

# Advanced OpenMP: Tools

Dirk Schmidl

[schmidl@rz.rwth-aachen.de](mailto:schmidl@rz.rwth-aachen.de)

Parallel Programming Summer Course

31.07.2013 / Aachen

# OpenMP Tools

## ■ Intel Inspector XE

→ Overview

→ Live Demo

## ■ Intel Amplifier XE

→ Overview

→ Live Demo

→ Case Study: Conjugate Gradient Solver

## ■ Hands On Exercises

## Race Condition

- **Data Race: the typical OpenMP programming error, when:**
  - two or more threads access the same memory location, and
  - at least one of these accesses is a write, and
  - the accesses are not protected by locks or critical regions, and
  - the accesses are not synchronized, e.g. by a barrier.
- **Non-deterministic occurrence: e.g. the sequence of the execution of parallel loop iterations is non-deterministic and may change from run to run**
- **In many cases *private* clauses, *barriers* or *critical regions* are missing**
- **Data races are hard to find using a traditional debugger**
  - Use the *Intel Inspector XE*

## ■ Detection of

- Memory Errors
- Dead Locks
- Data Races

## ■ Support for

- Linux (32bit and 64bit) and Windows (32bit and 64bit)
- WIN32-Threads, Posix-Threads, Intel Threading Building Blocks and OpenMP

## ■ New Features

- Binary Instrumentation gives full functionality
- Independent stand-alone GUI for Windows and Linux
- memory error detection
- static security analysis

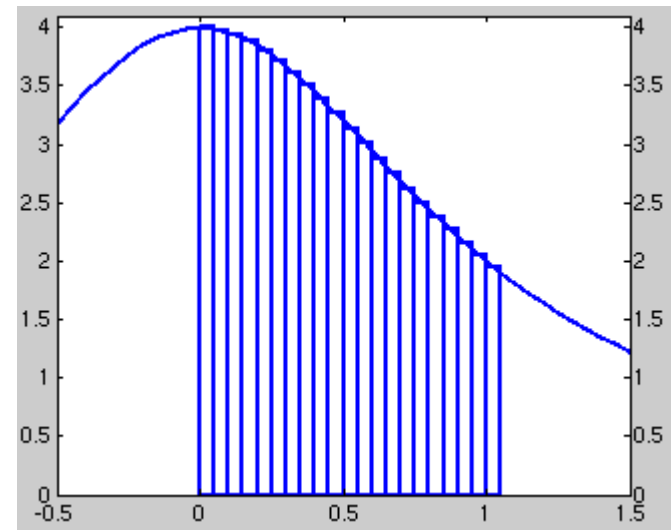
# PI Example Code

```
double f(double x)
{
    return (4.0 / (1.0 + x*x));
}
```

```
double CalcPi (int n)
{
    const double fH = 1.0 / (double) n;
    double fSum = 0.0;
    double fX;
    int i;

    for (i = 0; i < n; i++)
    {
        fX = fH * ((double)i + 0.5);
        fSum += f(fX);
    }
    return fH * fSum;
}
```

$$\pi = \int_0^1 \frac{4}{1+x^2}$$



# PI Example Code

```
double f(double x)
{
    return (4.0 / (1.0 + x*x));
}

double CalcPi (int n)
{
    const double fH = 1.0 / (double) n;
    double fSum = 0.0;
    double fX;
    int i;

    #pragma omp parallel for private(fX,i)
    for (i = 0; i < n; i++)
    {
        fX = fH * ((double)i + 0.5);
        fSum += f(fX);
    }
    return fH * fSum;
}
```

What did we  
forget?  
What will be  
the result?

# PI Example Code

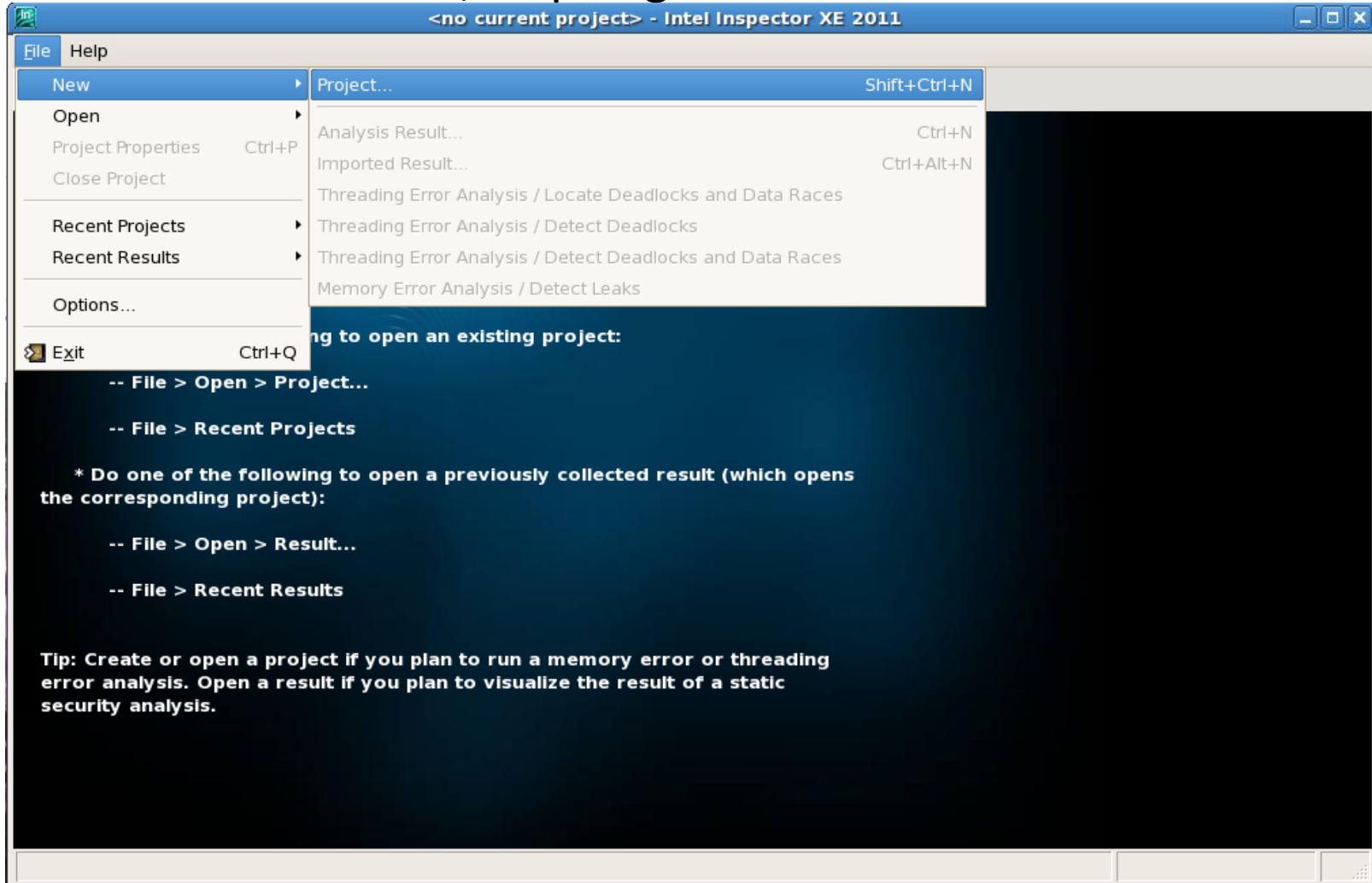
```
double f(double x)
{
    return (4.0 / (1.0 + x*x));
}
```

```
double CalcPi (int n)
{
    const double fH = 1.0 / (double) n;
    double fSum = 0.0;
    double fX;
    int i;
```

```
#pragma omp parallel for reduction(+:fSum)
for (i = 0; i < n; i++)
{
    fX = fH * ((double)i + 0.5);
    fSum += f(fX);
}
return fH * fSum;
}
```

What about  
this version?

