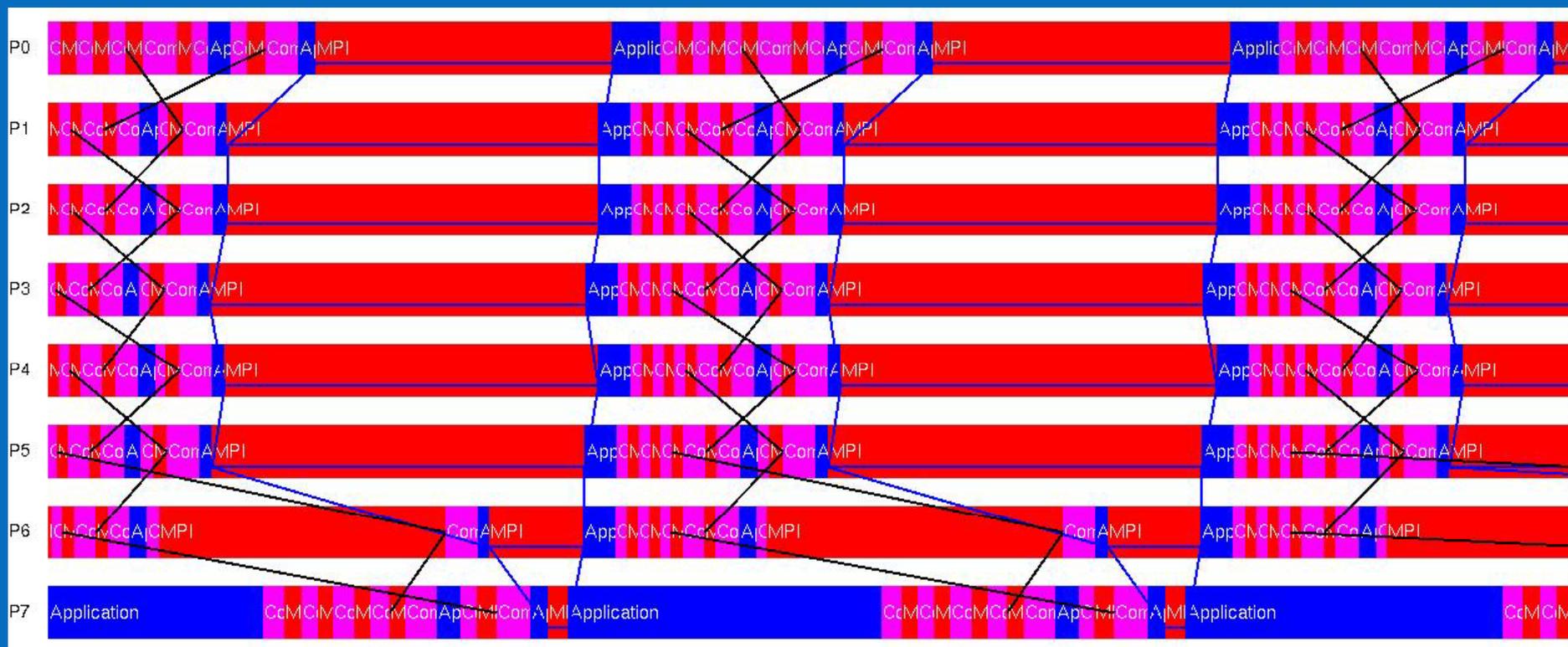
Chart

Chart



A **Chart** is a numerical or graphical diagram

Timelines: Event Timeline

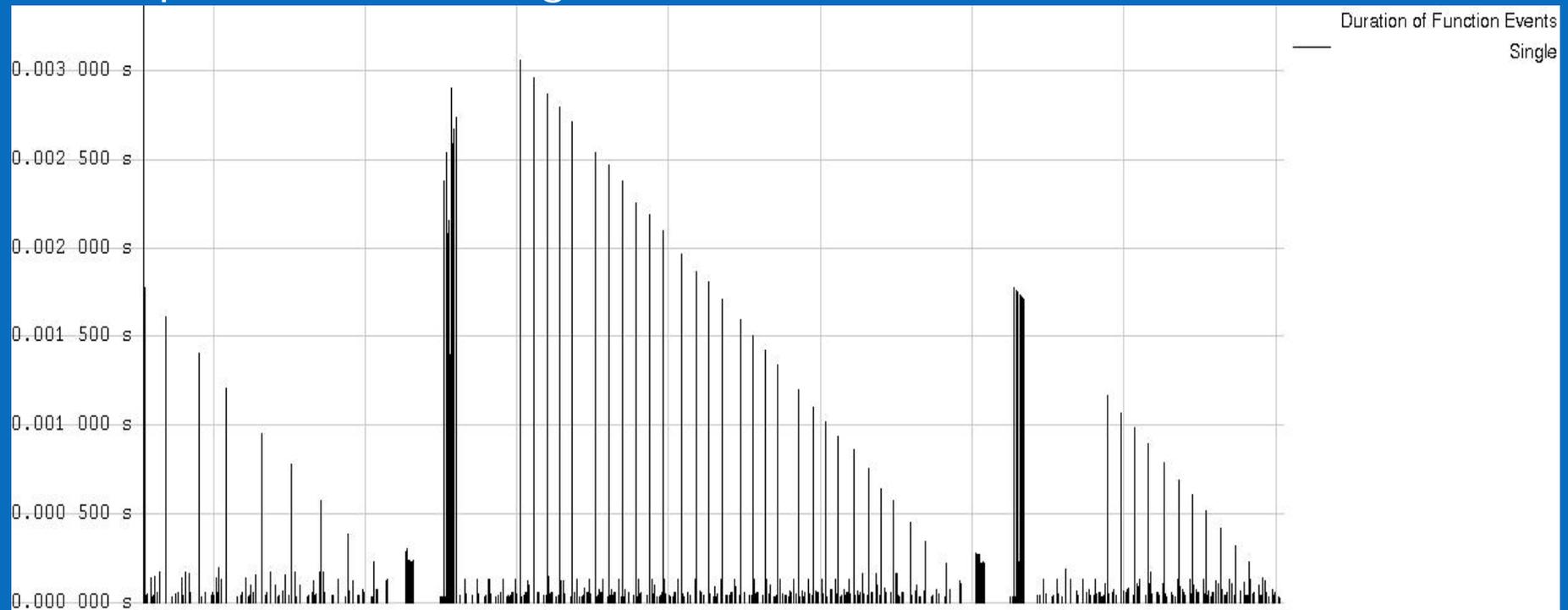


MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Timelines: Qualitative Timeline

Find patterns and irregularities



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



27

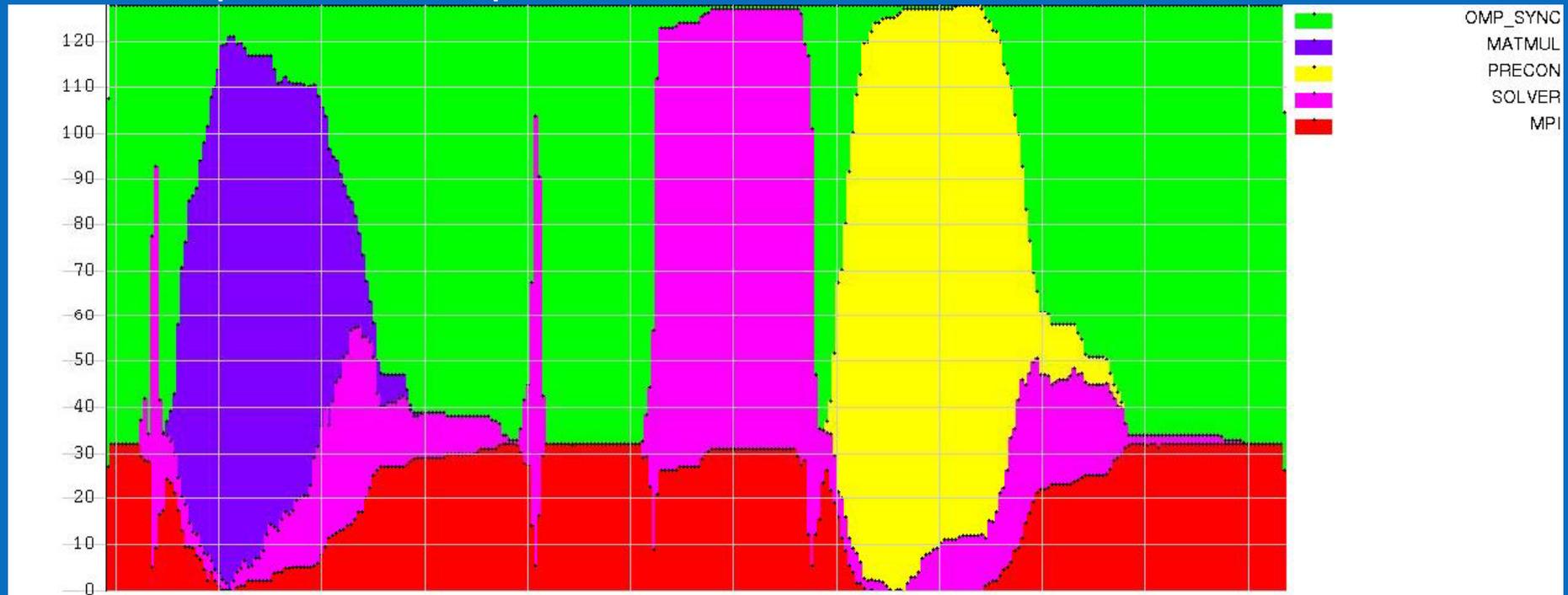
Copyright © 2006, Intel Corporation. All rights reserved.

