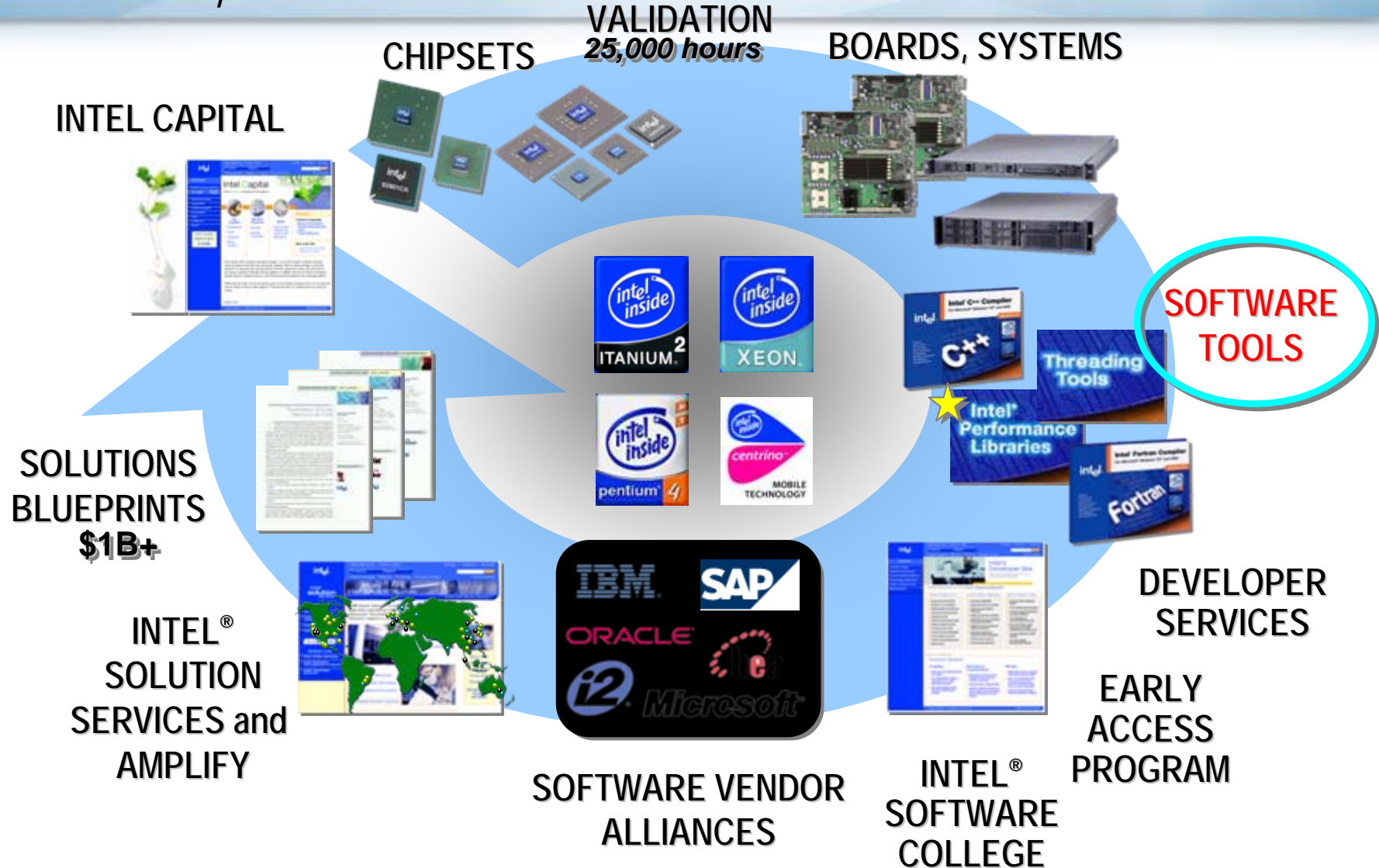# Intel Cluster Tools

Michael Hebenstreit

Graz, May 2006
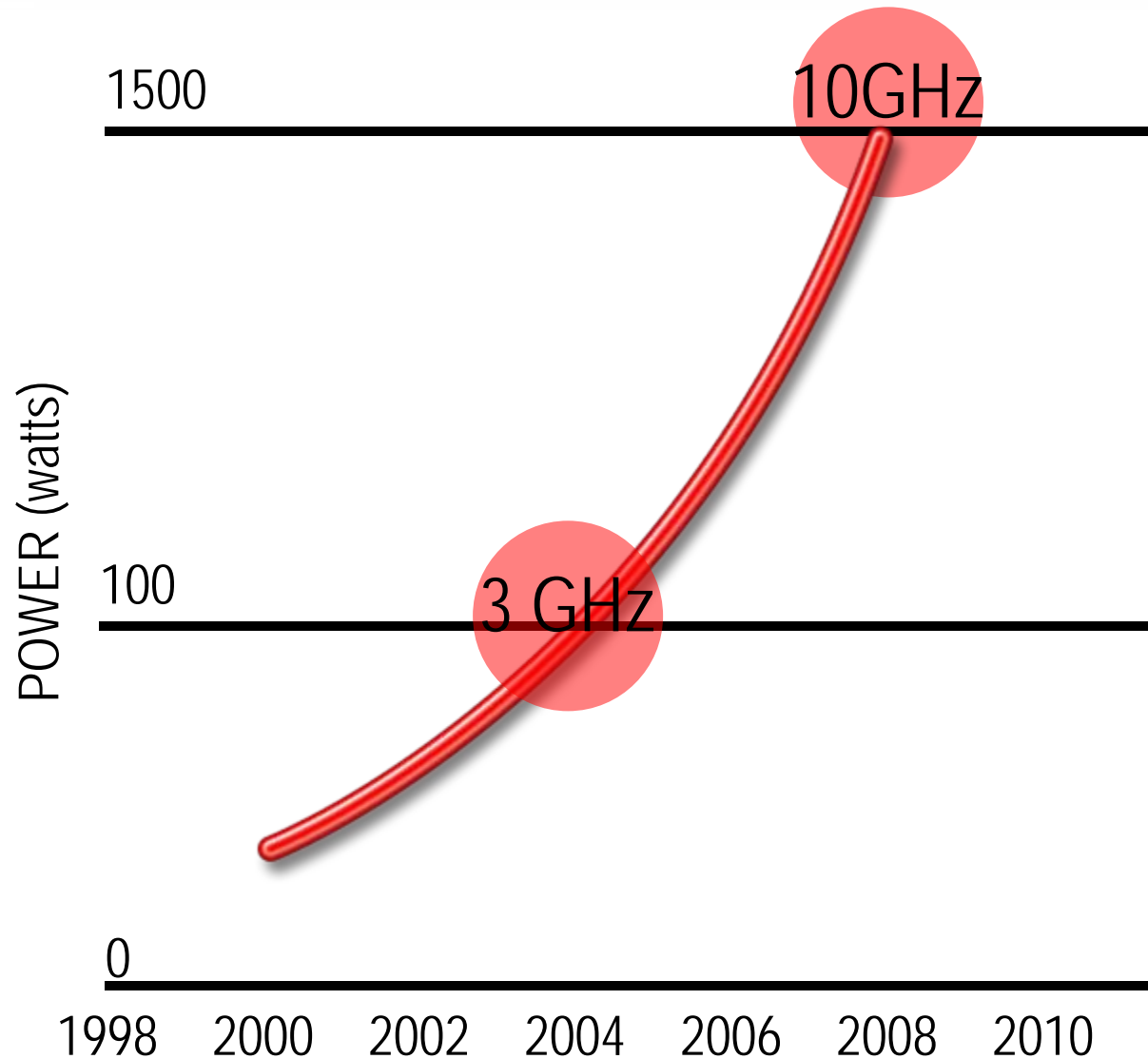
**SSG
Software Solutions Group**