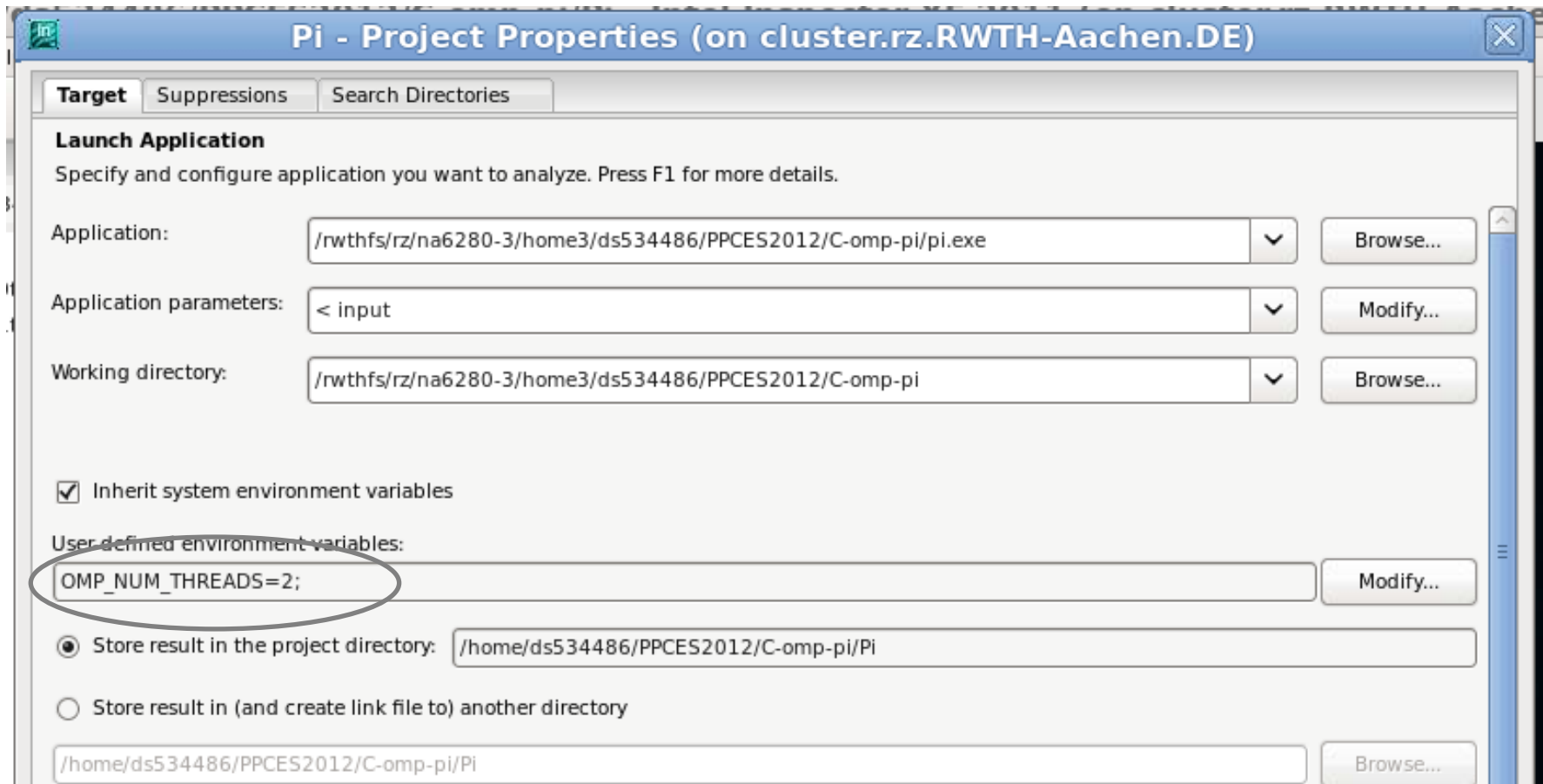
\$ module load intelixe ; inspxe-gui





# Inspector XE – Create Project

- ensure that multiple threads are used
- choose a real small dataset, execution time can grow 10X – 1000X

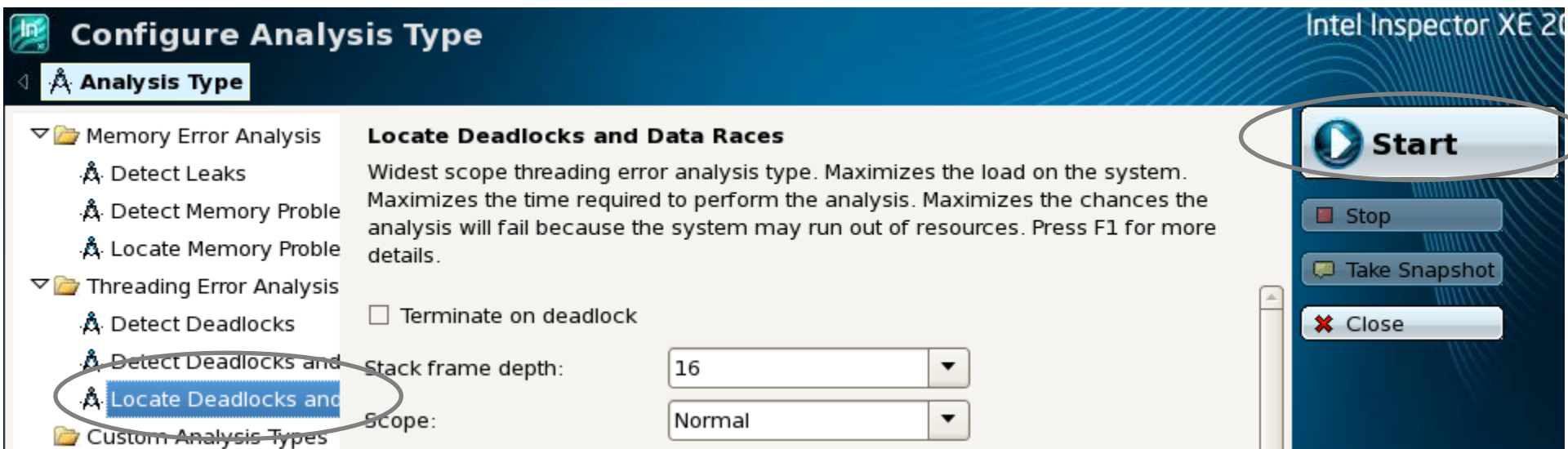


## Threading Error Analysis Modes

1. Detect Deadlocks
2. Detect Deadlocks and Data Races
3. Locate Deadlocks and Data Races



more details,  
more overhead



# Inspector XE – Results

- 1 detected problems
- 2 filters
- 3 code location

The screenshot shows the Intel Inspector XE 2011 interface. The main window is titled "Locate Deadlocks and Data Races". The "Problems" pane on the left shows a single problem, P1, which is a "Data race" detected in "pi.c" within the "pi.exe" module, with a state of "New". A yellow circle with the number "1" is placed over this problem entry.

The "Code Locations" pane below shows two locations for the problem: "X1 Read" and "X2 Write", both occurring in "pi.c:71" within the "CalcPi" function of "pi.exe". The code snippet for both locations is:

```

69 {
70     fX = fH * ((double)i + 0.5);
71     fSum += f(fX);
72 }
73     return fH * fSum;
    
```

The "X2 Write" location is highlighted with a yellow circle and the number "3".

The "Filters" pane on the right shows a summary of the results:

Severity	Count
Error	1 item(s)
<b>Problem</b>	
Data race	1 item(s)
<b>Source</b>	
pi.c	1 item(s)
<b>Module</b>	
pi.exe	1 item(s)
<b>State</b>	
New	1 item(s)
<b>Suppressed</b>	
Not suppressed	1 item(s)
<b>Investigated</b>	
Not investigated	1 item(s)

A yellow circle with the number "2" is placed over the "Investigated" section of the filters pane.

## 1 Timeline view

The screenshot displays the Intel Inspector XE 2011 interface. The main window title is "r001ti3". The top bar shows "Locate Deadlocks and Data Races" and "Intel Inspector XE 2011". Below the bar are tabs for "Target", "Analysis Type", "Collection Log", and "Summary".