Intel and the Intel logo are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States or other countries. *Other brands and names are the property of their respective owners.



Timelines: Quantitative Timeline

Get impression on parallelism and load balance



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



28

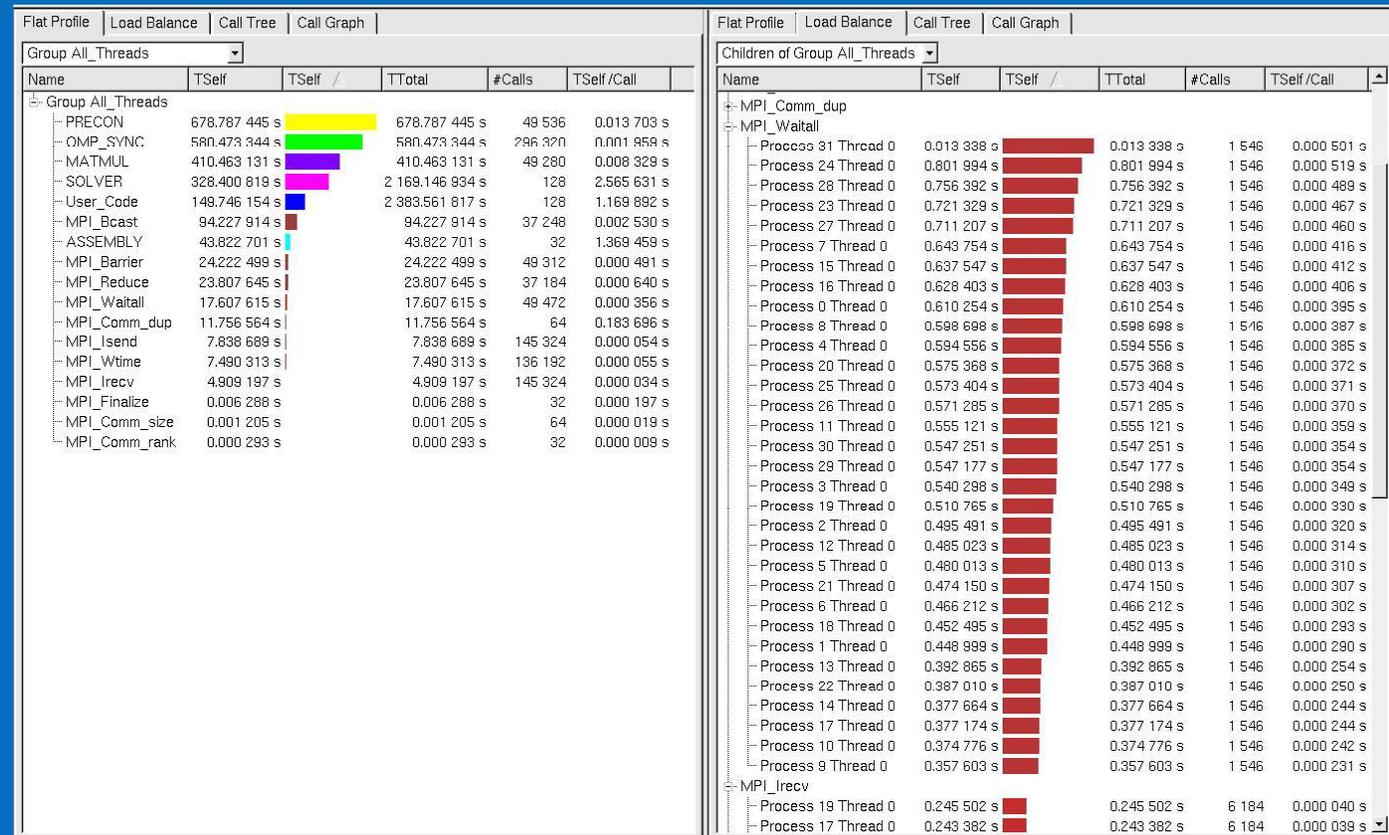
Copyright © 2006, Intel Corporation. All rights reserved.

Intel and the Intel logo are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States or other countries. *Other brands and names are the property of their respective owners.



Profiles: Flat Function Profile

Statistics about functions

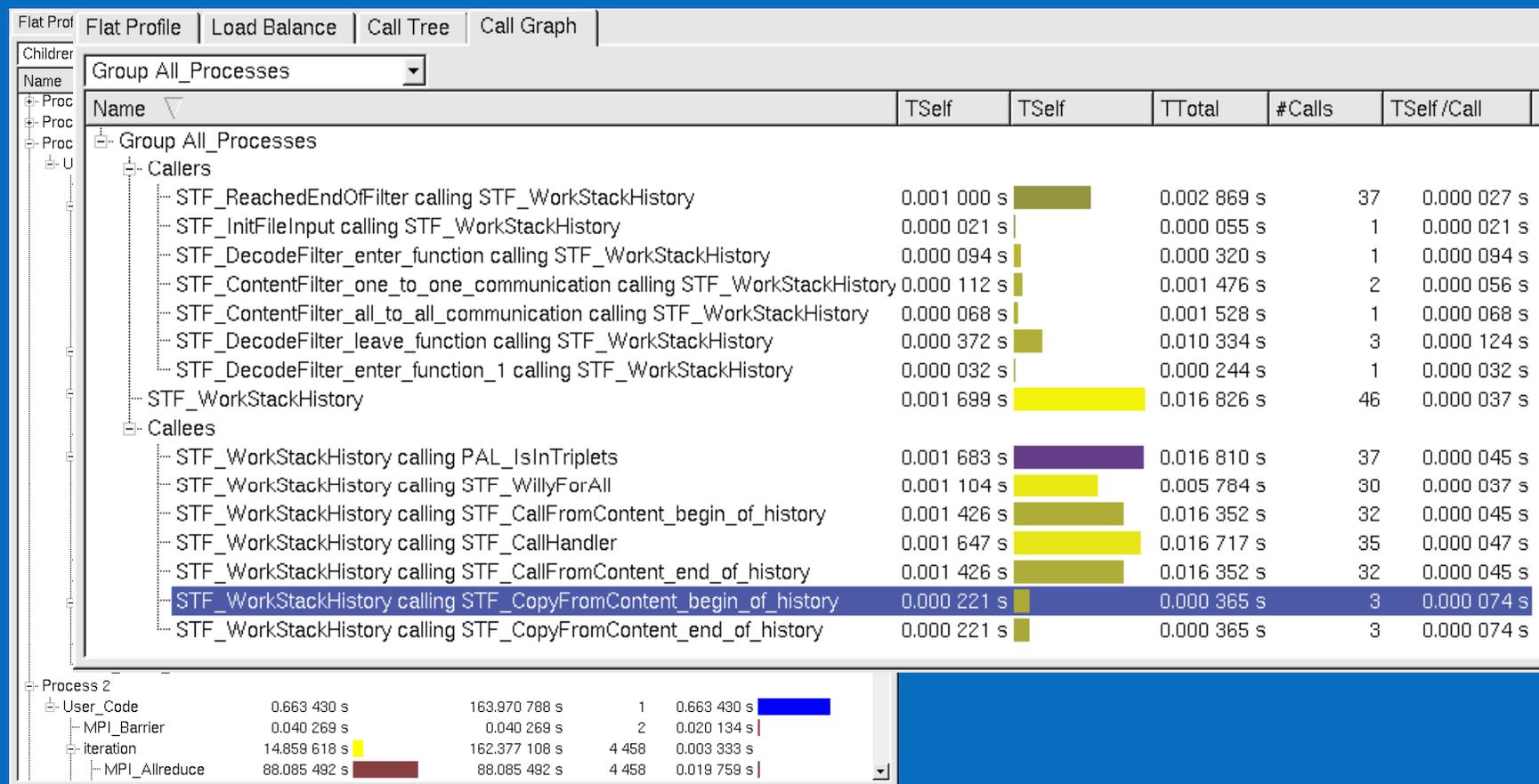


MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Profiles: Call-Tree and Call-Graph