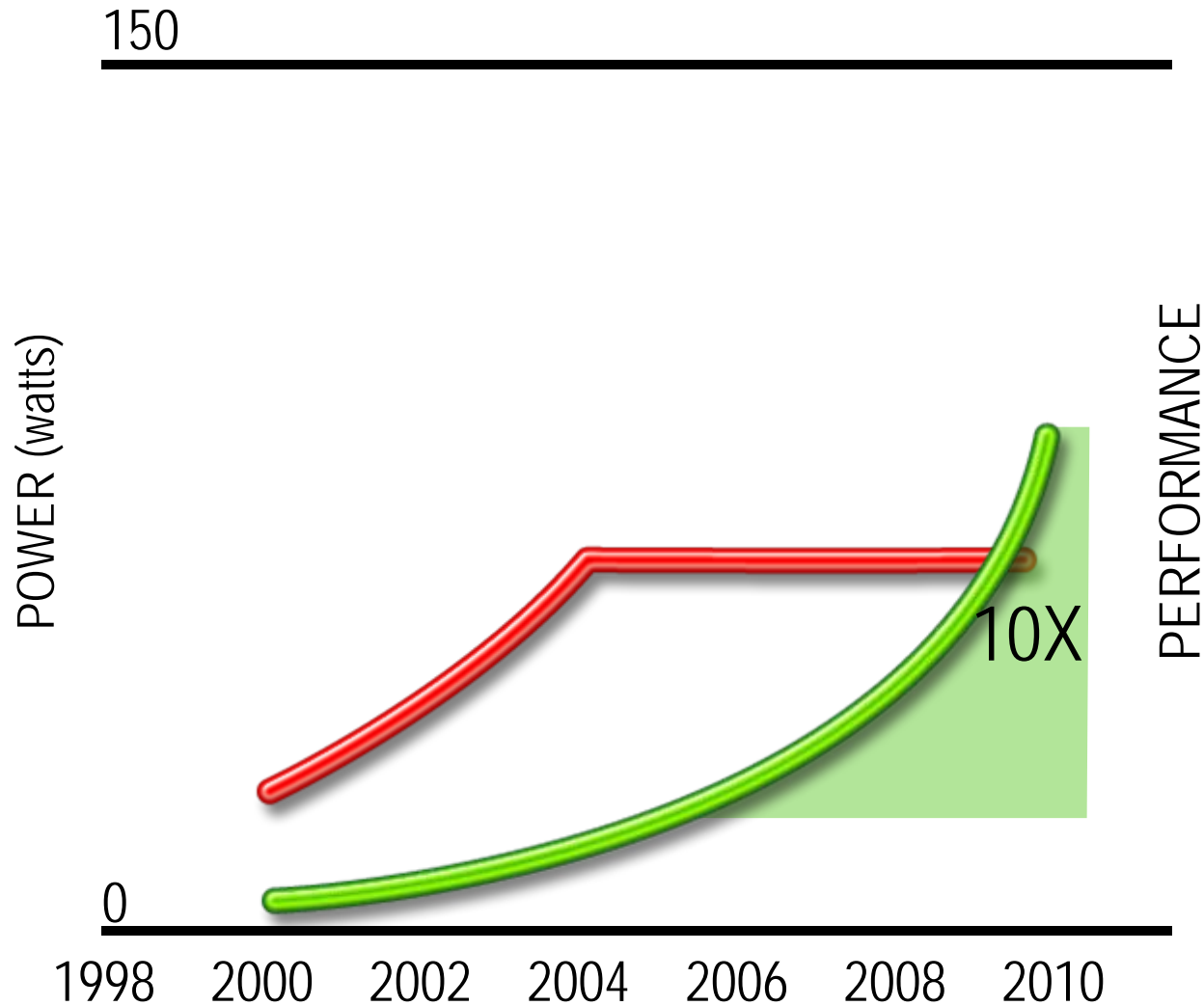# Beyond Silicon:
*Intel's Enterprise Software and Solution Focus*



VALIDATION
*25,000 hours*

CHIPSETS

BOARDS, SYSTEMS

INTEL CAPITAL

SOFTWARE
TOOLS

SOLUTIONS
BLUEPRINTS
$1B+

DEVELOPER
SERVICES

INTEL®
SOLUTION
SERVICES and
AMPLIFY

EARLY
ACCESS
PROGRAM

SOFTWARE VENDOR
ALLIANCES

INTEL®
SOFTWARE
COLLEGE

# The Path We WERE On

# The Path We ARE On

**10's to 100's of cores**

*Era of Tera-Scale Computing*

**Quad-Core**

*More performance Using less energy (See Keynote Tomorrow)*

**Dual Core**

**Hyper-Threading**

*The days of single-core chips*

**Instruction level parallelism**
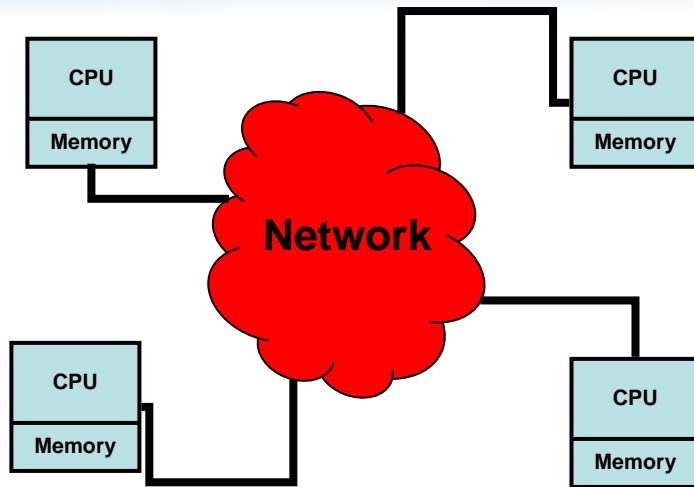
ENERGY-EFFICIENT PERFORMANCE

TIME
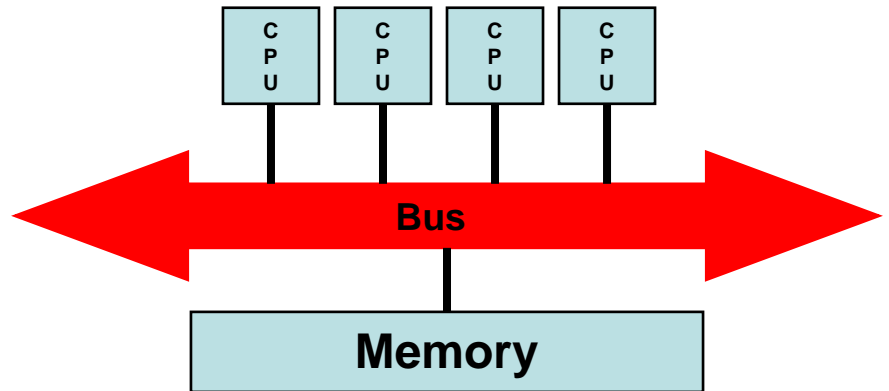
# Parallel Programming: Algorithms

**SSG**
**Software Solutions Group**

# Distributed Versus Shared Memory



| Message Passing |
|---|
| • Multiple processes<br>• Share data with messages<br>• MPI |

| Threads |
|---|
| • Single process<br>  – Concurrent execution<br>• Shared memory and resources<br>• Explicit threads, OpenMP* |

# Parallelize Loops

- Find your most time consuming loops.
- Split them up between threads.

```
void main()
{
    double Res[1000];
    for(int i=0;i<1000;i++){
        do_huge_comp(Res[i]);
    }
}
```

**Sequential Program**

→

```
#include "omp.h"
void main()
{
    double Res[1000];
#pragma omp parallel for
    for(int i=0;i<1000;i++){
        do_huge_comp(Res[i]);
    }
}
```

**Parallel  Program**

# Pi Integration parallel

```c
#include <stdio.h>
#include "mpi.h"
#define INTERVALS 10000000
int main (int argc, char* argv[])
{
  double partialSum, myPi = 0.0;int nProc, myID, source, dest, tag = 0;
   int i;double h, x, pi = 0.0; MPI_Status status;

  MPI_Init (&argc, &argv);
  MPI_Comm_size (MPI_COMM_WORLD, &nProc);
  MPI_Comm_rank (MPI_COMM_WORLD, &myID);

  h = 1.0 / (double)INTERVALS;

  for (i = myID + 1; i <= INTERVALS; i += nProc)  {
    x = h * ((double)i - 0.5);
    myPi += 4.0 / (1.0 + x * x);
  }
  MPI_Reduce (&myPi, &pi, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);

  if (myID == 0) {pi *= h; printf ("Pi = %.9f\n", pi);}
  MPI_Finalize ();
}
```

**SSG**

(intel)

# Conclusions

- Hardware architecture already complicated – and will not get easier to use in the future.