The "Problems" table lists a single issue:

ID	Problem	Sources	Modules	State
P1	Data race	pi.c	pi.exe	New

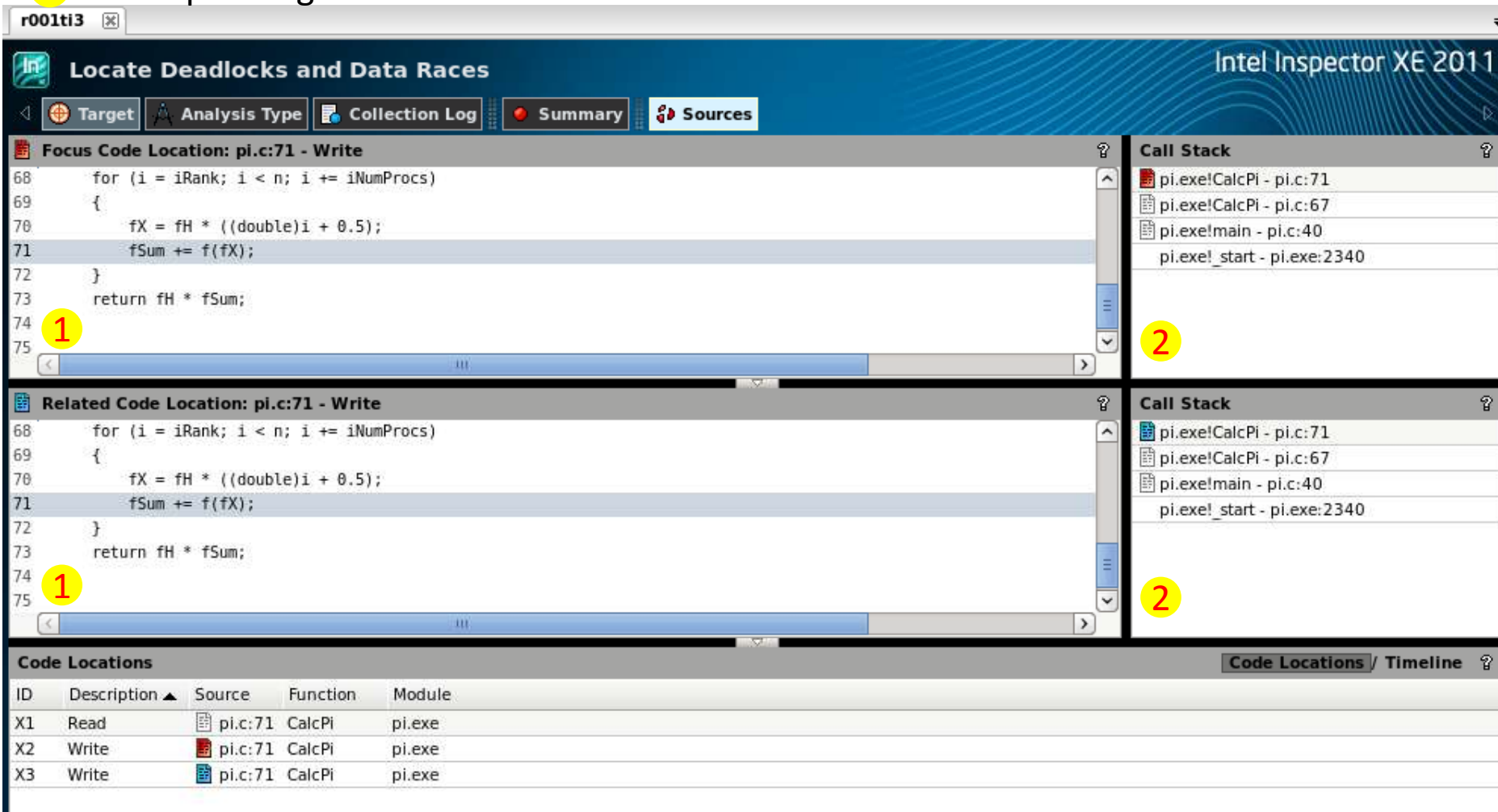
The "Filters" panel on the right shows the following counts:

Filter	Count
Severity: Error	1 item(s)
Problem: Data race	1 item(s)
Source: pi.c	1 item(s)
Module: pi.exe	1 item(s)
State: New	1 item(s)
Suppressed: Not suppressed	1 item(s)
Investigated: Not investigated	1 item(s)

The "Timeline" view at the bottom shows a horizontal axis from 0.05% to 0.4%. Two threads are visible: "main (16217)" and "OMP Worker ...". Blue bars represent the execution of each thread, with diamond markers indicating the data race event. A yellow circle with the number "1" is overlaid on the timeline view.

# Inspector XE – Results

- 1 Source Code producing the issue – double click opens an editor
- 2 Corresponding Call Stack



The screenshot shows the Intel Inspector XE 2011 interface. The top bar indicates the target is 'r001ti3' and the analysis type is 'Locate Deadlocks and Data Races'. The 'Summary' tab is selected, showing a red indicator for a deadlock.

The main window is divided into two panes, each showing source code and a call stack. Both panes are labeled with a yellow '1' next to the source code and a yellow '2' next to the call stack.

**Focus Code Location: pi.c:71 - Write**

```

68   for (i = iRank; i < n; i += iNumProcs)
69   {
70       fX = fH * ((double)i + 0.5);
71       fSum += f(fX);
72   }
73   return fH * fSum;
74
75

```

**Call Stack**

- pi.exe!CalcPi - pi.c:71
- pi.exe!CalcPi - pi.c:67
- pi.exe!main - pi.c:40
- pi.exe!\_start - pi.exe:2340

**Related Code Location: pi.c:71 - Write**

```

68   for (i = iRank; i < n; i += iNumProcs)
69   {
70       fX = fH * ((double)i + 0.5);
71       fSum += f(fX);
72   }
73   return fH * fSum;
74
75

```

**Call Stack**

- pi.exe!CalcPi - pi.c:71
- pi.exe!CalcPi - pi.c:67
- pi.exe!main - pi.c:40
- pi.exe!\_start - pi.exe:2340

**Code Locations**

ID	Description	Source	Function	Module
X1	Read	pi.c:71	CalcPi	pi.exe
X2	Write	pi.c:71	CalcPi	pi.exe
X3	Write	pi.c:71	CalcPi	pi.exe

# Inspector XE – Results

- 1 Source Code producing the issue – double click opens an editor
- 2 Corresponding Call Stack

The missing reduction is detected.

The screenshot displays the Intel Inspector XE 2017 interface. The top bar shows the target 'r001t13' and navigation options: Target, Analysis Type, Collection Log, Summary, and Sources. The main area is divided into three sections:

- Focus Code Location: pi.c:71 - Write**: Shows source code with line 71 highlighted. A yellow circle with '1' is next to line 71. The code includes a loop from line 68 to 72, and a return statement at line 73.
- Related Code Location: pi.c:71 - Write**: Shows the same source code as the focus location, with line 71 highlighted. A yellow circle with '1' is next to line 71.
- Call Stack**: Two call stacks are shown, one for each code location. Both call stacks are identical and show the following frames:
  - pi.exe!CalcPi - pi.c:71
  - pi.exe!CalcPi - pi.c:67
  - pi.exe!main - pi.c:40
  - pi.exe!\_start - pi.exe:2340
 A yellow circle with '2' is next to the top frame in both call stacks.

At the bottom, the **Code Locations** table is visible:

ID	Description	Source	Function	Module
X1	Read	pi.c:71	CalcPi	pi.exe
X2	Write	pi.c:71	CalcPi	pi.exe
X3	Write	pi.c:71	CalcPi	pi.exe

# Command Line Tool – inspxe-cl

## Threading Error Analysis Modes

1. Detect Deadlocks (ti1)
2. Detect Deadlocks and Data Races (ti2)
3. Locate Deadlocks and Data Races (ti3)

```
$ inspxe-cl -collect ti3 -- pi.exe < input
```

Data collection without GUI  
allows to use batch jobs.

```
$ inspxe-cl -report problems ...
```

Viewing results in text mode  
is helpful, when remote  
connections are slow.

# Amplifier XE

## ■ Performance Analyses for

- Serial Applications
- Shared Memory Parallel Applications

## ■ Sampling Based measurements

## ■ Features:

- Hot Spot Analysis
- Concurrency Analysis
- Wait
- Hardware Performance Counter Support

## Prerequisites:

- ssh -Y cluster-linux-tuning (to use Hardware Performance Counters)
- module load intelvtune



- Standard Benchmark to measure memory performance.
- Version is parallelized with OpenMP.