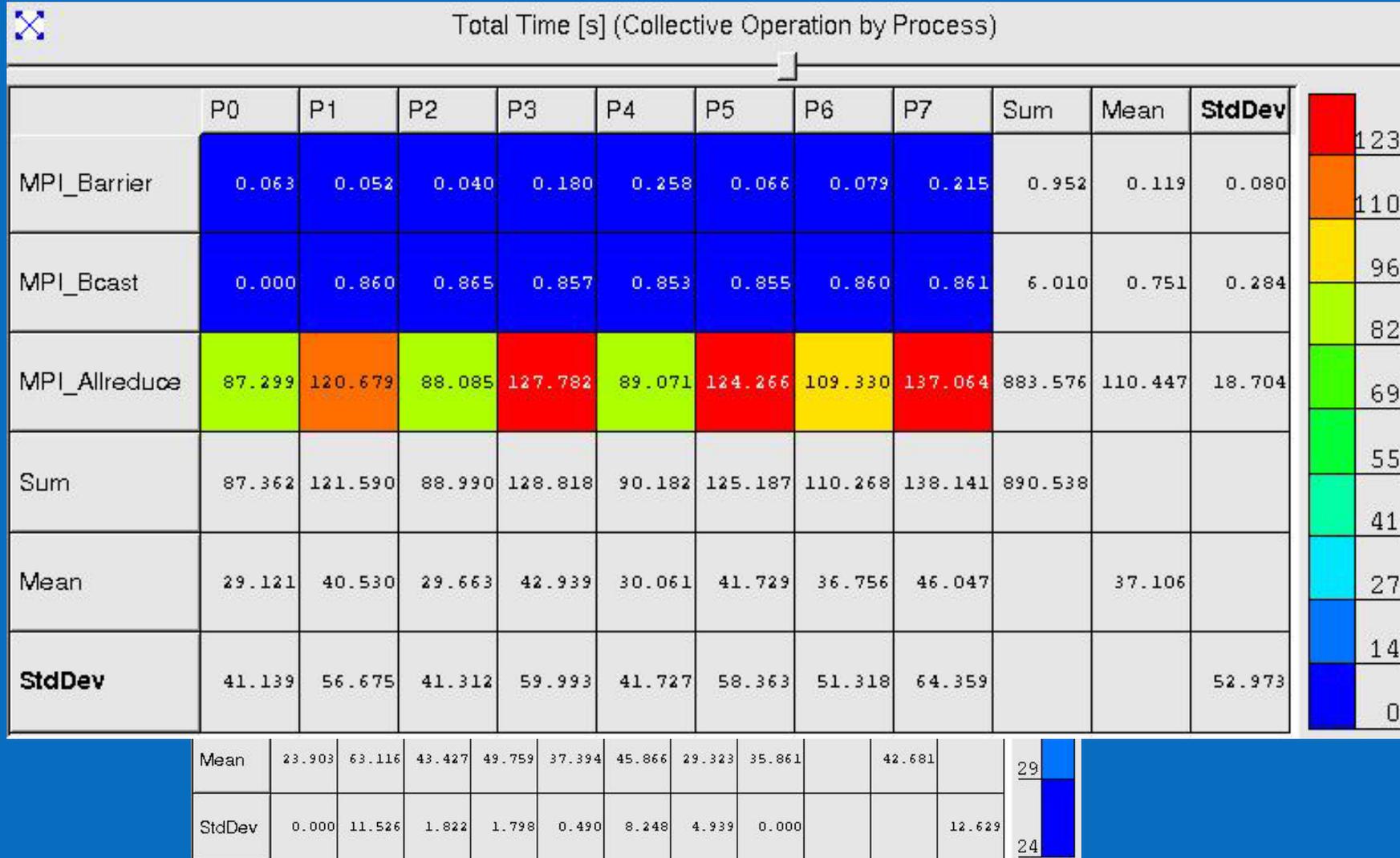
Function statistics including calling hierarchy



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Communication Profiles



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



View

Helps navigating through the trace data and keep orientation

Every **View** can contain several **Charts**

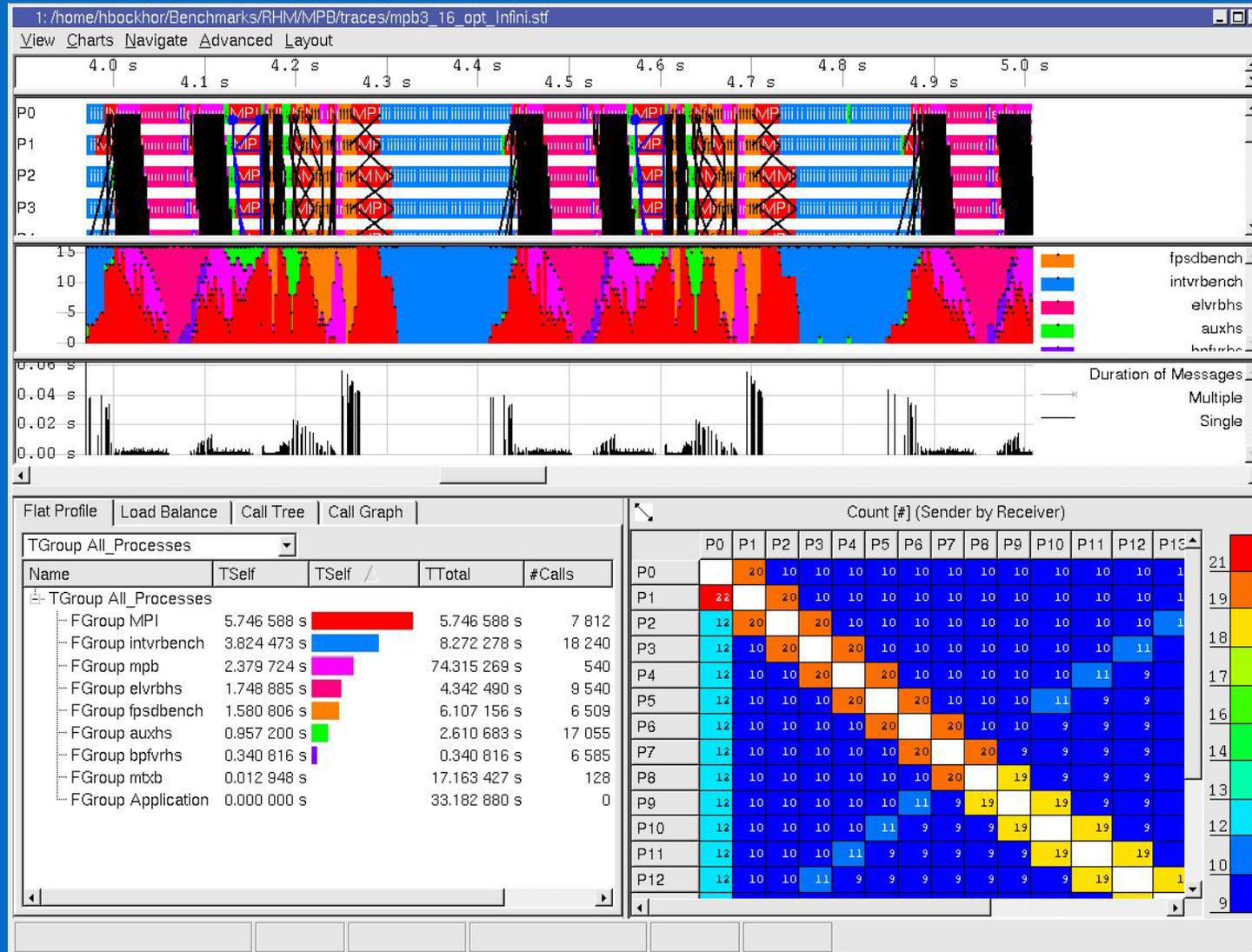
A View on a file is defined by a triplet of

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

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

Timelines are correctly aligned along time

View - zooming



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Flexibility of Views

Several Views can be opened (on the same or on different files)

Location, orientation and size of charts can easily be changed

Entire Views can and individual charts can be cloned and closed

Individual charts can be cloned in own View

Lab 4 – Views and Charts



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Charting and Viewing

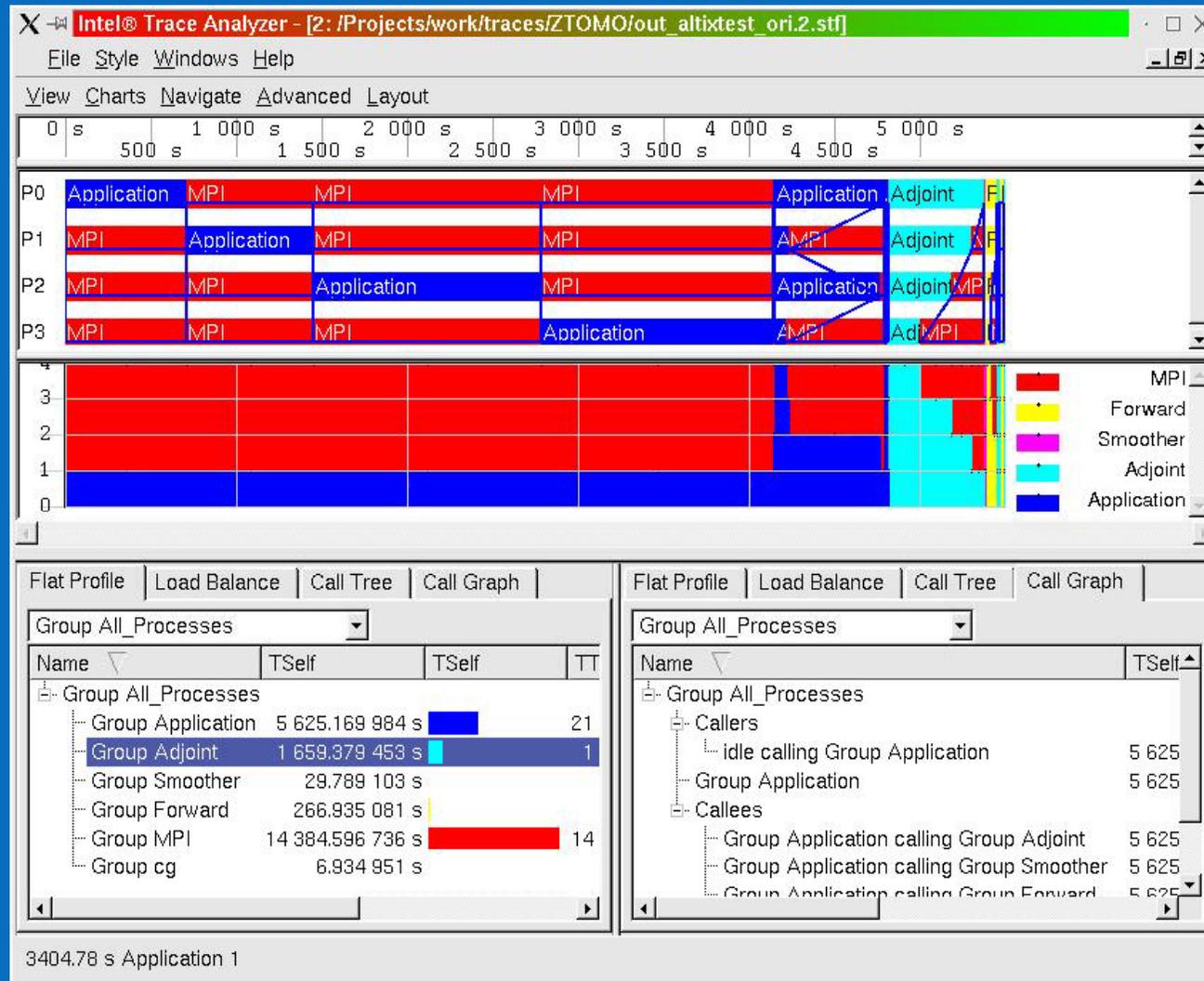
Start with

```
traces/ztomo/out_altixtest_ori.2.stf
```