## Use tools wherever possible to reduce errors and enhance performance!

# Introduction to VTune
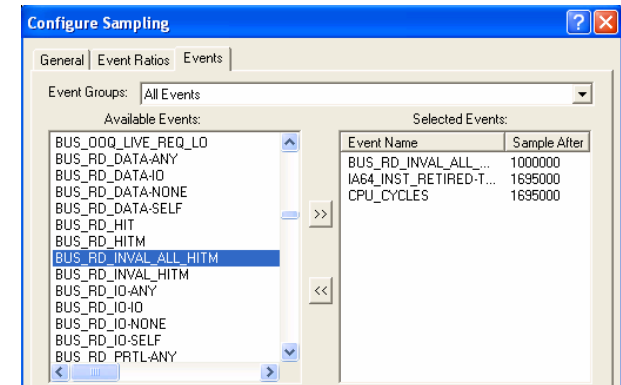**Optimization on a single Node**
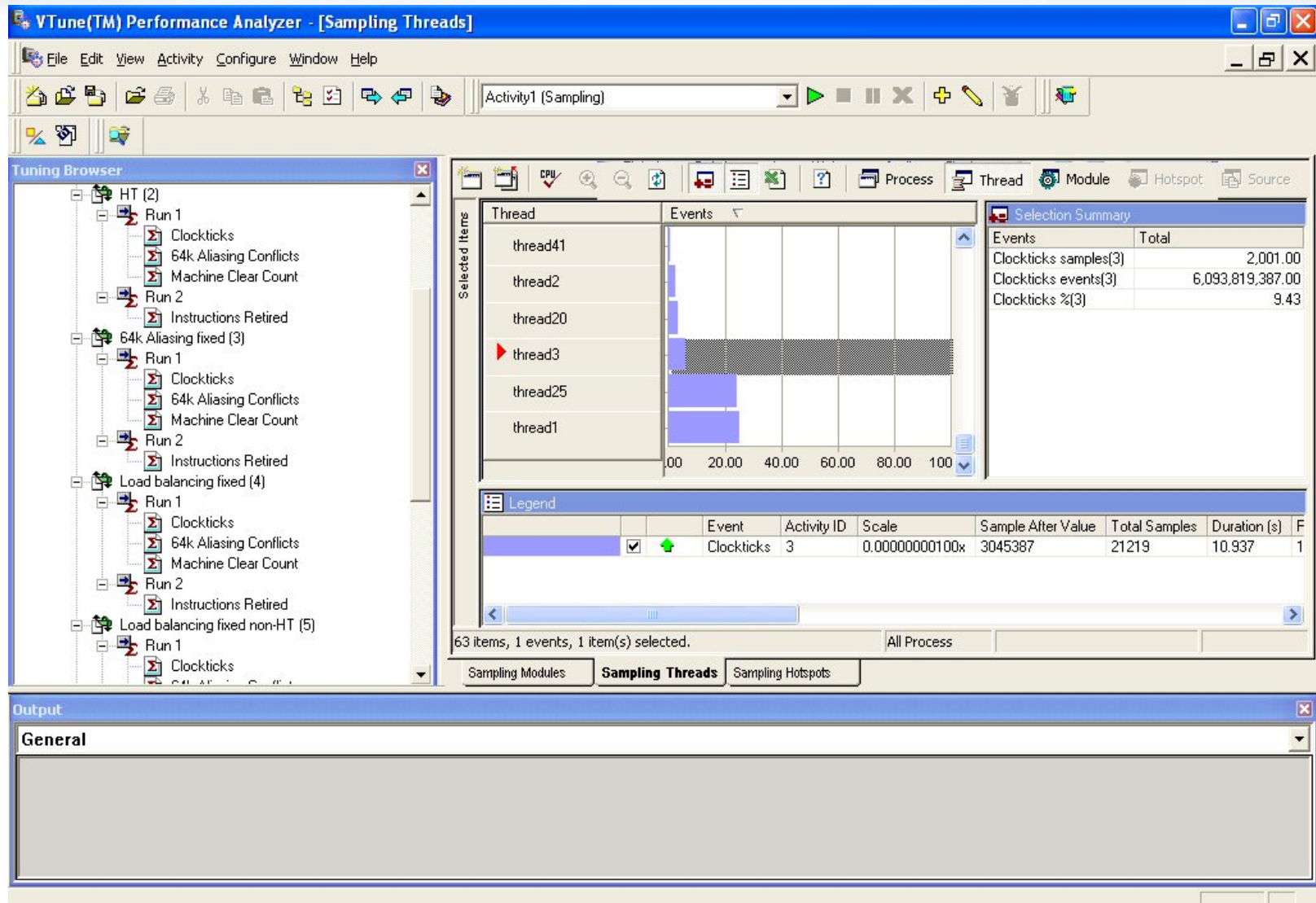
**SSG**
**Software Solutions Group**

# Events

- Clockticks

- Instructions Retired

- VTune™ Analyzer combines the counters to find the event ratio
  Cycles Per Instruction (CPI)
  - CPI = clockticks / instructions retired

# Sampling - GUI

# Dig down to the source view

# Benefits of Sampling

- Traditional Instrumentation:

  - Modified code

  - Higher overhead

  - Function-specific

- VTune™ analyzer Sampling:

  - Unmodified code

  - Lower overhead

  - System wide coverage

# Intel® Thread Checker

A "correctness" tool for threaded apps

- Finds errors in Win32*- and OpenMP*-threaded apps

- Detects incorrect memory usage – data races, deadlocks, ...more

- Plug-in for VTune™



| Context [Best] | ID | Seve... | Description | Counts | 1st Access [Routine] | 1st Access [Variable] | 1st Access [ |
|---|---|---|---|---|---|---|---|
| Total | | | | | | | |
| Group 1 : "Pi.cpp" : 12 | | | | | | | |
| "Pi.cpp" : 12 | 0 | 🔴 | Memory read of pArg at "Pi.cpp" : 13 conflicts with a prior memory write of i at "Pi.cpp" : 51 (flow dependence) | 16 | main | i | "Pi.cpp" : 5 |
| "Pi.cpp" : 12 | 1 | 🔴 | Memory write of dSum at "Pi.cpp" : 24 conflicts with a prior memory read dSum at "Pi.cpp" : 24 (anti dependenc... | 1240 | PiFunc | dSum | "Pi.cpp" : 2 |
| "Pi.cpp" : 12 | 2 | 🔴 | Memory read of dSum at "Pi.cpp" : 24 conflicts with a prior memory write of dSum at "Pi.cpp" : 24 (flow dependen... | 1240 | PiFunc | dSum | "Pi.cpp" : 2 |
| "Pi.cpp" : 12 | 3 | 🔴 | Memory write of dSum at "Pi.cpp" : 24 conflicts with a prior memory write of dSum at "Pi.cpp" : 24 (output depend... | 1240 | PiFunc | dSum | "Pi.cpp" : 2 |
| Group 2 : Whole Program 1 | | | | | | | |
| Whole Program 1 | 4 | 🔵 | Thread Info at "Pi.cpp" : 55 - includes stack allocated of 0 and use of 1044 | 1 | main | unknown | "Pi.cpp" : 55 |

Intel® Thread Checker Results - Thu Sep 19 18:04:37 2002: Error List

# Intel® Thread Profiler 2.0

- Examines processor utilization to determine parallel activity of the application

- *Concurrency* is the number of threads that are *active*



**Note: Categorization shown for a system configuration with 2 processors**

# Compiler optimization

**SSG**
**Software Solutions Group**

# General Optimizations

| Linux* | Windows* | |
|--------|----------|---|
| -O0 | /Od | Disables optimizations |
| -g | /Zi | Creates symbols |
| -O1 | /O1 | Optimizes for speed without increasing code size (disables library function inlining) |
| -O2 | /O2 | Optimizes for speed (default) |
| -O3 | /O3 | High-level optimizations |

# Benchmarks

## Results from „square_charge" dataset Linux, single Northwood 1.7GHz

| | |
|---|---:|
| -O2 | 20.34 s |
| -ipo | 20.34 s |
| -xk | 13.26 s |
| -ipo -xk | 1.71 s |
| -ipo –xk -openmp | 1.67 s |

# Vector-processing: convert this

```
for (I=0;I<=MAX;I++)
    c[I]=a[I]+b[I];
```

| A[1] | not used | not used | not used |

\+

| B[1] | not used | not used | not used |

| C[1] | not used | not used | not used |

# SSE Registers Usage

- xP switch is used for vectorizing loops as well as floating point and scalar operations.