Measures Memory bandwidth for:

**y=x (copy)**

**y=s\*x (scale)**

**y=x+z (add)**

**y=x+s\*z (triad)**

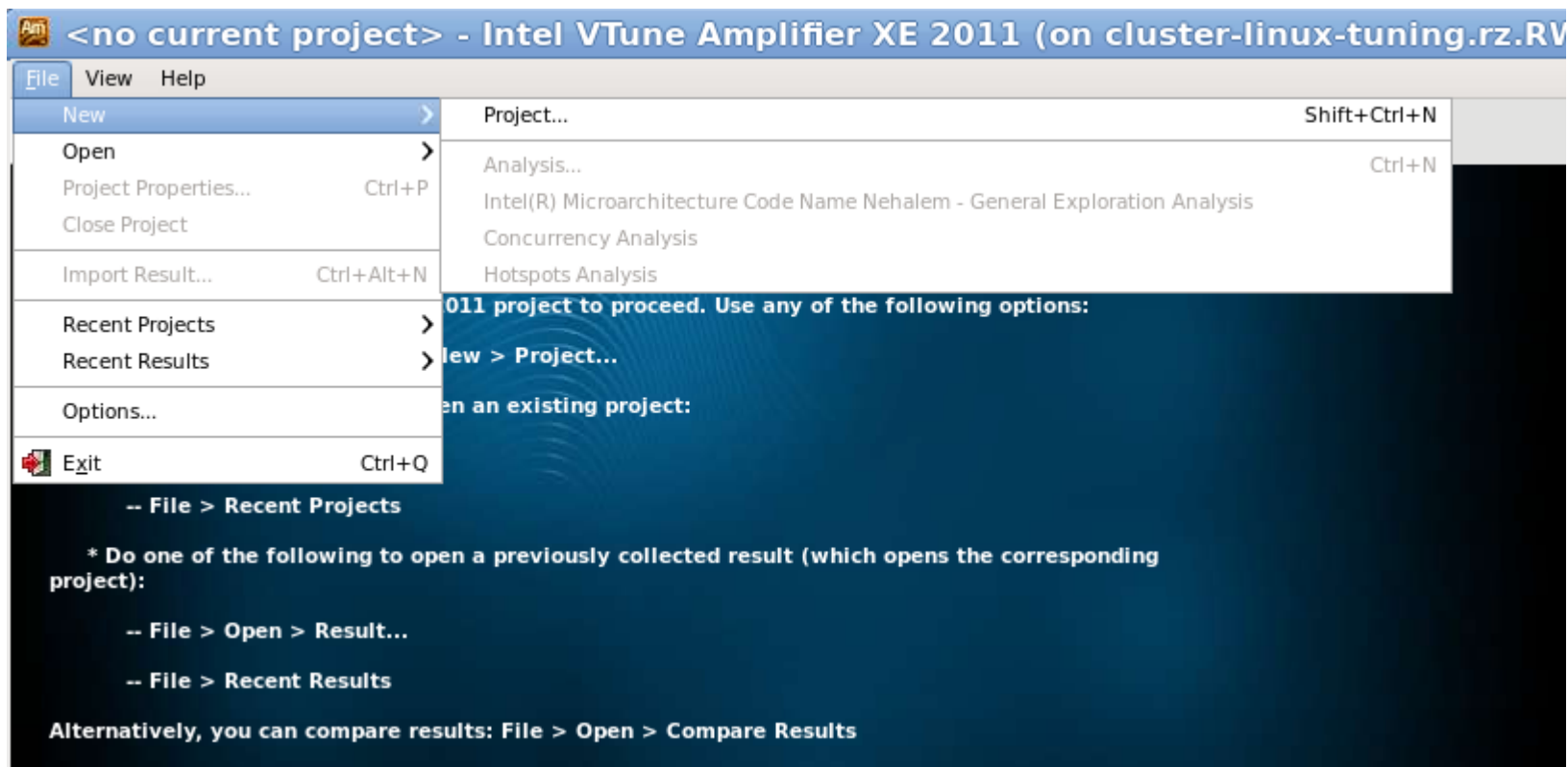
```
#pragma omp parallel for
for (j=0; j<N; j++)
    b[j] = scalar*c[j];
```

for double vectors x,y,z and scalar double value s

Function	Rate (MB/s)	Avg time	Min time	Max time
Copy:	33237.0185	0.0050	0.0048	0.0055
Scale:	33304.6471	0.0049	0.0048	0.0059
Add:	35456.0586	0.0070	0.0068	0.0073
Triad:	36030.9600	0.0069	0.0067	0.0072

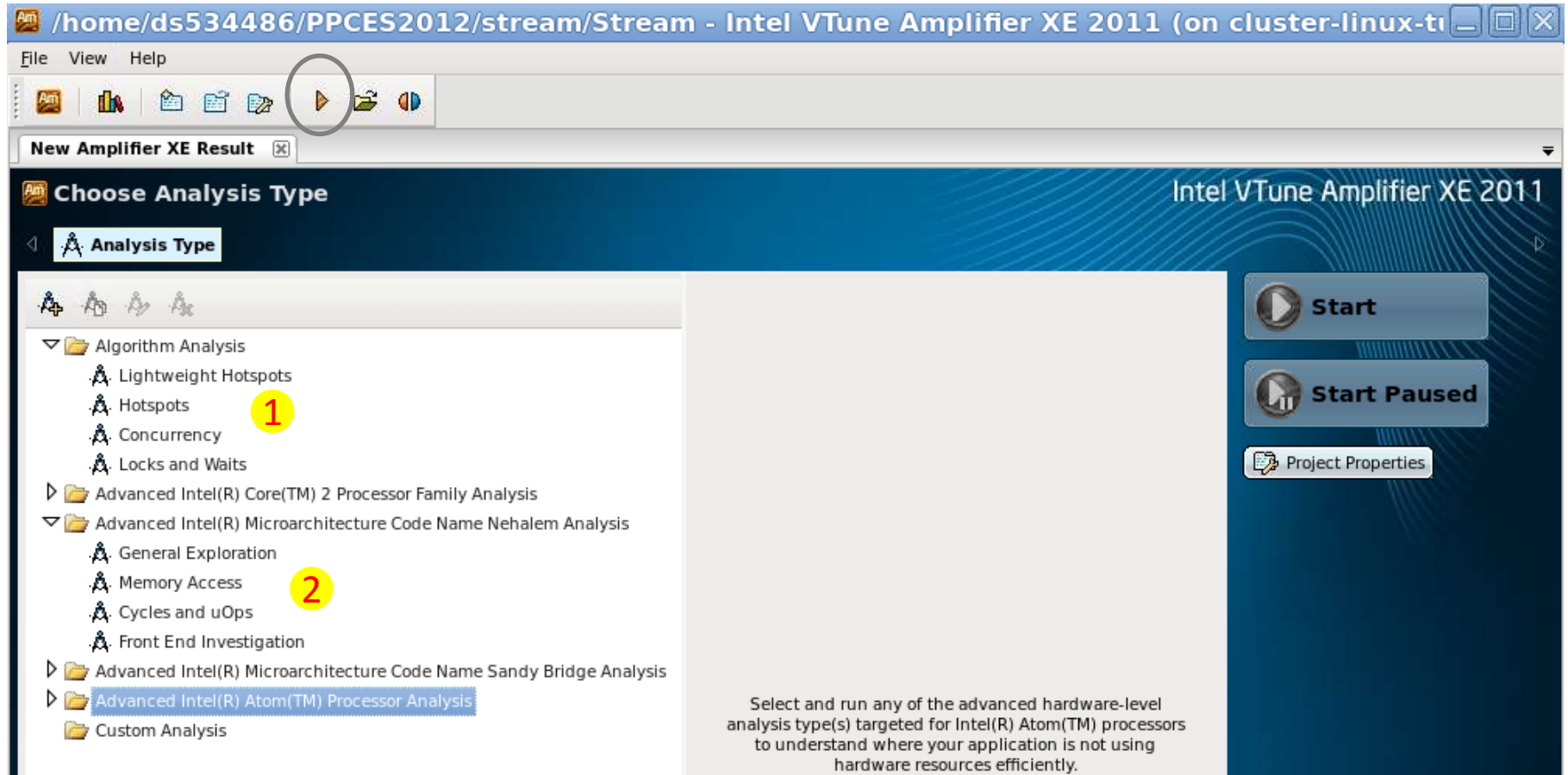
# Amplifier XE – Create Project

- Create a Project in the same way as with the inspector.
- Executable should be build with optimization.
- Use a reasonable sized data set.