Use Chart menu to open

- EventTimeline
- Quantitative Timeline
- Function Profile (Flat Profile)
- Function Profile (CallGraph)

Lab 4 - Results



Charting and Viewing

Start with

```
traces/ztomo/out_altixtest_ori.2.stf
```

Use Chart menu to open

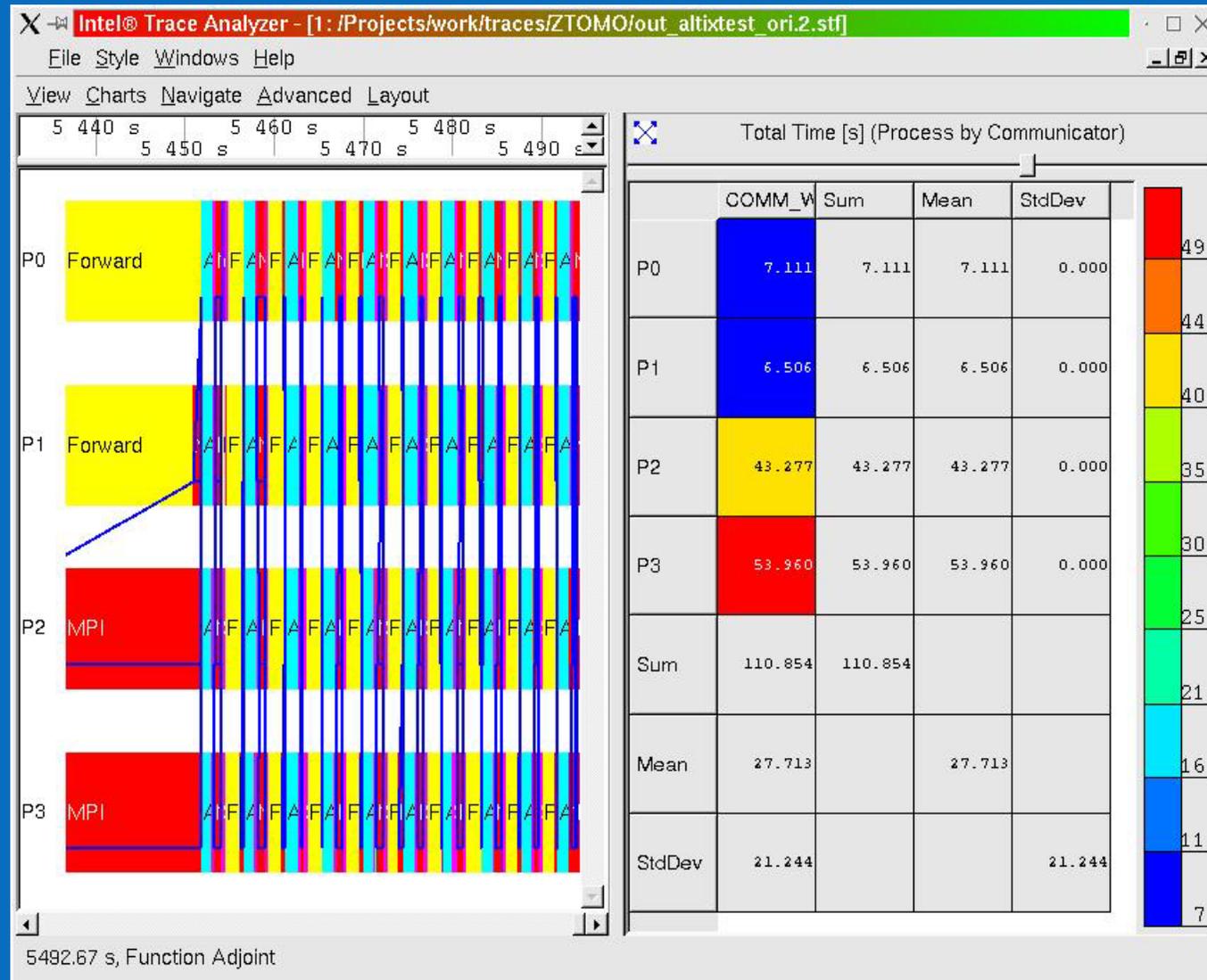
- EventTimeline (zoom to very end (left))
- Collective Operation Profile

Use Layout menu to get time line to left

Use Context menu in CollOpChart

- Columns to show Communicator
- Rows to show Processes

Lab 4b - Results



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



End of Module

MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Grouping and Aggregation

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

Functions and threads can be grouped hierarchically

- **Function Groups and Thread Groups**

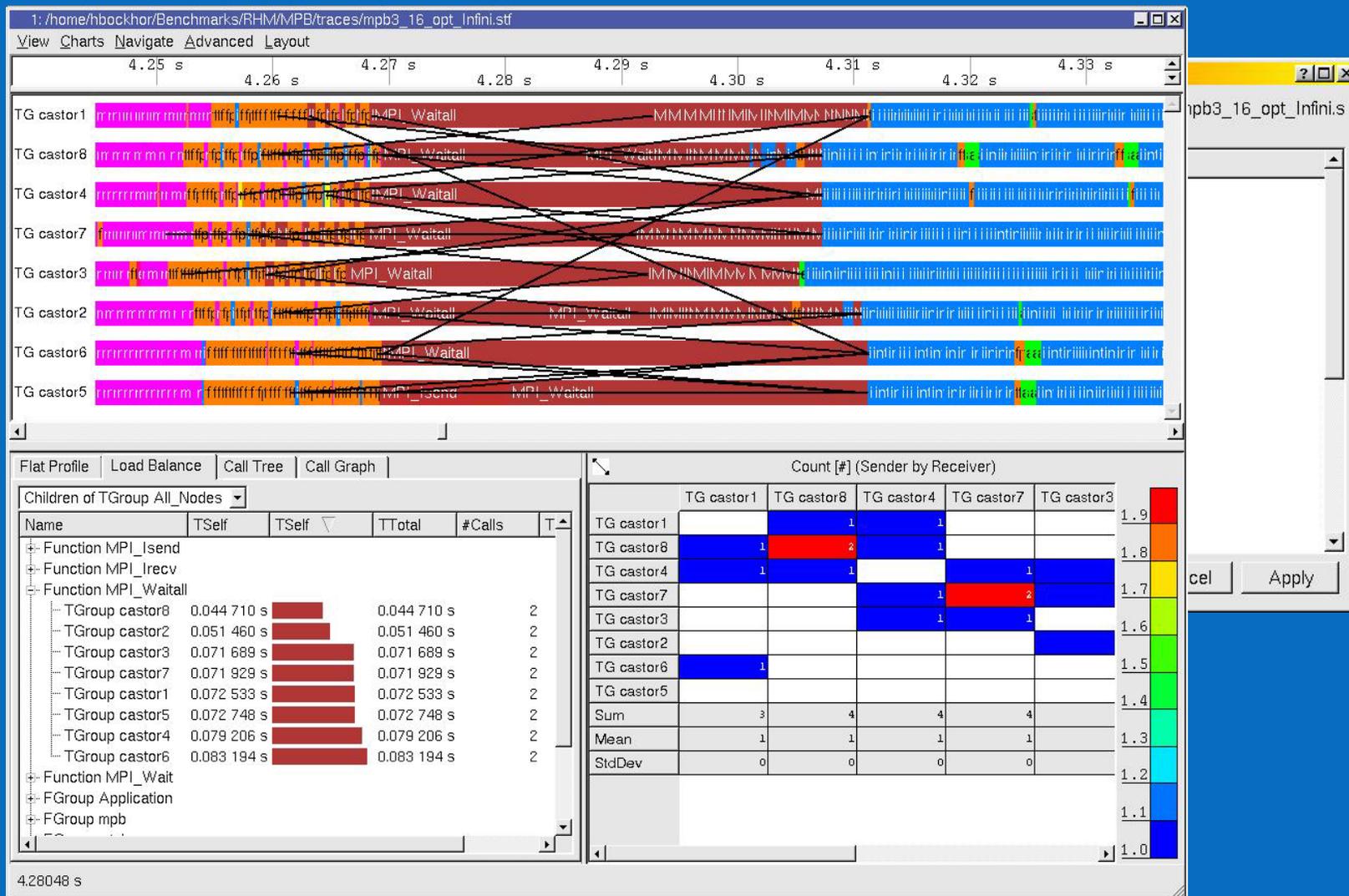
Arbitrary nesting is supported

- 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



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Lab 5 - Aggregation



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Two ways to determine time spent not in MPI for instrumented codes