- Example:

```
for (I=0;I<=MAX;I++)
    c[I]=a[I]+b[I];
```

- Usage:      -xP

# Why Didn't My Loop Vectorize?

Linux                 Windows

`-vec_report`$n$     `-Qvec_report`$n$

Set diagnostic level dumped to stdout

$n=0$: No diagnostic information

$n=1$: **(Default)** Loops successfully vectorized

$n=2$: Loops not vectorized – and the reason why not

$n=3$: Adds dependency Information

$n=4$: Reports only non-vectorized loops

$n=5$: Reports only non-vectorized loops and adds dependency info

# Intel MPI

**SSG
Software Solutions Group**

# Features

- MPI-2 specification conformance
  - Standardized job startup mechanism (mpiexec)
  - Process spawning/attachment (sock device only)
  - One-sided communication
  - Extended collective operations
  - File I/O
- Multiple communication fabrics selectable at run-time
- Dynamic or static linking, plus debugging and tracing libraries
- Support for additional tools

# MPI Benchmark

**SSG**
**Software Solutions Group**

# MPI transfer measured at DC



MPI all to all - 32 CPUs

# Current switch configuration

| 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 |
|----|----|----|----|----|----|----|----|
| 38 | 39 | 40 | 41 | 42 | 43 | 44 | 37 |

| 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 |
|----|----|----|----|----|----|----|----|
| 23 | 24 | 25 | 26 | 27 | 28 | 22 |    |

| 47 | 48 | 49 | 50 | 51 | 52 | 53 | 54 |
|----|----|----|----|----|----|----|----|
| 56 | 57 | 58 | 59 | 60 | 61 | 62 | 55 |

| 01 | 02 | 03 | 04 | 05 | 06 | 07 |    |
|----|----|----|----|----|----|----|----|
| 08 | 09 | 10 | 11 | 12 | 13 | 14 |    |

|    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|
| 45 | 46 |    |    |    |    | 63 | 64 |

- Each number signifies one specific node.
- User typically request a block 01-16 or 49-64.
- Communication uses typically only one or two modules.

# Striped switch configuration

| 01 | 06 | 11 | 16 | 21 | 26 | 31 | 36 |
|----|----|----|----|----|----|----|----|
| 41 | 46 | 51 | 56 | 61 |    |    |    |

| 02 | 07 | 12 | 17 | 22 | 27 | 32 | 37 |
|----|----|----|----|----|----|----|----|
| 42 | 47 | 52 | 57 | 62 |    |    |    |

| 03 | 08 | 13 | 18 | 23 | 28 | 33 | 38 |
|----|----|----|----|----|----|----|----|
| 43 | 48 | 53 | 58 | 63 |    |    |    |

| 04 | 09 | 14 | 19 | 24 | 29 | 34 | 39 |
|----|----|----|----|----|----|----|----|
| 44 | 49 | 54 | 59 | 64 |    |    |    |

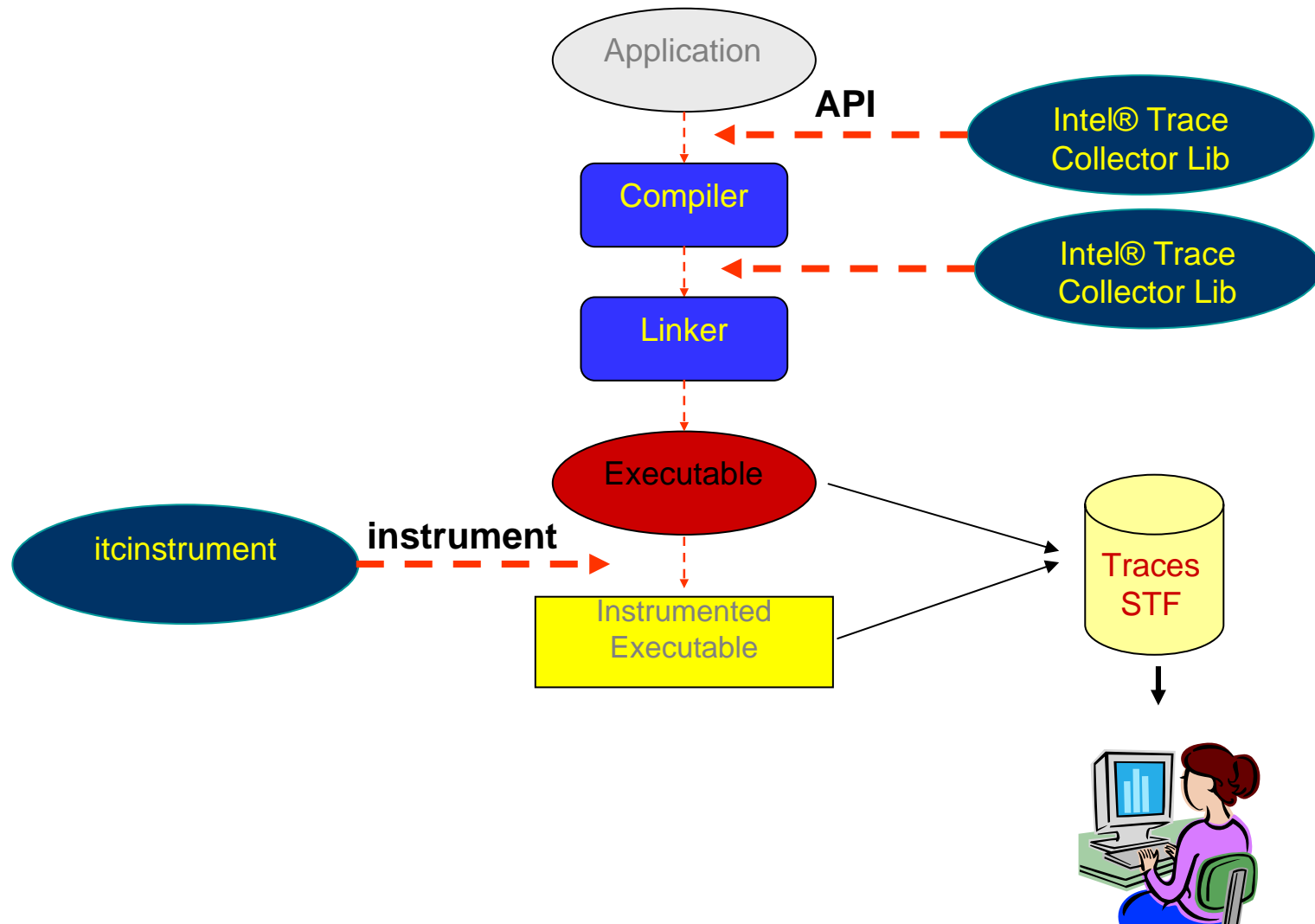| 05 | 10 | 15 | 20 | 25 | 30 | 35 | 40 |
|----|----|----|----|----|----|----|----|
| 45 | 50 | 55 | 60 |    |    |    |    |

- Each number signifies one specific node.
- User typically request a block 01-16 or 49-64.
- Communication uses typically all modules.