# Amplifier XE – Measurement Runs

- 1 Basic Analysis Types
- 2 Hardware Counter based Analysis Types, choose Nehalem Architecture, here.  
(only available on cluster-linux-tuning)



The screenshot shows the 'Hotspots' application interface with the following sections:

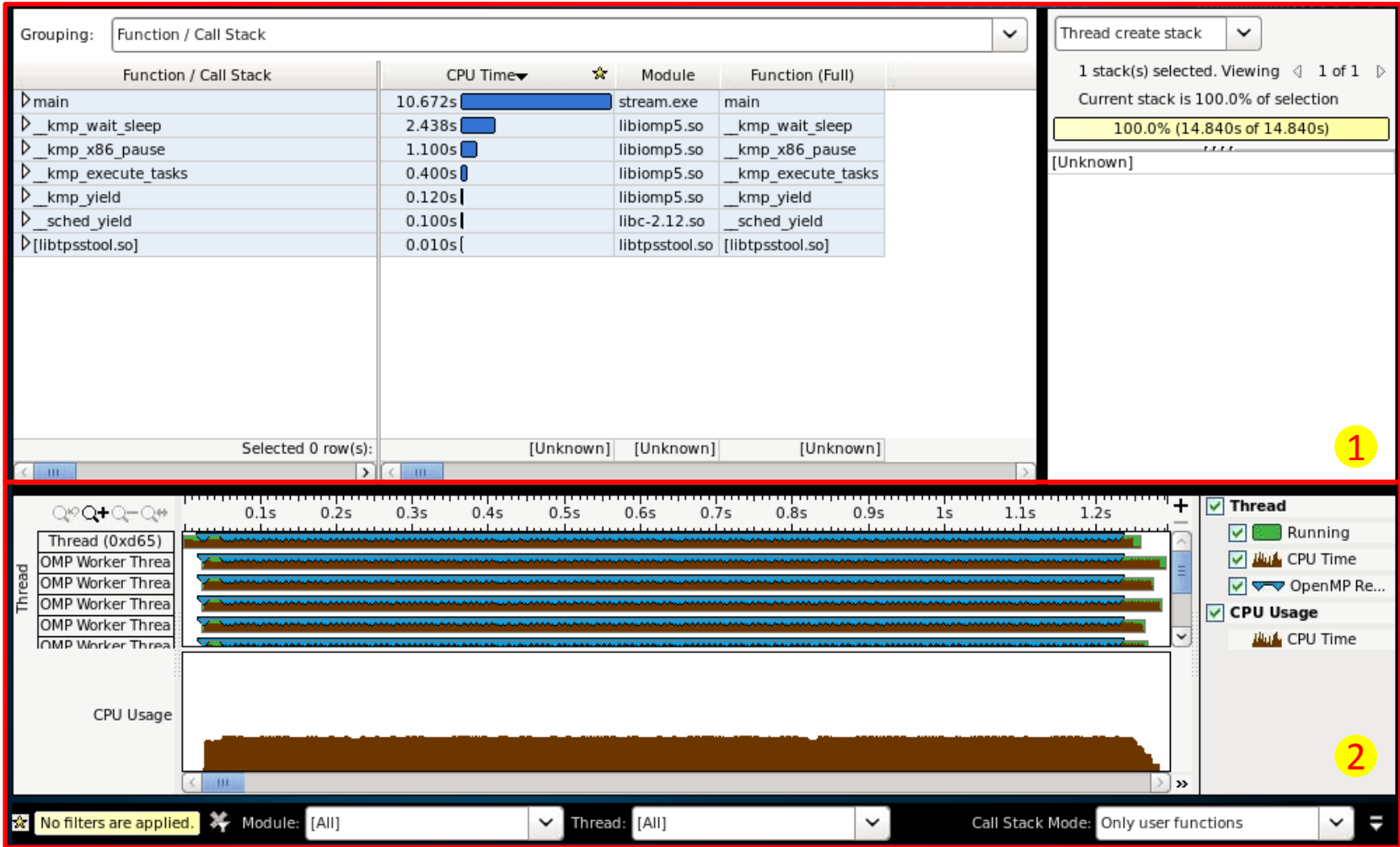
- Elapsed Time:** 1.294s. Sub-metrics: Total Thread Count: 12, CPU Time: 14.840s, Paused Time: 0s.
- Top Hotspots:** A table listing the most active functions.

Function	CPU Time
main	10.672s
__kmp_wait_sleep	2.438s
__kmp_x86_pause	1.100s
__kmp_execute_tasks	0.400s
__kmp_yield	0.120s
[Others]	0.110s
- Collection and Platform Info:** Command Line: /rwrhfs/rz/cluster/home/ds534486/PPCES2012/stream/stream.exe; Environment Variables: OMP\_NUM\_THREADS=12; Frequency: 3.07 GHz; Logical CPU Count: 24; Operating System: Linux; Computer Name: cluster-linux-tuning.rz.RWTH-Aachen.DE; Result Size: 71 KB.

Summary:

- 1 General Timing Information
- 2 Top Hotspots
- 3 Platform Information

## 1 Function Summary      2 Timeline View



# Amplifier XE – Hotspot Analysis

Double clicking on a function opens source code view.

- 1 Source Code View (only if compiled with -g)
- 2 Hotspot: Add Operation of Stream
- 3 Metrics View

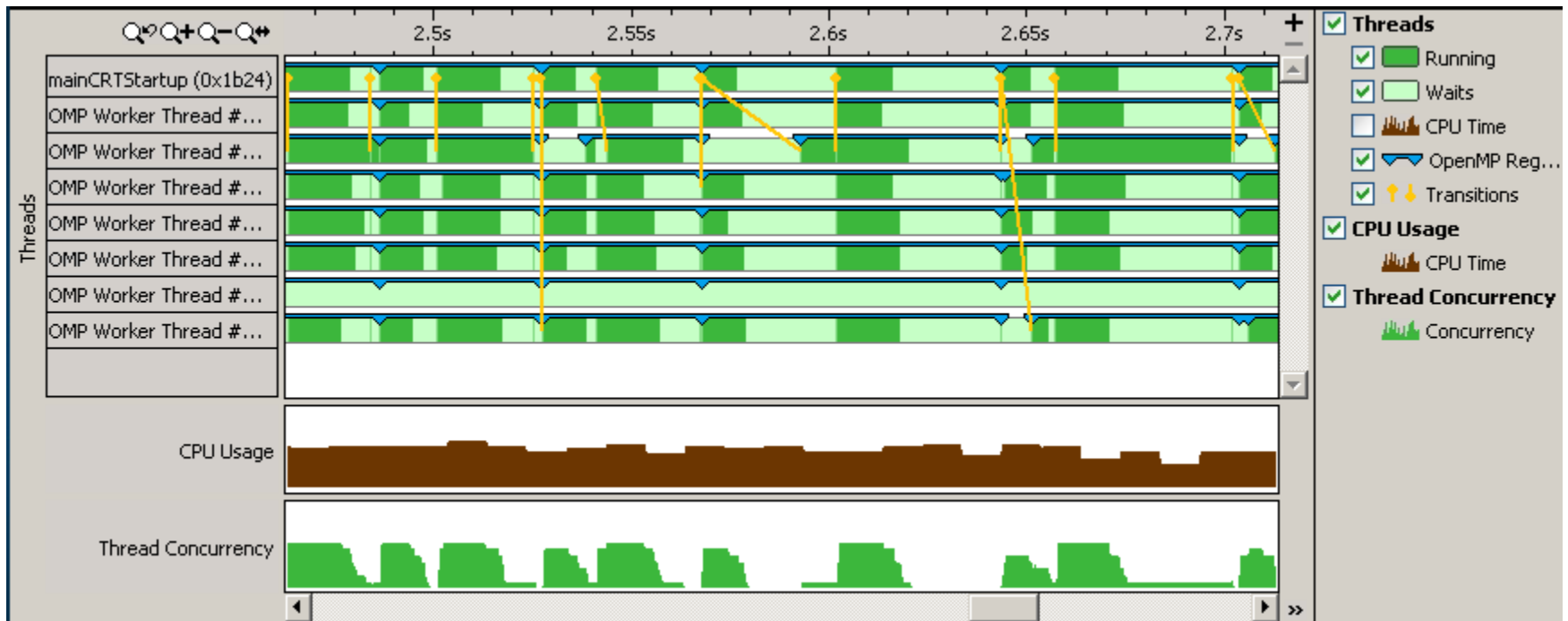
The screenshot displays the Amplifier XE interface with two main panels. The left panel shows the source code with line numbers and function calls. The right panel shows a 'CPU Time' table with a bar chart. A red box highlights the source code and the CPU time table. Three yellow circles with numbers 1, 2, and 3 are placed on the interface to indicate specific actions: 1 is on the source code, 2 is on the highlighted line, and 3 is on the CPU time table. Arrows point to the 'Hotspots' label on the right side of the CPU time table.