Start with

```
traces/ztomo/out_altixtest_ori.2.stf
```

Quick

- Select MPI for Function Aggregation

MPI <-> non-MPI

- Create new group
- Drag & Drop MPI
- Select new group for Function Aggregation

Tagging & Filtering

Help concentrating on relevant parts

Avoid getting lost in huge amounts of trace data

Define a set of interesting data

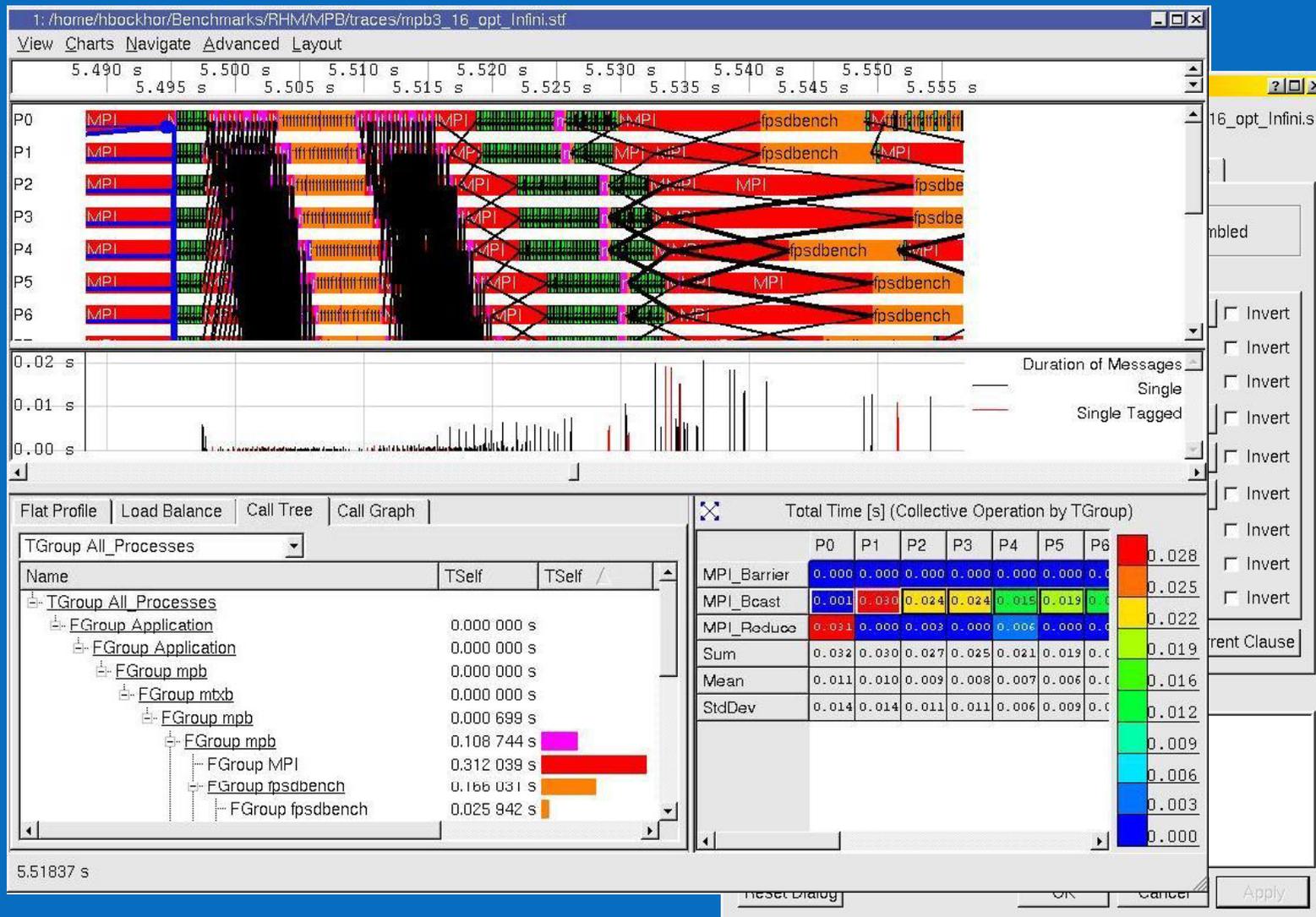
- E.g. all occurrences of function x
- E.g. all messages with tag y on communicator z

Combine several filters:
Intersection, Union, Complement

Apply it

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

Tagging Example



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Filtering Example

The screenshot shows the Intel Trace Analyzer interface with a filtered MPI trace. The main window displays a Gantt chart for processes P0-P6, showing various MPI operations such as mpb, mtbx, and fpsdbench. A 'Duration of Messages' plot is visible below the Gantt chart. The bottom panel shows a 'Flat Profile' tree and a 'Total Time [s] (Collective Operation by TGroup)' table.

	P0	P1	P2	P3	P4	P5	P6	
MPI_Bcast	0.001	0.030	0.024	0.024	0.015	0.019	0.028	0.028
MPI_Reduce	0.031	0.000	0.003	0.000	0.006	0.000	0.025	0.025
Sum	0.032	0.030	0.027	0.024	0.020	0.019	0.022	0.022
Mean	0.016	0.015	0.014	0.012	0.010	0.010	0.019	0.019
StdDev	0.015	0.015	0.011	0.012	0.005	0.009	0.016	0.016
								0.012
								0.009
								0.006
								0.003
								0.000

MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Lab 6 – Tagging and Filtering



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Find threads performing a task

Start with

```
traces/ita/ita_fg1_threaded.stf
```

Search for `FuncProfiler::*`

Load Balance & Event Timeline

Tagging

- Advanced|Tagging
 - Select `Function` tab and `Custom`
 - Type `FuncProfiler` and press `ok`

Suppress GUI idle time

Start with

```
traces/ita/ita_fg1_threaded.stf
```

Assumption: GUI is idle time

Use Filter (Function, "GUI", invert)