# Intel® Trace Analyzer and Collector 6.0

**SSG**
**Software Solutions Group**

# Components and Interaction
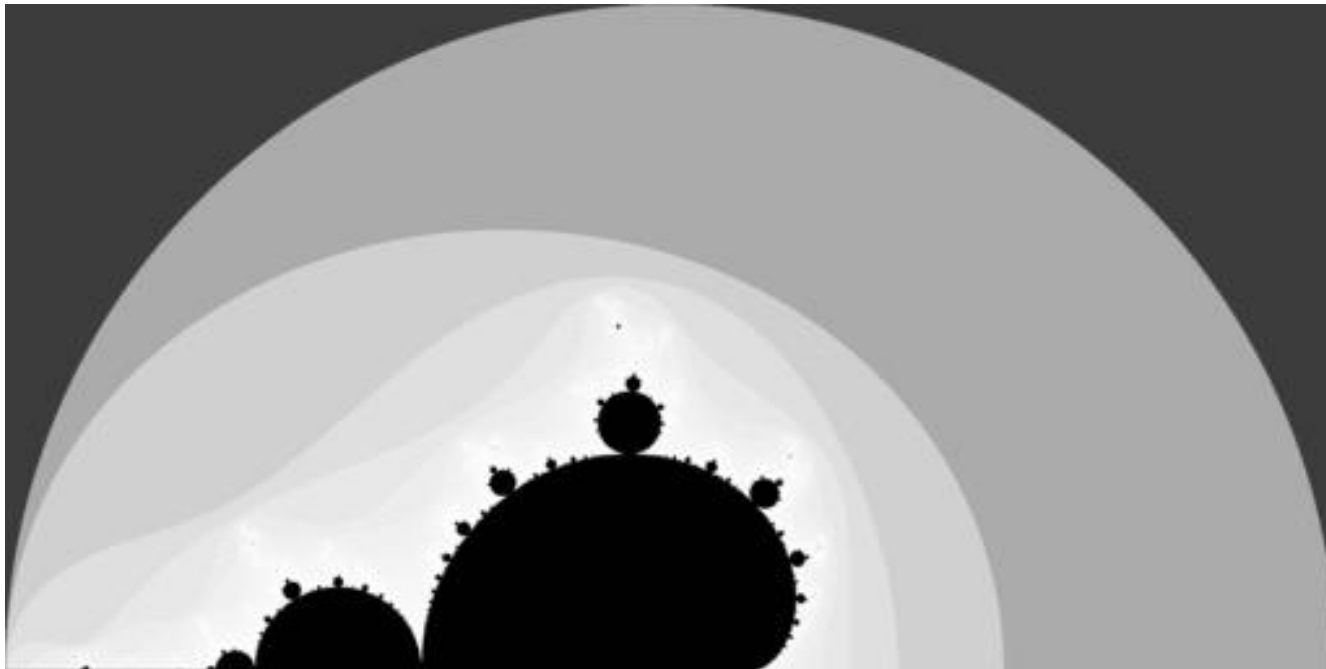
# Case Study 1 – performance measurements

- CFD code
  - setup/iteration part
  - usullay > 100 iterations, but only 10 for performance benchmark
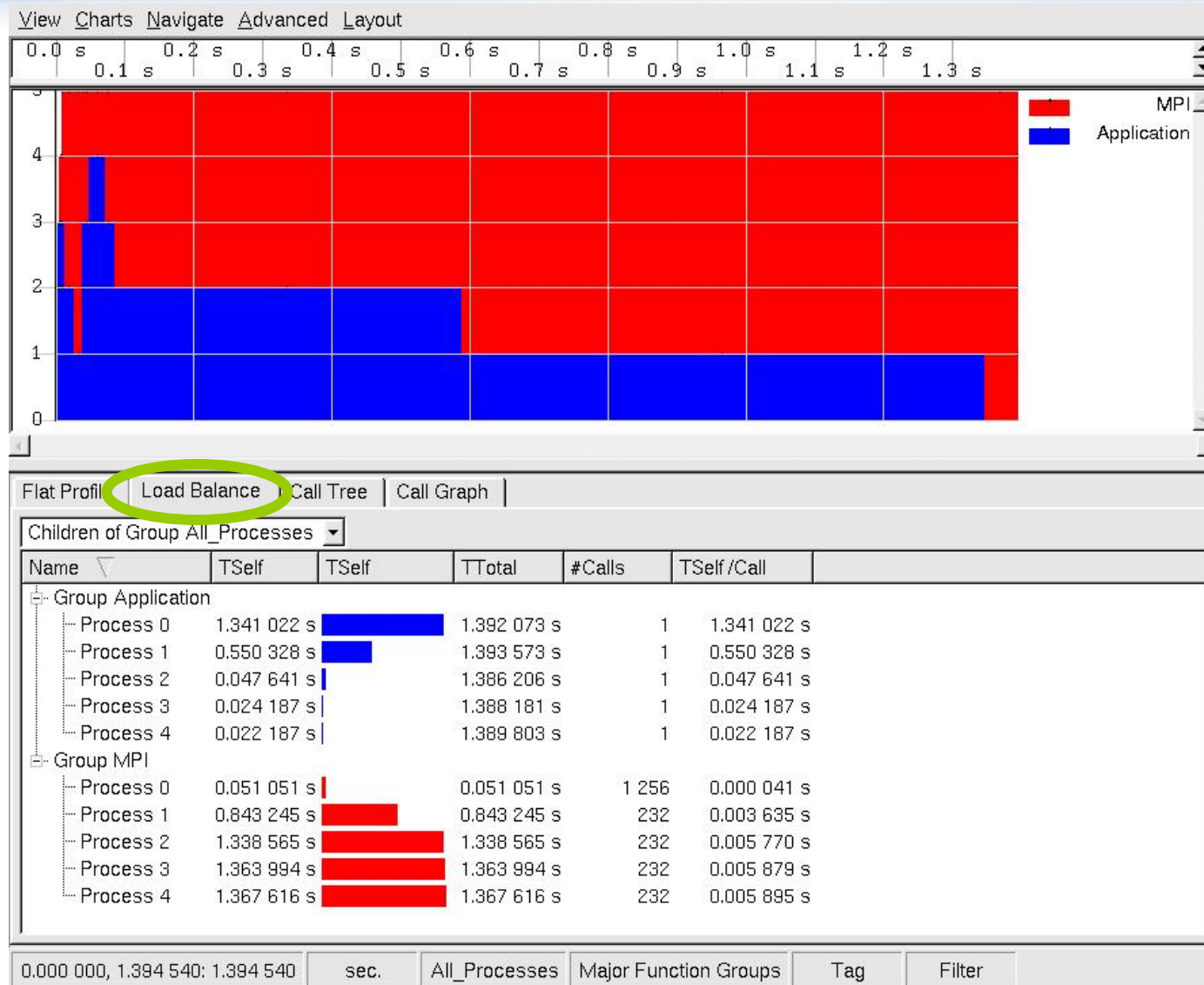
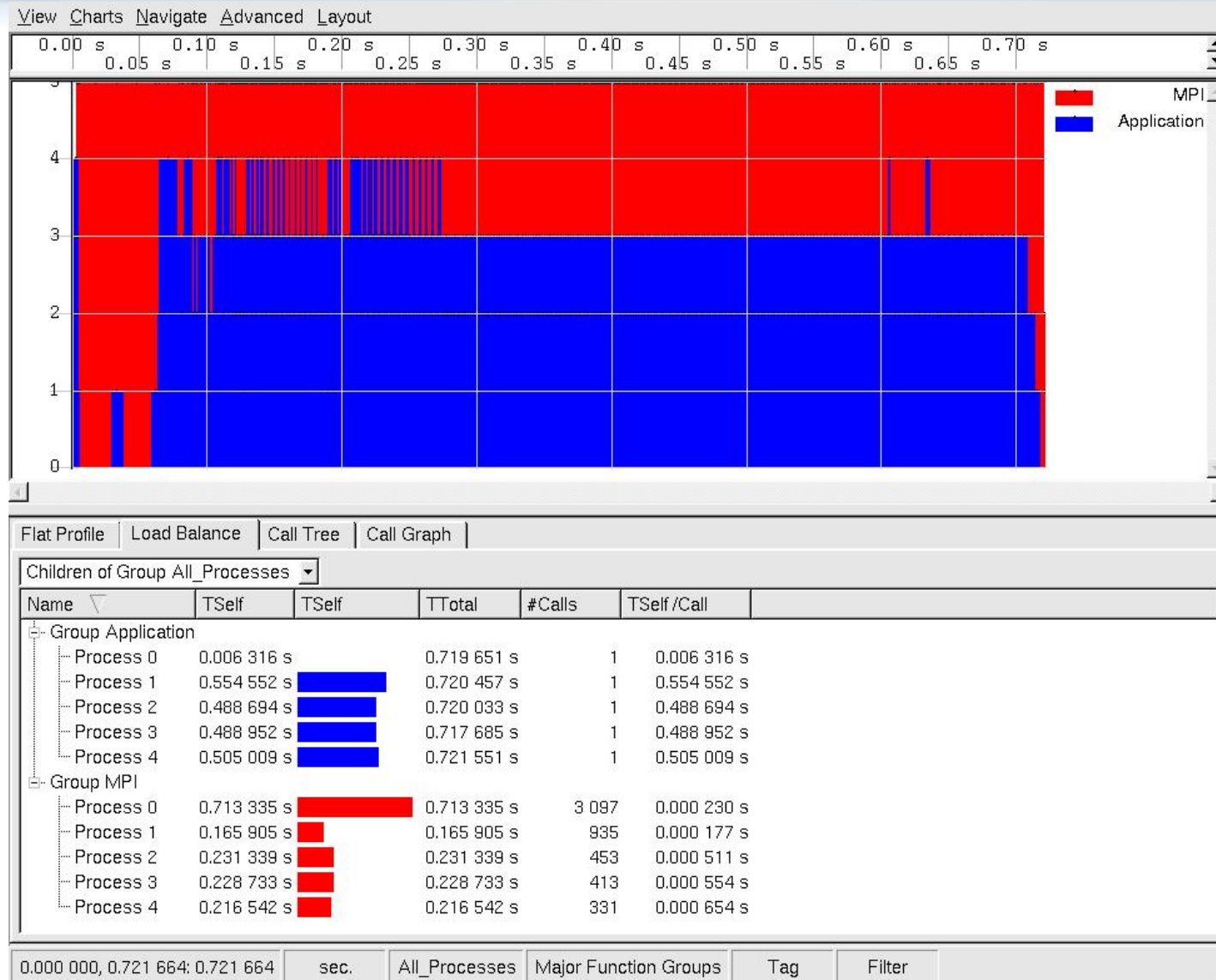# Case Study 2 - Load Imbalance

- Static partitioning assigns contiguous blocks of rows
  - Workload changes gradually with row-index!