Line	Source	CPU Time
238	#else	
239	#pragma omp parallel for	0.010s
240	for (j=0; j<N; j++)	0.140s
241	c[j] = a[j]+b[j];	2.790s
242	#endif	
243	times[2][k] = mysecond() - times[2][k];	
244		
245	times[3][k] = mysecond();	
246	#ifdef TUNED	
247	tuned_STREAM_Triad(scalar);	
248	#else	
249	#pragma omp parallel for	0.160s
250	for (j=0; j<N; j++)	2.751s
251	a[j] = b[j]+scalar*c[j];	
252	#endif	
253	times[3][k] = mysecond() - times[3][k];	
254	}	
255		

Hotspots

# Amplifier XE – Locks and Waits Analysis

- Waiting time is shown in light green.
- Execution time is shown in dark green.
- CPU Usage and Thread Concurrency differ because waiting threads utilize a CPU.







## Derived Metrics

### ■ **Clock cycles per Instructions (CPI)**

- CPI indicates if the application is utilizing the CPU or not
- Take care: Doing “something” does not always mean doing “something useful”.

### ■ **Floating Point Operations per second (FLOPS)**

- How many arithmetic operations are done per second?
- Floating Point operations are normally really computing and for some algorithms the number of floating point operations needed can be determined.

## Amplifier XE – Hardware Counter

- **Only single user at a time per system.**
- **different predefined Analysis types (use “Nehalem” on cluster-linux-tuning)**
  - General Exploration
  - Read Bandwidth
  - Write Bandwidth
  - Memory Access
  - Cycles and uOps
  - Front End Investigation

**Hardware Counters provide very detailed information, but they are often hard to understand and interpret for non experts.**

# Amplifier XE – Hardware Counter

**1** **CPI rate (Clock cycles per instruction):** In theory modern processors can finish 4 instructions in 1 cycle, so a CPI rate of 0.25 is possible. A value between 0.25 and 1 is often considered as good for HPC applications.

**Elapsed Time: 1.872s**

Hardware Event Count:	125,574,000,000
CPU_CLK_UNHALTED.THREAD:	6.3462e+10
INST_RETIRED.ANY:	6.2112e+10

**1** **CPI Rate: 1.022**

The CPI may be too high. This could be caused by issues such as memory instructions. Explore the other hardware-related metrics to identify what

**Retire Stalls: 0.570s**

A high number of retire stalls is detected. This may result from branch mispredict issues. Use this metric to find where you have stalled instructions. Once

**LLC Miss: 0.013s**

**LLC Load Misses Serviced By Remote DRAM: 0.001s**

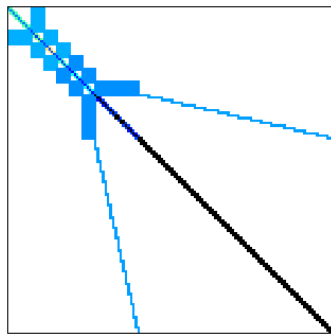
**Instruction Starvation: 0.098s**

**Branch Mispredict: 0.001s**

**Execution Stalls: 0.288s**

## ■ Sparse Linear Algebra

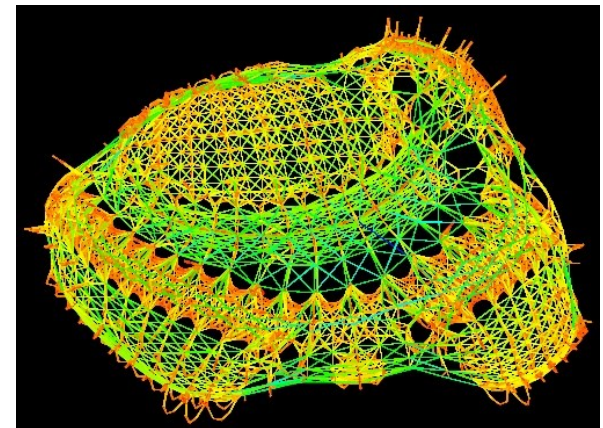
- Sparse Linear Equation Systems occur in many scientific disciplines.
- Sparse matrix-vector multiplications (SpMxV) are the dominant part in many iterative solvers (like the CG) for such systems.
- number of non-zeros  $\ll n*n$



### Beijing Botanical Garden

Oben Rechts: Original Gebäude  
Unten Rechts: Modell  
Unten Links: Matrix

(Quelle: Beijing Botanical Garden and University of Florida, Sparse Matrix Collection)



## Case Study: CG

- $A = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & 2 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 4 & 0 & 4 & 4 \end{pmatrix}$

- **Format: two dimensional array**
- **double A[4][4];**
- **stores n\*n doubles (16 here)**

## Case Study: CG

- $A = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & 2 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 4 & 0 & 4 & 4 \end{pmatrix}$

- **Format: coordinate list**

- **stores for all non-zeros (row, column, value)**

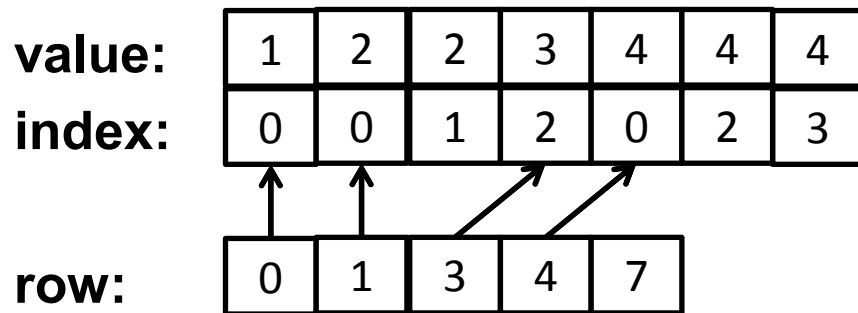
- **(0,0,1) (1,0,2) (1,1,2) (2,2,3) (3,0,4) (3,3,4) (3,4,4)**

- **nnz \* (2\*int + 1\*double)**

## Case Study: CG

■ 
$$A = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & 2 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 4 & 0 & 4 & 4 \end{pmatrix}$$

- **Format: compressed row storage**
- **store all values and columns in arrays (length nnz)**
- **store beginning of a new row in a third array (length n+1)**

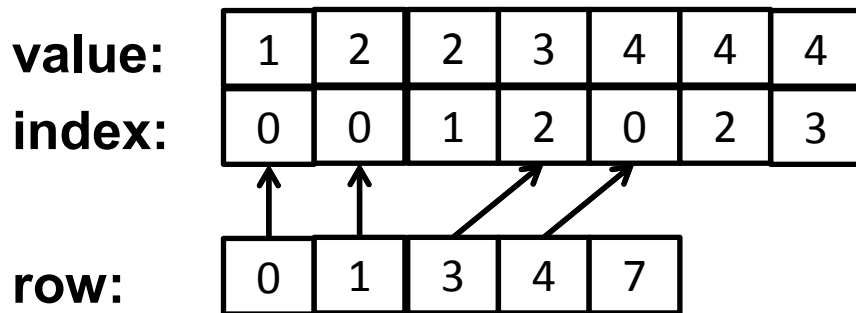


## Case Study: CG

$$A = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & 2 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 4 & 0 & 4 & 4 \end{pmatrix}$$