Detecting Problems

Link to material

Verify important hypotheses

Slide demo

- Live demo



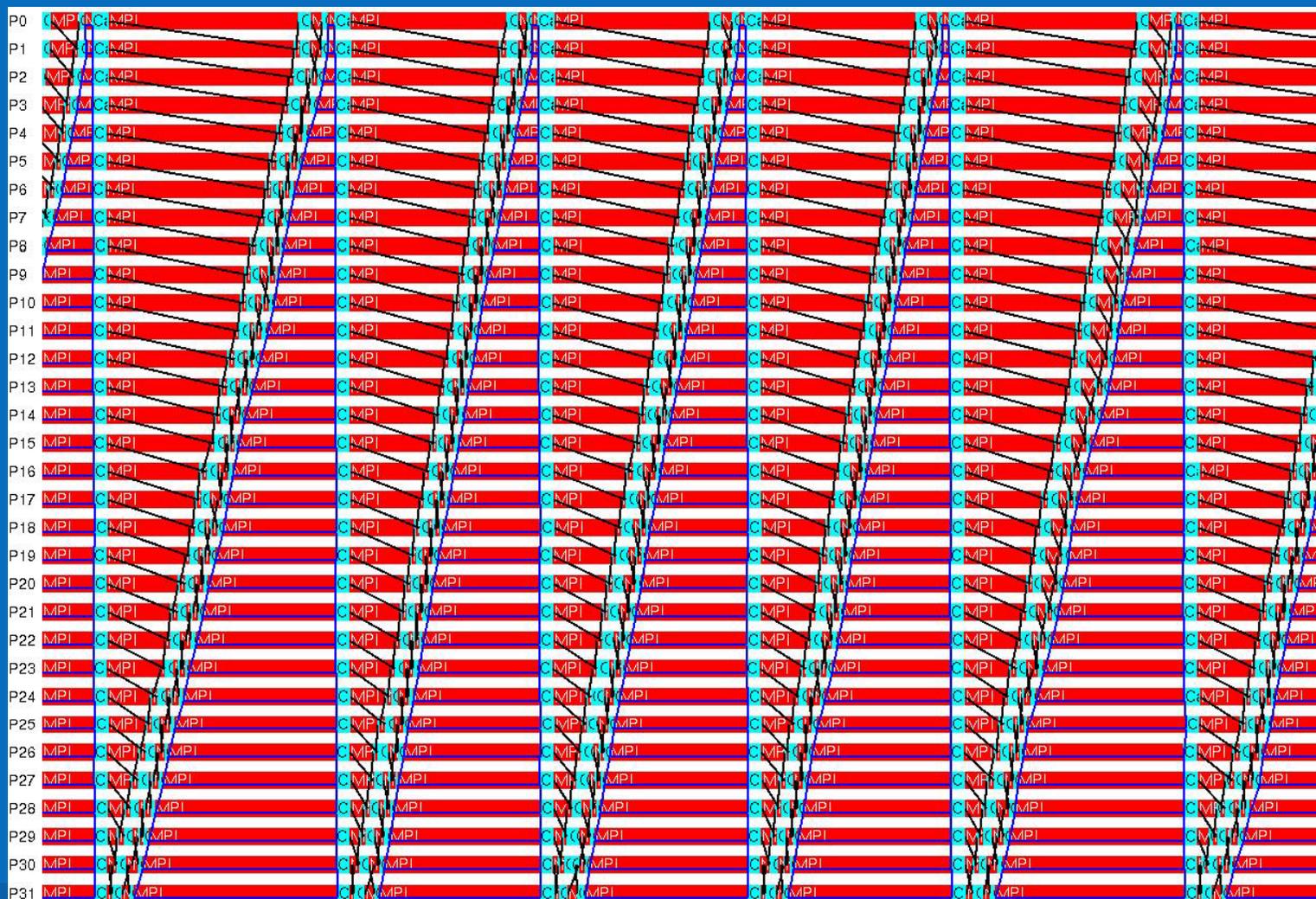
Understanding your code

Parallel Poisson Solver

Example of intuitive parallelization with disadvantageous communication pattern



Partial Trace of Poisson Solver



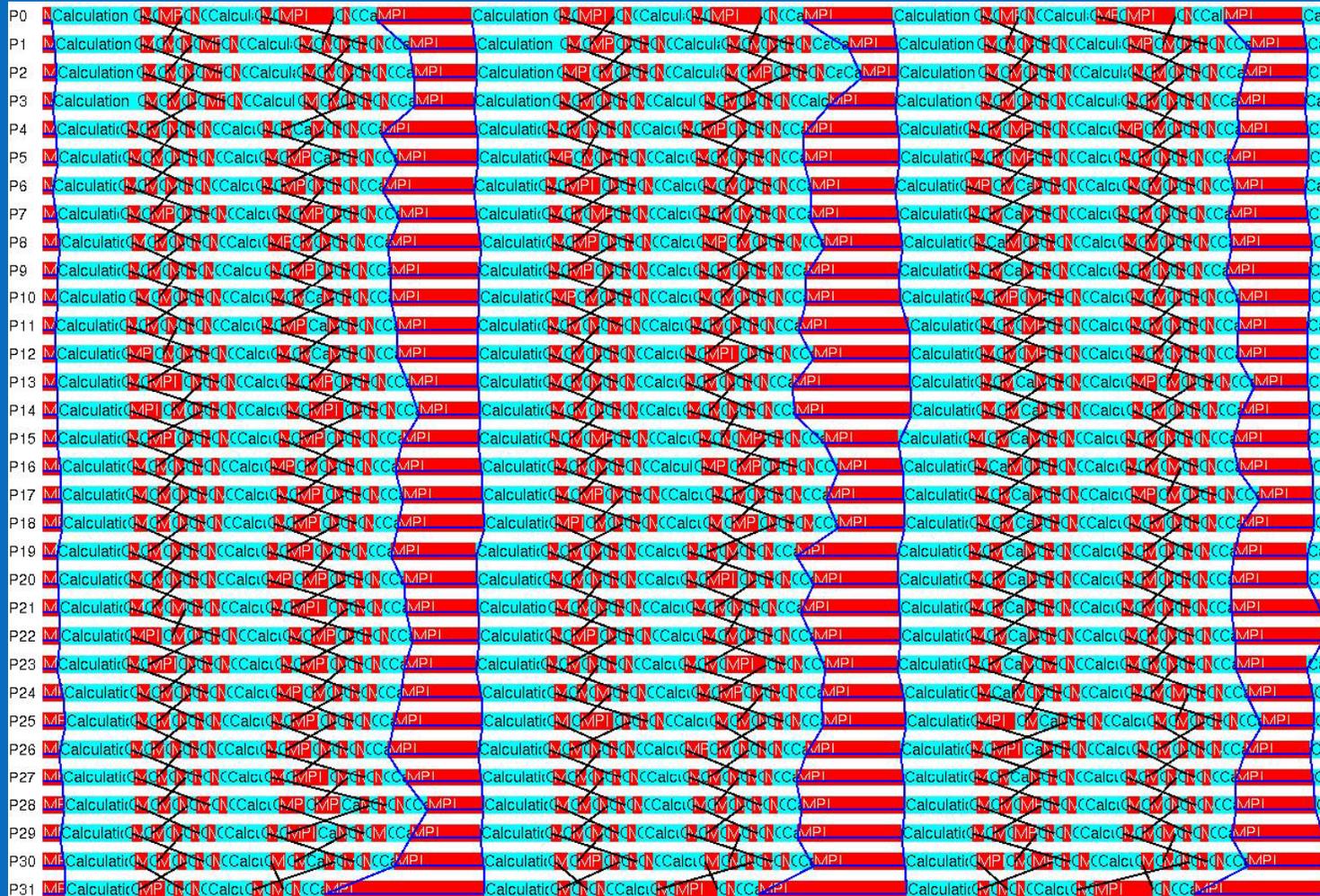
Understanding the problem

Blocking border exchange

P_n has blocks until communication between P_{n+1} and P_{n+2} was completed

Solution: Non blocking communication

Non-Blocking Version



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



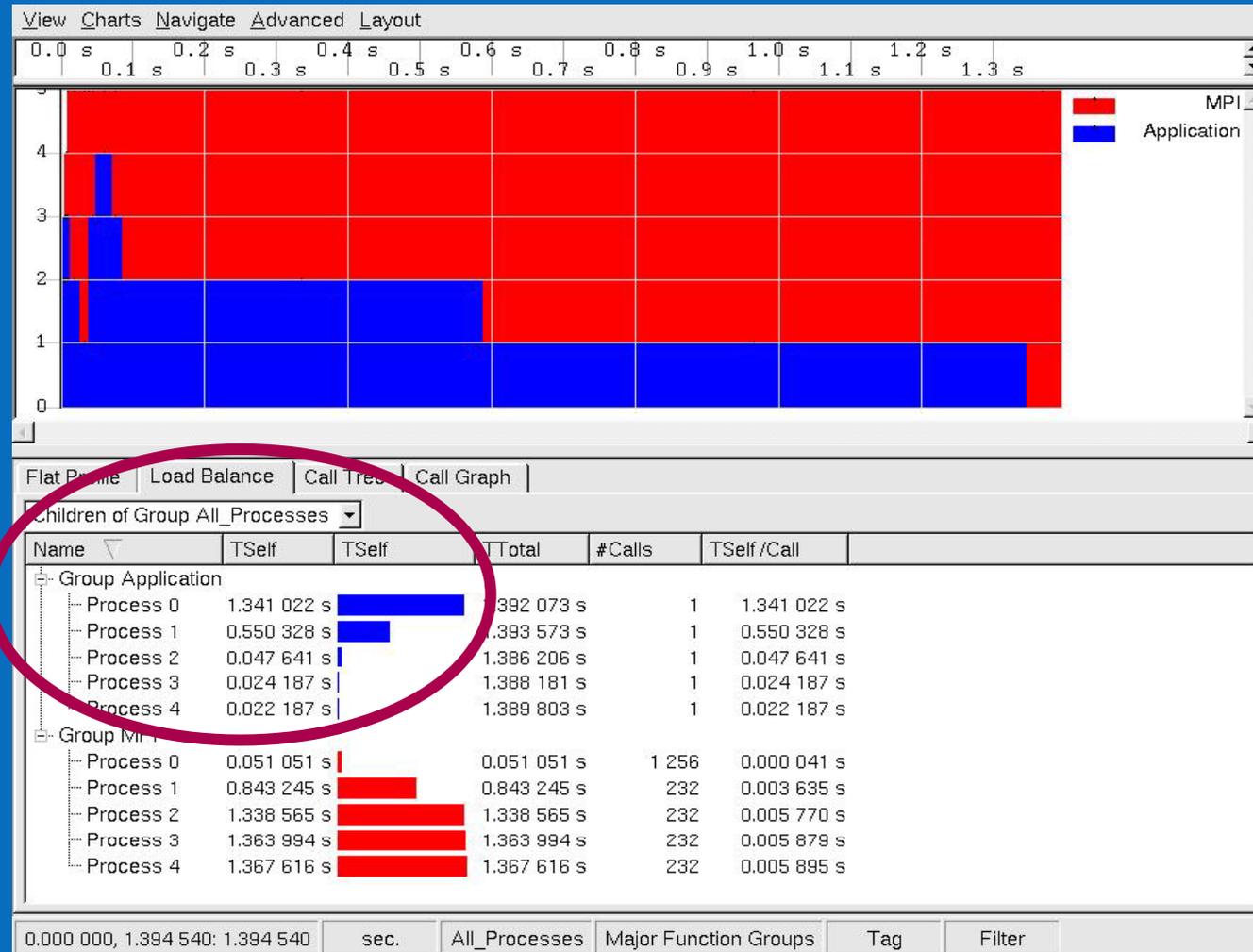
Detecting Load Imbalance

Mandelbrot set (MPI-tutorial)