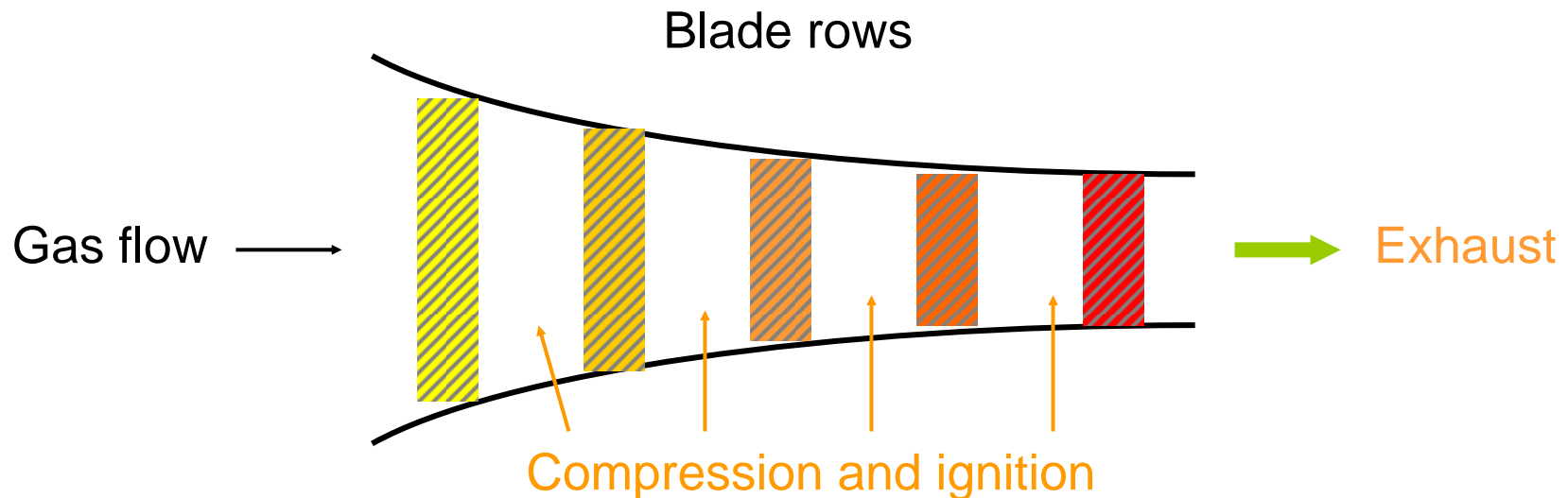# Detecting Load Imbalance

# Solving Load Imbalance

# Case Study 3 – programming error

- Each blade row is independent MPI simulation
- Blade row simulations run concurrently with file-level synchronization
- Goal: Optimize blade row simulations

Blade rows

Gas flow ⟶

Exhaust

Compression and ignition

# Three Customer Workloads

- Test problems 1 and 2
  - Large and small blade row simulations
  - Scale well
  - Represent current production requirements

- Test problem 3
  - Bisected blade row simulation
  - Scales poorly
  - Represents future production requirements

# ITC Result Case 1 & 2



vampir.10.bpv: Message Statistics (Sum, Length, 0.0 s-4:01.938)

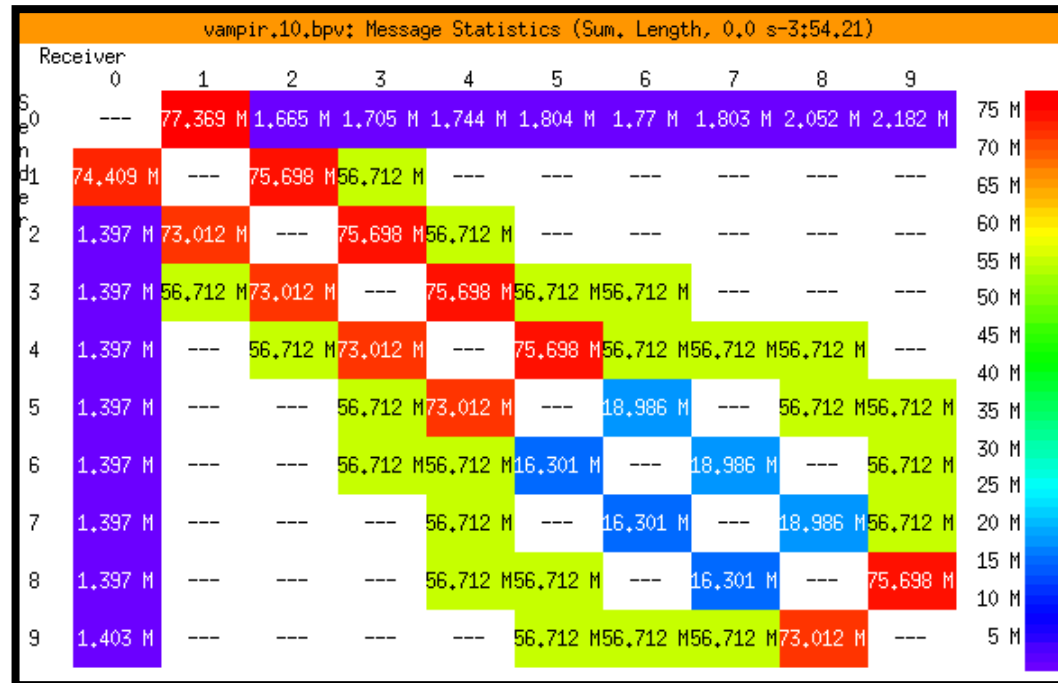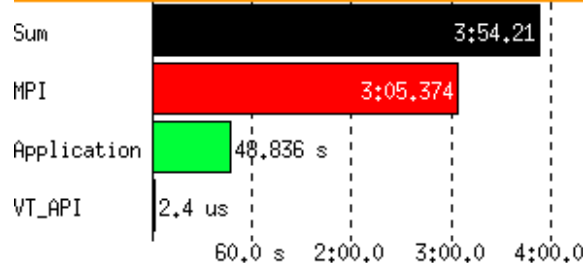| Receiver | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | --- | 72.007 M | 4.884 M | 4.797 M | 4.728 M | 4.544 M | 4.542 M | 4.619 M | 4.334 M | 3.686 M |
| 1 | 61.184 M | --- | 67.052 M | --- | --- | --- | --- | --- | --- | --- |
| 2 | 2.921 M | 58.263 M | --- | 67.052 M | --- | --- | --- | --- | --- | --- |
| 3 | 2.921 M | --- | 58.263 M | --- | 67.052 M | --- | --- | --- | --- | --- |
| 4 | 2.921 M | --- | --- | 58.263 M | --- | 67.052 M | --- | --- | --- | --- |
| 5 | 2.921 M | --- | --- | --- | 58.263 M | --- | 67.052 M | --- | --- | --- |
| 6 | 2.921 M | --- | --- | --- | --- | 58.263 M | --- | 67.052 M | --- | --- |
| 7 | 2.921 M | --- | --- | --- | --- | --- | 58.263 M | --- | 67.052 M | --- |
| 8 | 2.78 M | --- | --- | --- | --- | --- | --- | 58.263 M | --- | 67.052 M |
| 9 | 3.137 M | --- | --- | --- | --- | --- | --- | --- | 58.263 M | --- |