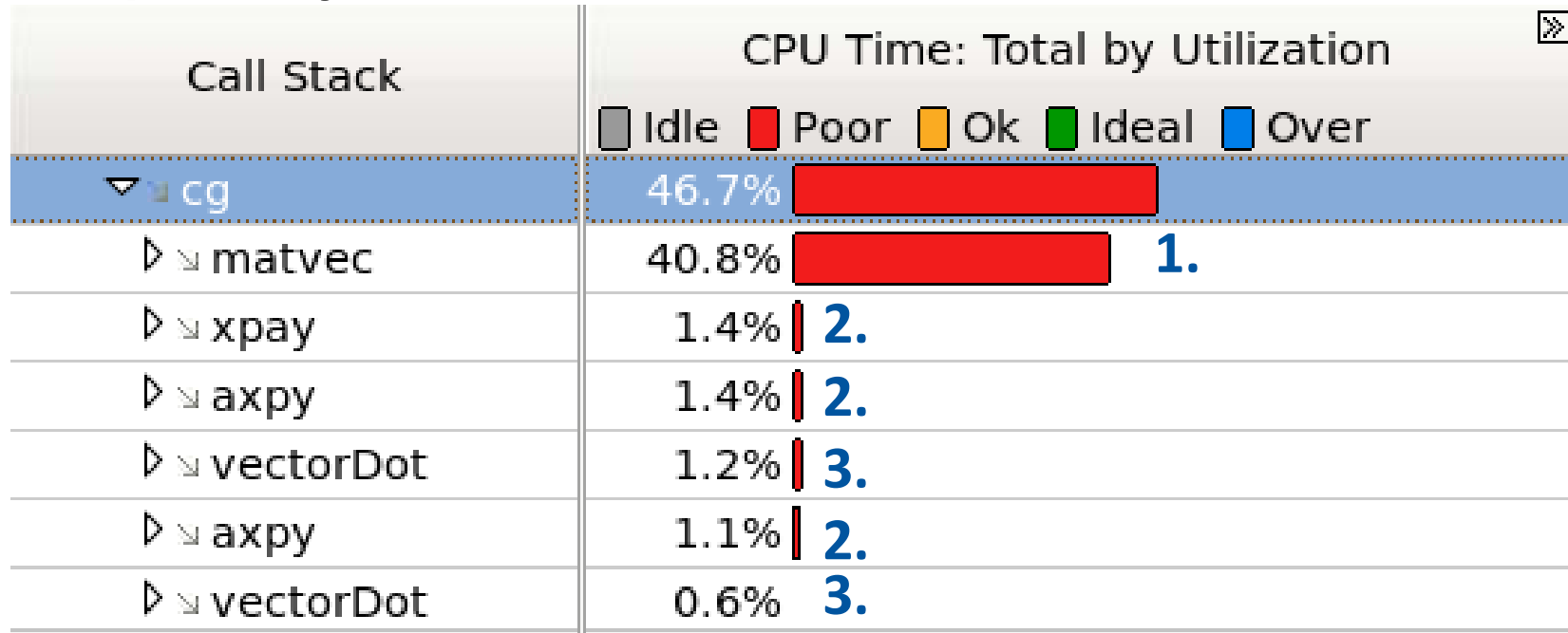
```
for (i = 0; i < num_rows; i++){
    sum = 0.0;
    for (nz=Arow[i]; nz<Arow[i+1]; ++nz){
        sum+= Aval[nz]*x[Acol[nz]];
    }
    y[i] = sum;
}
```

- **Format: compressed row storage**
- **store all values and columns in arrays (length nnz)**
- **store beginning of a new row in a third array (length n+1)**





## Hotspot analysis of the serial code:

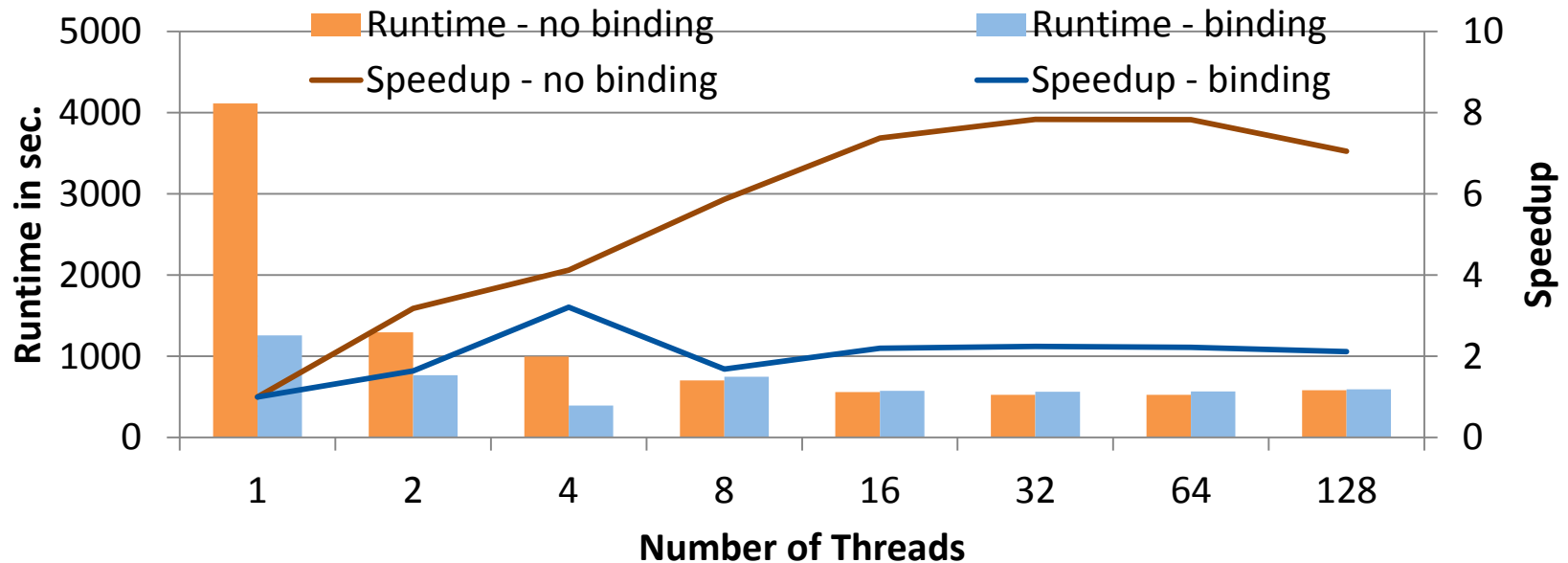


## Hotspots are:

1. matrix-vector multiplication
2. scaled vector additions
3. dot product

## Tuning:

- parallelize all hotspots with a parallel for construct
- use a reduction for the dot-product
- activate thread binding



## Hotspot analysis of naive parallel version:

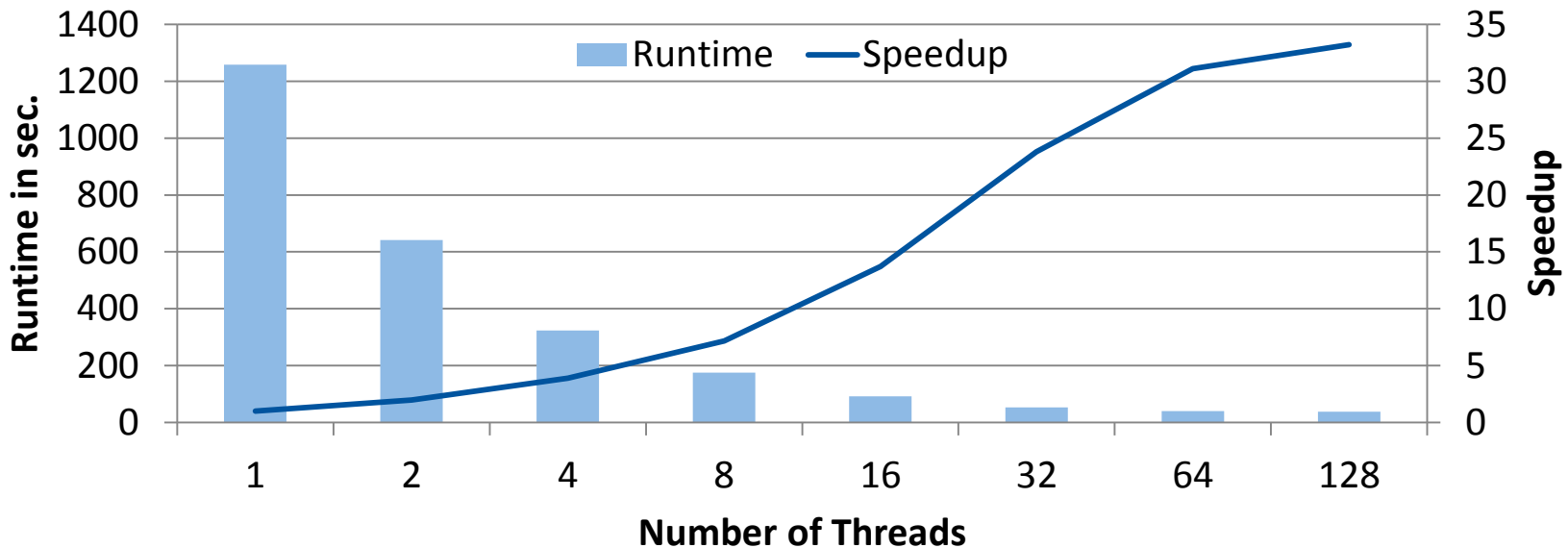
Event Name
MEM_UNCORE_RETIRED.LOCAL_DRAM_AND_REMOTE_CACHE_HIT
MEM_UNCORE_RETIRED.REMOTE_DRAM

A lot of remote accesses occur in nearly all places.