- `mpitutorial.tar.gz`

Example of intuitive parallelization with huge load imbalance

Mandelbrot Example Trace



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



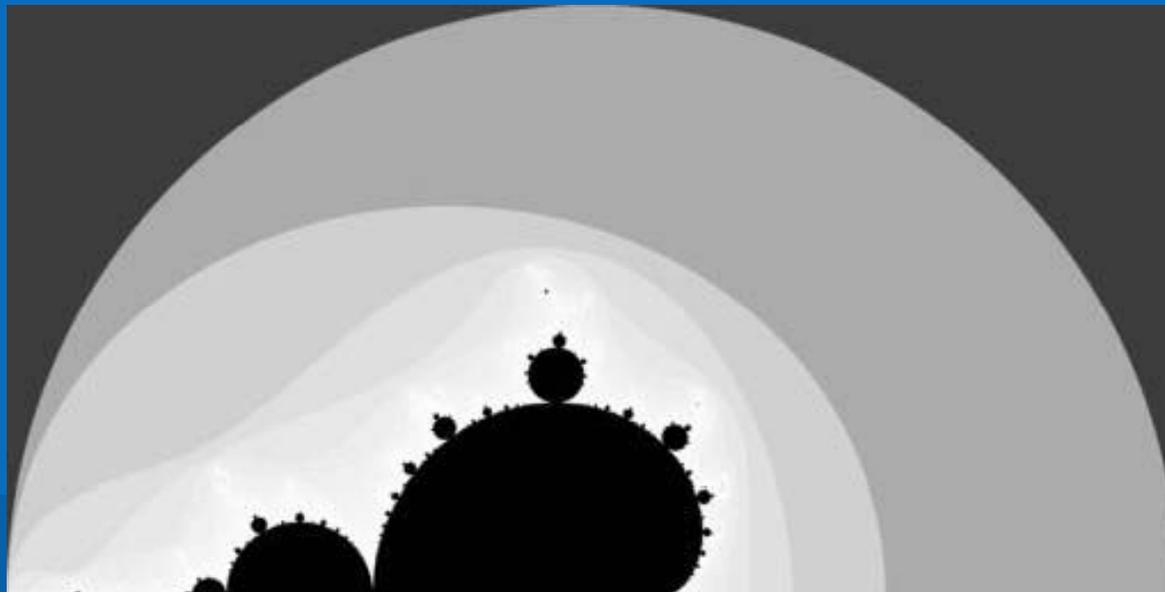
Understanding Load Imbalance

Static partitioning assigns contiguous blocks of rows

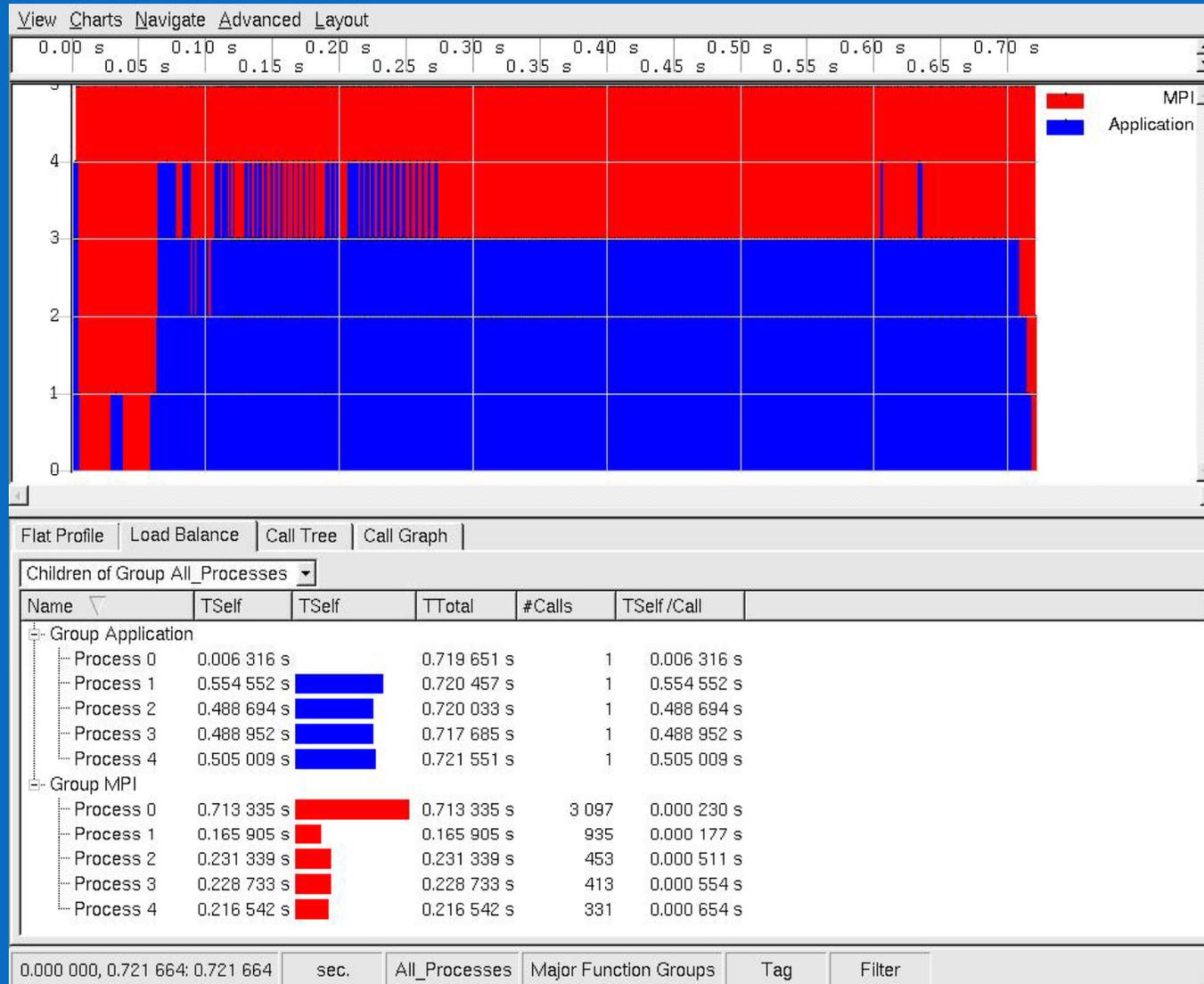
- Workload changes gradually with row-index!

What is the/a solution?

- Use modulo operator
- Use dynamic master/slave approach



Understanding Load Imbalance



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Investigating Scalability - Amdahl's Law

Non-parallel sections limit the possible speedup

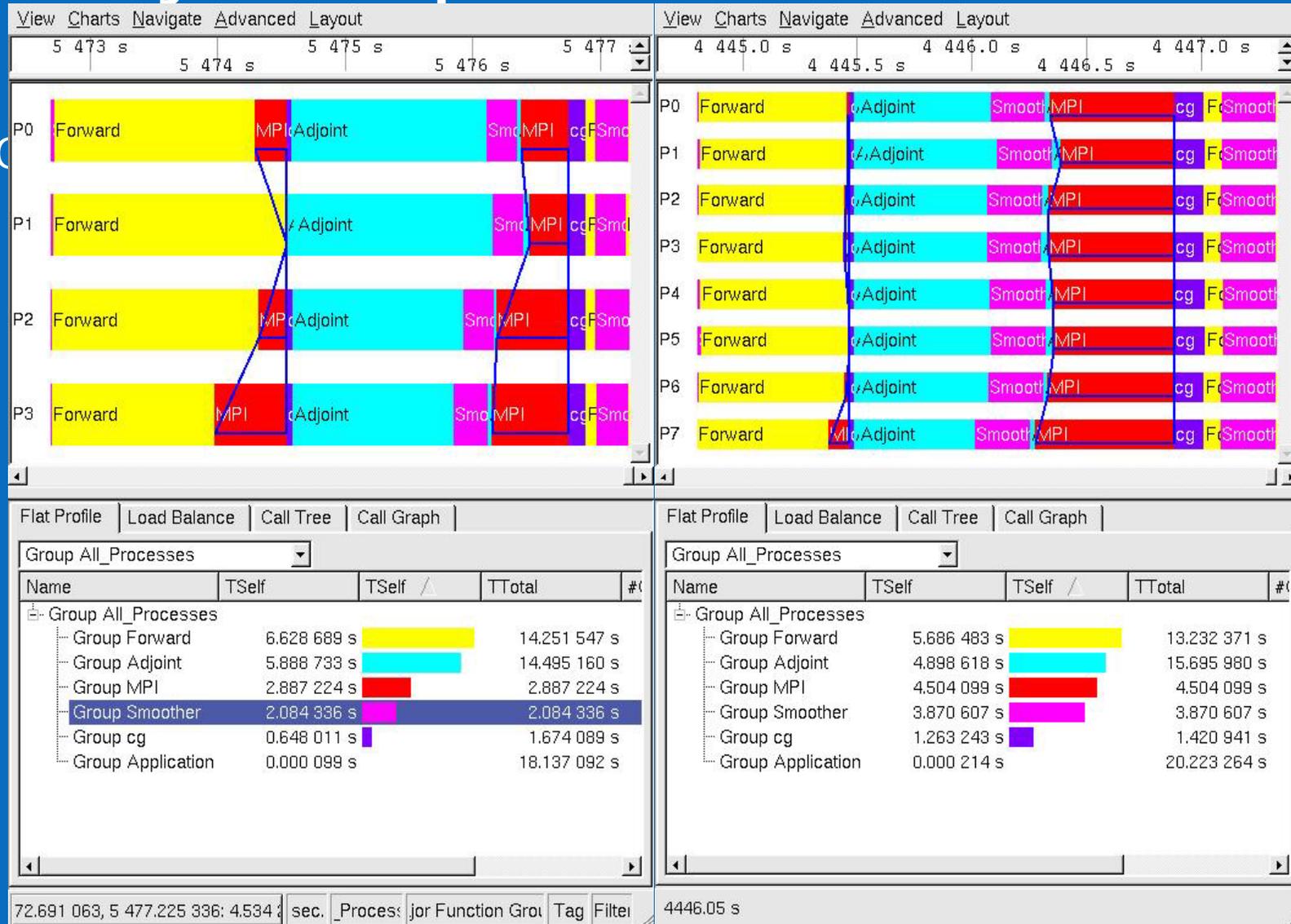
Amdahl's Law: If F is the fraction of a calculation that is sequential, and $(1-F)$ is the fraction that can be parallelized, then the maximum speedup that can be achieved by using P processors is $1/(F + (1-F)/P)$

2 runs with different number of processes:

- If sum of time spent in function increases with increasing number of processes
=> non-parallel

Scalability - Example

ZTOMC



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Detecting Communication Inefficiency

Message Profile

Use minimum transfer rate

Thread Groups and Aggregation

- All_nodes group
- Intra <-> extra node

Compare with HW expectation

Determining Overhead of Modules

Modules can contain more than one function

Use **Function Aggregation** to determine the module's overhead

Trace of Intel® Trace Analyzer

- Does ZLib produce significant overhead?