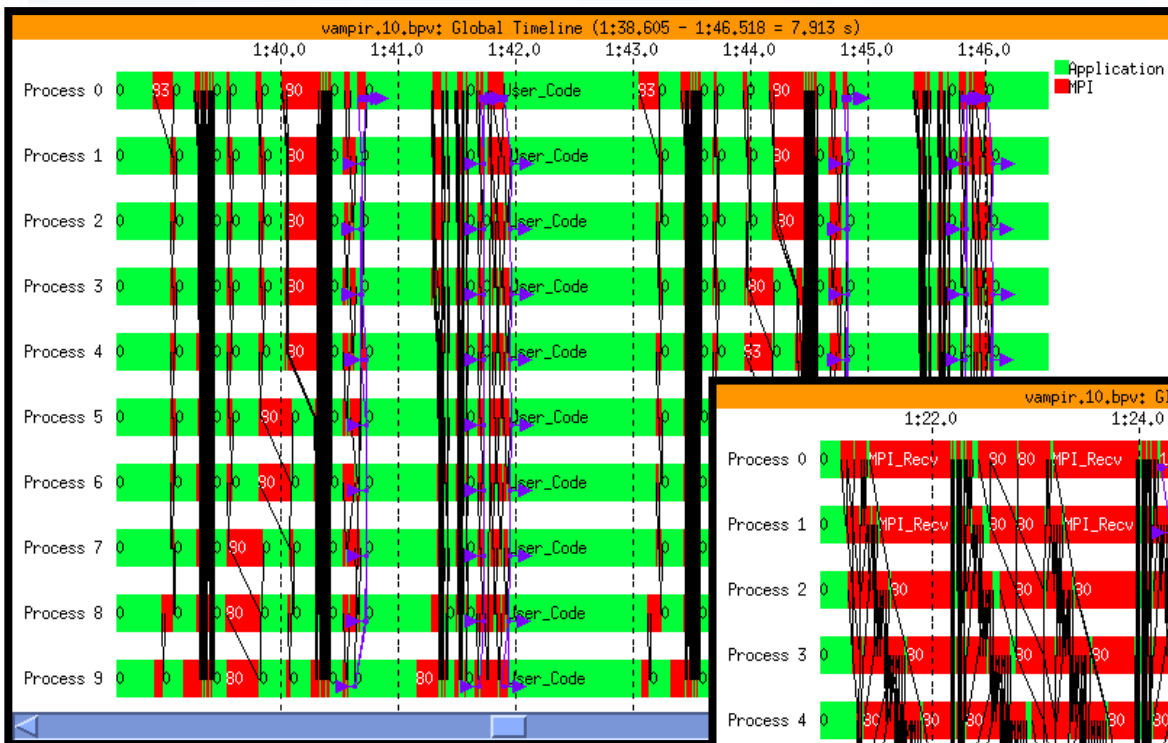Cases 1 and 2, which scale well, have a simple nearest-neighbor communication scheme

# ITC Result Case 3

Case 3, which scales poorly, has a more complex communication scheme

# MPI Communication Timeline

# Optimized Communication

**Before**

**Case 3 with communication bottleneck**



vampir.10.bpv: Summary Chart (Times, 0.0 s–3:54.21)

| | |
|---|---|
| Sum | 3:54.21 |
| MPI | 3:05.374 |
| Application | 48.836 s |
| VT_API | 2.4 us |

60.0 s   2:00.0   3:00.0   4:00.0



pir.10.async.bpv: Summary Chart (Times, 0.0 s–1:17.5

| | |
|---|---|
| Sum | 1:17.569 |
| Application | 48.715 s |
| MPI | 28.854 s |
| VT_API | 2.3 us |