Focusing on the relevant

Goal: determine overhead of ZLib in the Intel Trace Analyzer

Create appropriate Function Group(s)

Use Tagging to zoom to representing interval

Use Function Profile

Step 1: Creating Group(s)

Function Group Editor ? X

View: 1: /Projects/psp/fschlimb/work/ict/tracing/mosel/IA32-LIN/bin/vtmos
elprof3.stf

Name	Depth	Children	Id
/	6	3	4294967284
[-] G All Functions	1	1612	2147483749
[-] G Major Function Groups	5	5	2147483648
[-] G Analysis	4	181	2147483692
[-] G GUI	4	62	2147483841
[-] G Graphics	4	5	2147484110
[-] G Misc	4	132	2147483995
[-] G STF	2	5	2147483935
[-] G New Group	3	1	4294967279
[-] G ZLib	2	9	4294967278

OK Cancel Apply

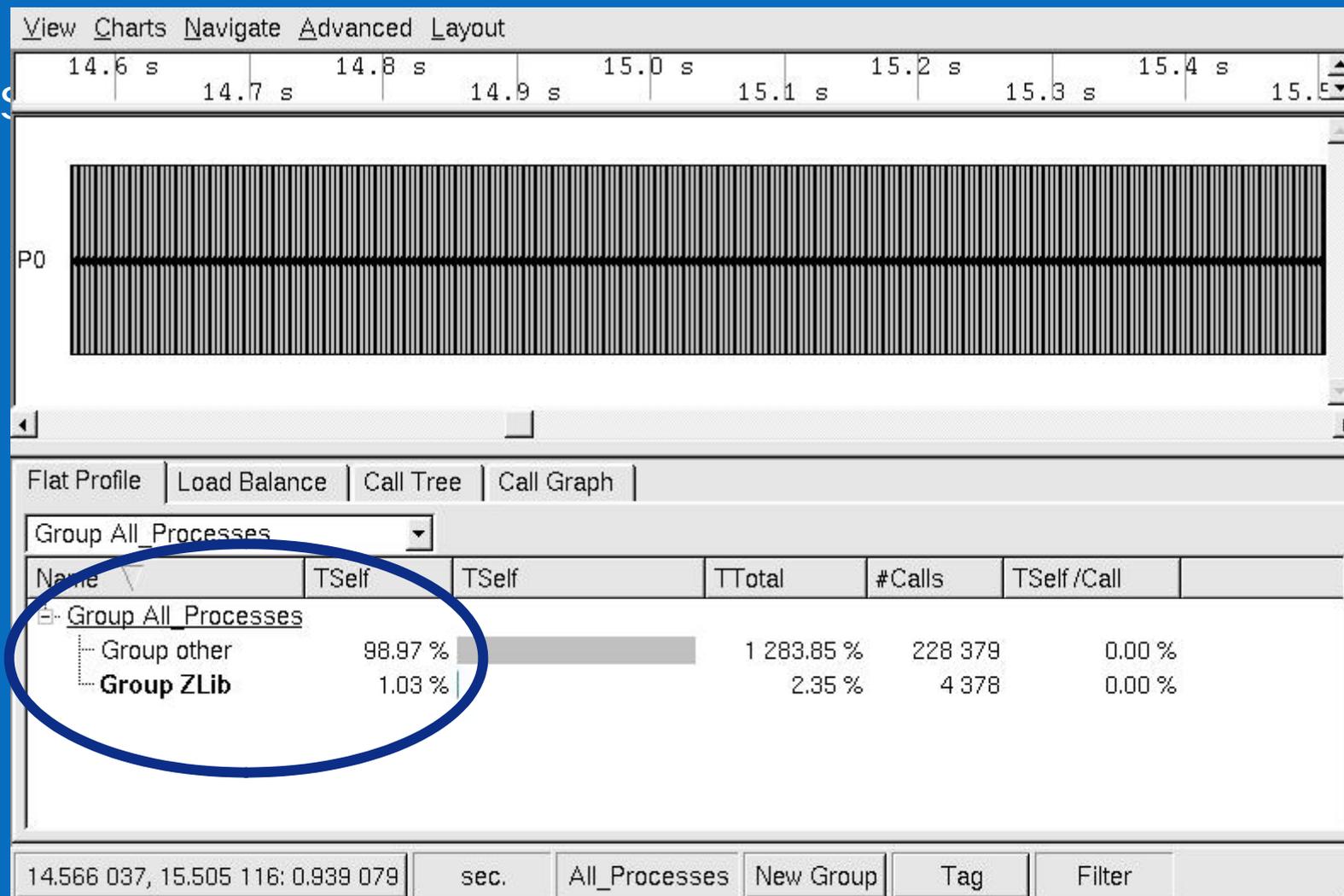
Step 2: Tagging and zooming

The image shows a software interface for tagging and zooming. At the top, a dialog box titled "Tagging" is open. It displays the current view path: "View: 1: /Projects/psp/fschlimb/work/ict/tracing/mosel/IA32-LIN/bin/vtmos elprof3.stf". Below this, there is a text field for the "Definition of Filter Expression". Two radio buttons are present: "Using GUI Interface" (which is selected) and "Manually".

Below the dialog box, a signal trace is shown. The trace consists of a dense series of vertical lines, with a vertical purple line indicating a specific point in time. The label "P0" is positioned to the left of the trace.

At the bottom, another dialog box is open, showing a list of filter clauses: "funcfilter(fg(ZLib))", "#", "p2pfilter(NONE)", "#", and "collfilter(NONE)". Buttons for "Add New Clause" and "Remove Current Clause" are visible above the list. At the bottom of this dialog are buttons for "Reset Dialog", "OK", "Cancel", and "Apply".

Step 3: Determine Overhead



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



67

Copyright © 2006, Intel Corporation. All rights reserved.

Intel and the Intel logo are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States or other countries. *Other brands and names are the property of their respective owners.



Lab 7 -



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Who was calling?

Start with

```
traces/ita/ita_fgt1_threaded.stf
```

Which caller of STF_WorkStackHistory caused the most costly calls to it?

Function Profile

- Advanced|Function Aggregation->All_Functions
- Search&Find STF_WorkStackHistory
- Show it in Call Graph
- Max TimeSelf/Call of Callers group

Instrumenting Functions

Generic API to instrument functions

Hierarchical grouping (classes)

Insert probes to

- declare function
- enter/leave the function

Lab 8 -



MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Lab 9 – Basic use of API

Objective: Use of Intel® Trace Collector API to instrument a function

Edit cpi.c to instrument function 'f' as of group 'Application'

- Use `VT_classdef()` and `VT_funcdef()` to declare class/function
- Use `VT_enter()` and `VT_leave()`, to record function begin and end

Compile-link-run and view tracefile in Intel® Trace Analyzer

- Find events 'f' in the timeline, in the function profile and in the call tree display

ITC API: Instrument a function

Allocate a handle for a class name (optionally)

```
int VT_classdef( const char * classname,  
                int * classhandle )
```

classname: name of the class

returns handle (int) referencing this class

Allocate a handle for a state

```
int VT_funcdef( const char * symname,  
               int classhandle,  
               int * statehandle )
```

symname: name of the symbol

classhandle: handle of the class (see above)

returns handle (int) referencing this function

Solution – Use of API

```
#include <VT.h>
[...]
```

```
int application_class;
```

```
double f( double a ){
    static int f_state = 0;
    if( f_state == 0 ) {
        VT_funcdef( "f", application_class, &f_state );
    }
    VT_enter( f_state, VT_NOSCL );
    [...]
    VT_leave( VT_NOSCL );
    return (result);
}
```

```
int main( int argc, char *argv[] )
[...]
```

```
MPI_Init(&argc,&argv);
VT_classdef( "Application", &application_class );
```

Lab 10 – Selective Tracing

Objective: Switch tracing off and on.

Edit **cpi.c** to instrument function 'f'.

- use
VT_traceoff(), VT_traceon()
do not record MPI_Reduce()

```
VT_traceoff();  
MPI_Reduce(&mympi, &pi, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);  
VT_traceon();
```

Lab 10 – Selective Tracing

Objective: Use of a VT_CONFIG configuration file

Compile-link-run-tracefile and check with Intel® Trace Analyzer.

- check timeline display for `TRACE_OFF` region where `MPI_Reduce` had been removed.
- Create a file '`my_vt_config`' with the following line in it `SYMBOL MPI_Bcast OFF`
- Set environment variable `VT_CONFIG` to refer to complete path of '`pwd` /my_vt_config`'
- Check timeline or statistics ... no more data of `MPI_Bcast`

Lab 11 – Instrumentation

Objective: Use of Intel® Trace Collector to instrument a binary

*Note: shared MPI libraries ***MUST*** be available*

Currently only available for IA32

Instrument binary:

```
itcinstrument --input ./cpi_vt --mpi /opt/intel/mpi_2.0/lib --  
profile --output ./cpi_vt_instr
```

Compile-link - Run the instrumented binary and check tracefile with Intel® Analyzer

Lab 12 – Compare applications

Objective: Analyze and compare solver application with Intel® Trace Analyzer

The laplace equation is solved by two different iterative methods

- Using Jacobi iteration method
- Using Successive Over Relaxation

Open both tracefiles in traceanalyzer (in subdirectory ./traces) - use Window|Tile to compare the traces.

Which solver method seems to be best regarding:

- Amount of MPI?
- Load balance?
- Performance?

Backup

MPI Tuning with Intel® Trace Analyzer & Intel® Trace Collector



Lab 12 – Create STF single file

Objective: Learn about stftool, use single trace file to simplify transfer of storage

Convert set of stf (structured trace format) files to single file:

```
stftool cpi_itc.stf --logfile-format stfsingle --convert  
cpi_itc.single.stf
```

Note: Can degrade performance of ITA

Note: ITC config allows creation of single-STF

Lab 13 – Statistics trace file

Objective: How to generate a tracefile with only statistic information

Generate file my_vt_config

- # enable statistics gathering
- STATISTICS ON
- #no need to gather trace data
- PROCESS 0:N OFF

Set environment variable and build example

- > setenv VT_CONFIG \$PWD/my_vt_config

Run example and look at stats-file cpi_itc.prot

- > stftool cpi_itc.stf --print-statistics