20.0 s   40.0 s   1:00.0   1:20.0

**Compiler and MPI tuning yielded 2.7X speedup**
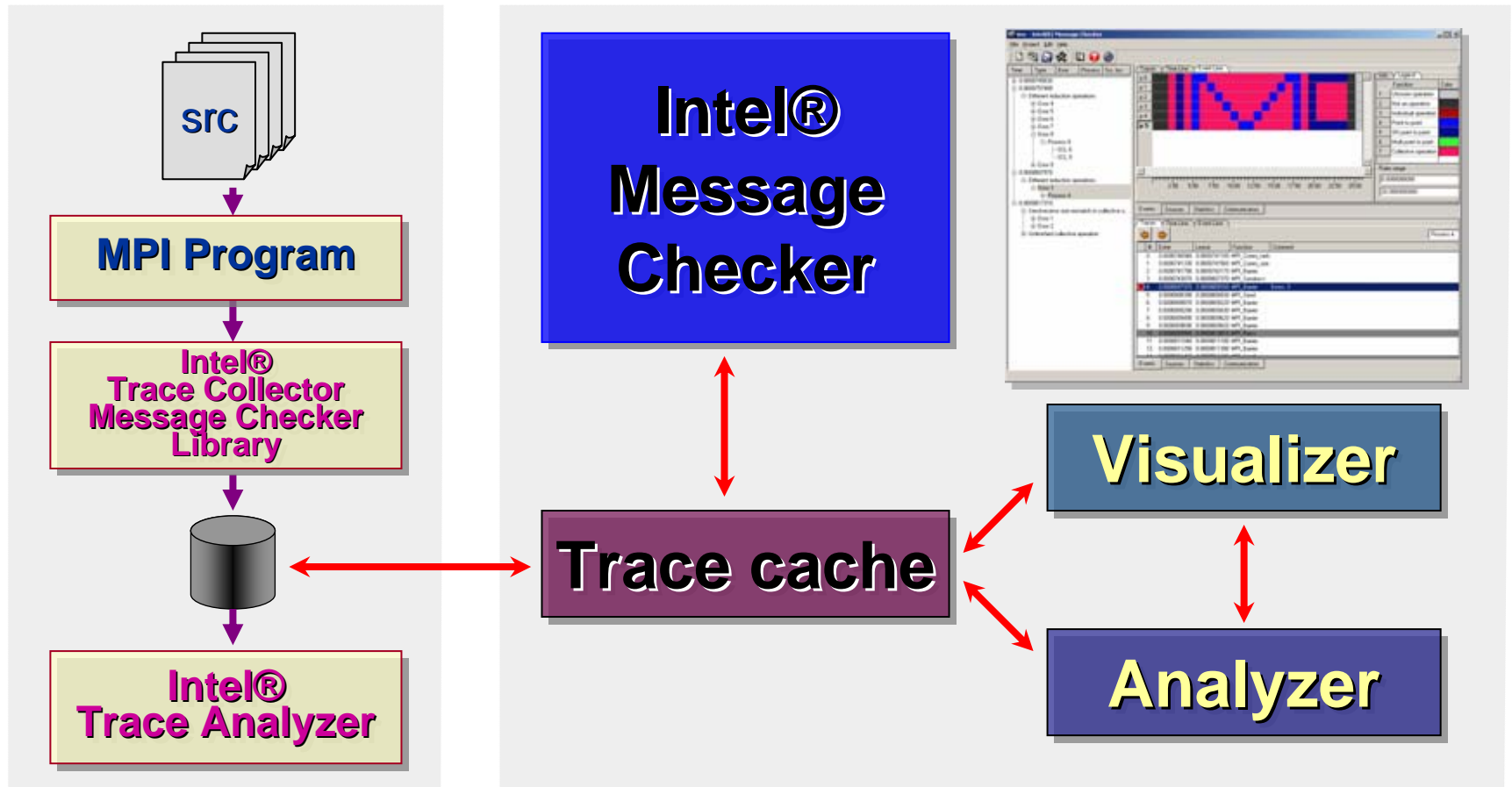
**After**

**Case 3 with optimized communication**

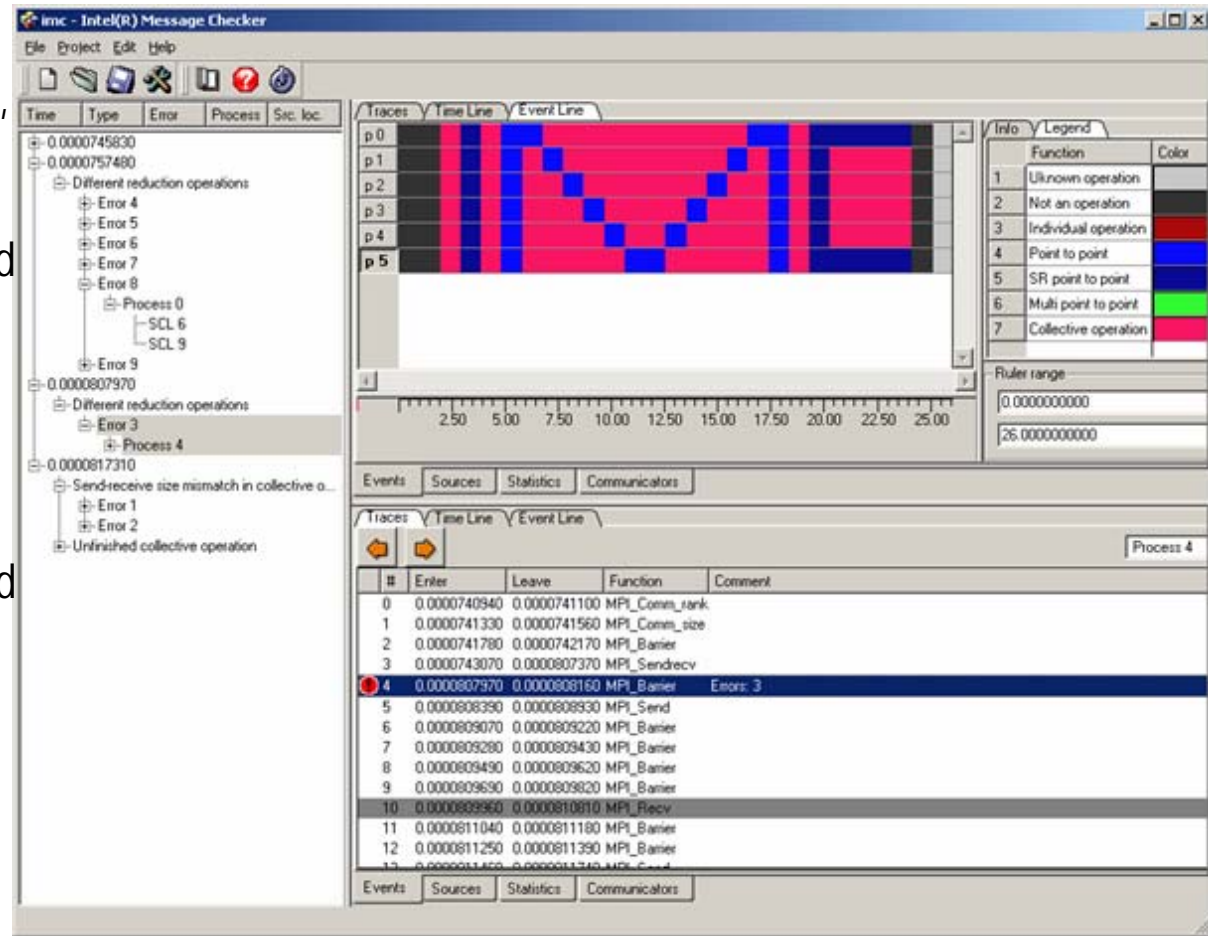# Intel® Message Checker
### A Unique Tool for MPI Correctness Analysis

## SSG
## Software Solutions Group

# Intel® Message Checker

src

**MPI Program**

↓

**Intel®
Trace Collector
Message Checker
Library**

↓

**Intel®
Trace Analyzer**

**Intel®
Message
Checker**

↕

**Trace cache** ↔ **Visualizer**
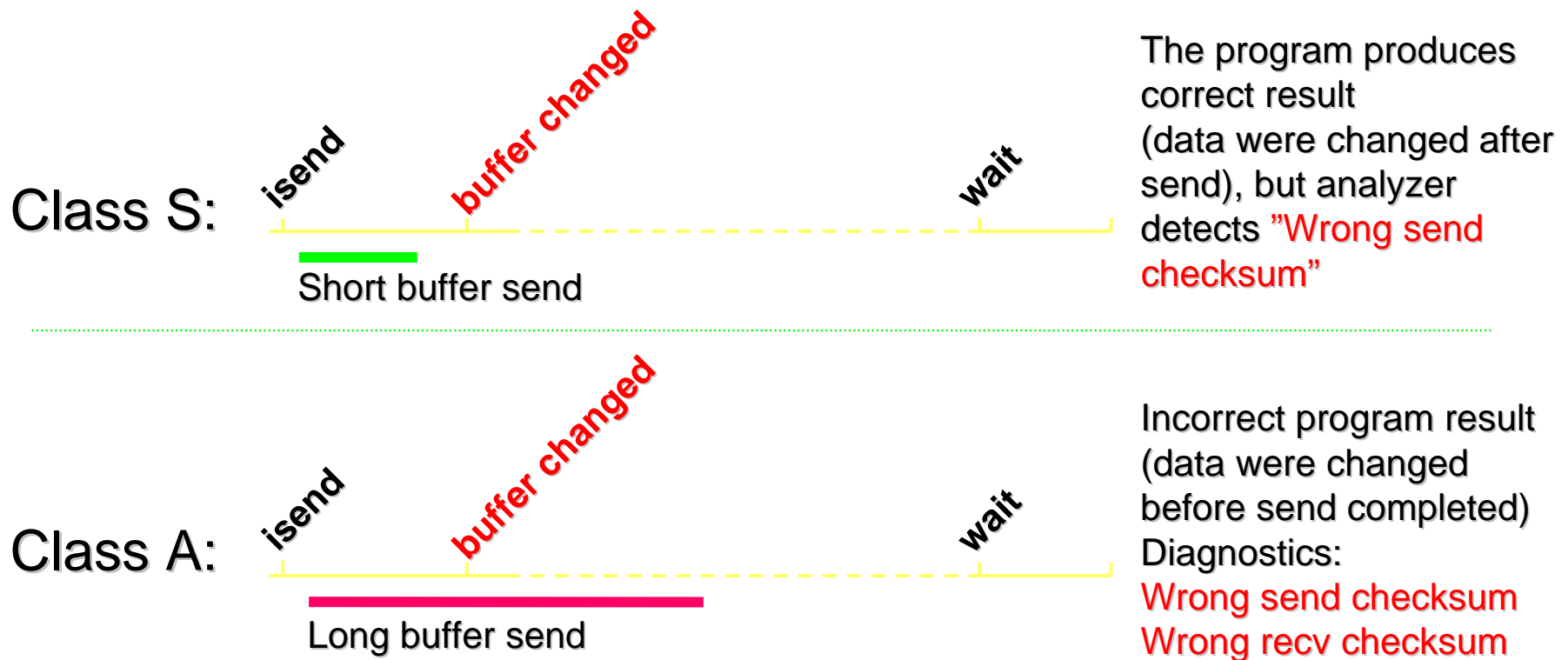
↕

**Analyzer**

**SSG**

(intel)

# Message Checker GUI

- Error view with different classifiers (#, type, process, time, source location)
- Source view with error context
- Trace view over time, events, and function calls
- Task statistics (general, communicator, and per-process)
- Communicator view
- Stack information
- Windows are synchronized based on user actions

# Example: NAS Parallel Benchmarks

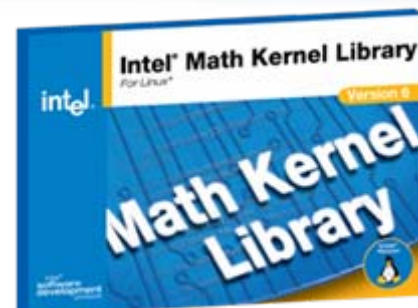Test was changed for buffer modification after non-blocking send

**Class S:**

isend — buffer changed — wait

Short buffer send

The program produces correct result (data were changed after send), but analyzer detects "Wrong send checksum"

**Class A:**

isend — buffer changed — wait

Long buffer send

Incorrect program result (data were changed before send completed)
Diagnostics:
Wrong send checksum
Wrong recv checksum

# Performance Libraries

**SSG**
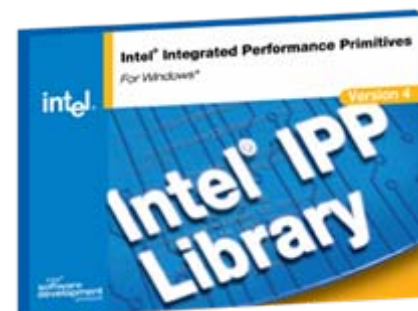**Software Solutions Group**

# Intel® Performance Libraries

- ## Intel® Math Kernel Library
  - **A set of highly optimized numerical processing functions for math, scientific, engineering and financial applications**
    - **Large-scale numerical applications**
    - **Workstations & server platform focus**

- ## Intel® Integrated Performance Primitives
  - **A library of signal, image, graphic, multimedia and numeric processing functions**
    - **Real-time and interactive applications**

- ## Intel® Optimized XML Library
  - **Optimized routines for XML parsing**
  - **To be released in Q1/2005**

- ## Intel® Computer Vision Library
  - **For Intel XScale® architecture**
  - **Available for free from developer.intel.com**