	MEM_UNCORE_RETIRED.LOCAL_...	MEM_UNCORE_RETIRED.REMOTE...
void matvec(const int n, const int int i,j;		
#pragma omp parallel for private(j)	20,000	0
for(i=0; i<n; i++){	0	0
y[i]=0;	0	0
for(j=ptr[i]; j<ptr[i+1]; j	6,740,000	3,720,000
y[i]+=value[j]*x[index[	17,580,000	6,680,000
}		
}		

## Tuning:

- Initialize the data in parallel
- Add parallel for constructs to all initialization loops



- Scalability improved a lot by this tuning on the large machine.

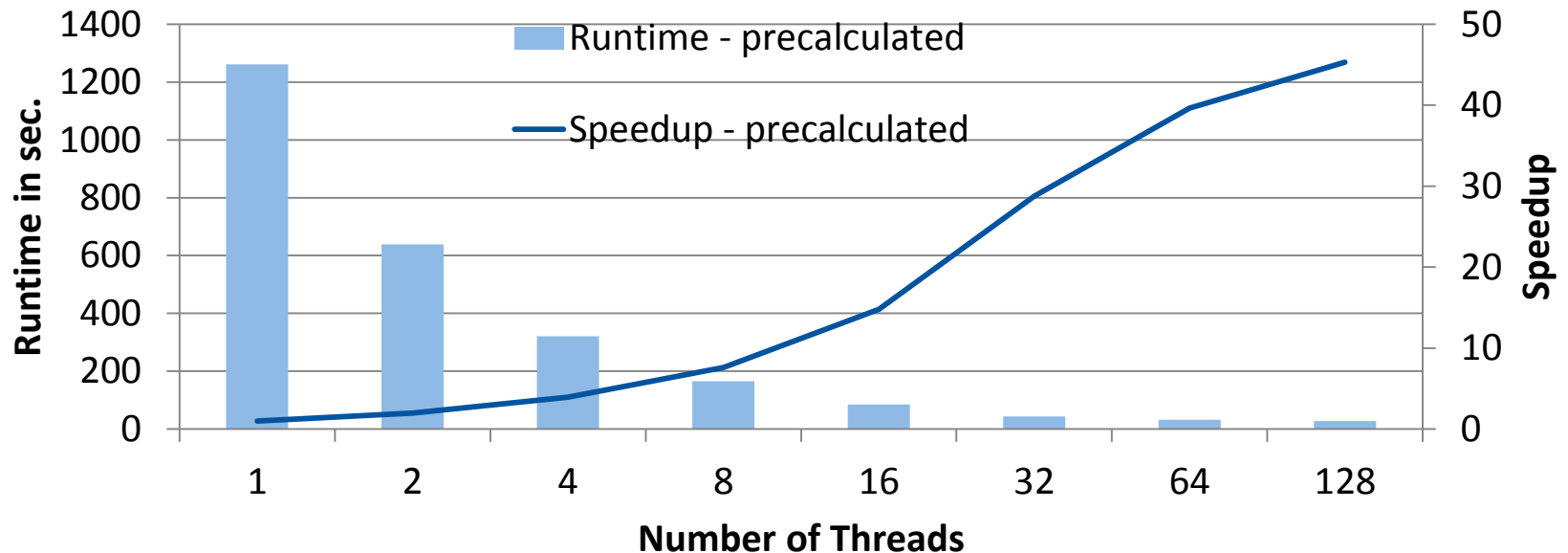
- Analyzing load imbalance in the concurrency view:

So.. Line	Source	CPU Time: Total by... Idle Poor Ok Ic	Ove... and...
49	void matvec(const int n, const int nnz,		
50	int i,j;		
51	#pragma omp parallel for private(j)	22.462s	10.612s
52	for(i=0; i<n; i++){	0.050s	0s
53	y[i]=0;	0.060s	0s
54	for(j=ptr[i]; j<ptr[i+1]; j++){	1.741s	0s
55	y[i]+=value[j]*x[index[j]];	9.998s	0s

- 10 seconds out of ~35 seconds are overhead time
- other parallel regions which are called the same amount of time only produce 1 second of overhead

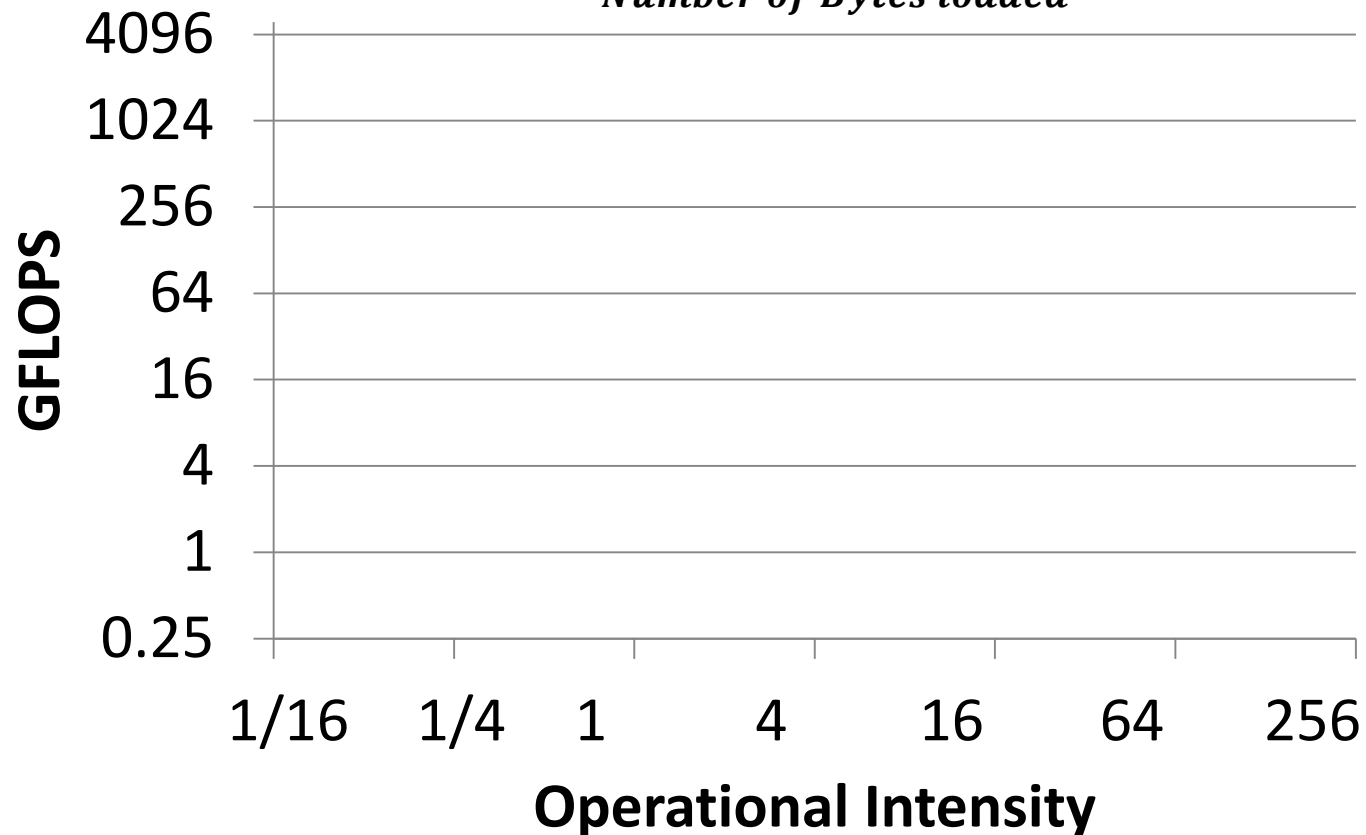
## ■ Tuning:

→ pre-calculate a schedule for the matrix-vector multiplication, so that the non-zeros are distributed evenly instead of the rows



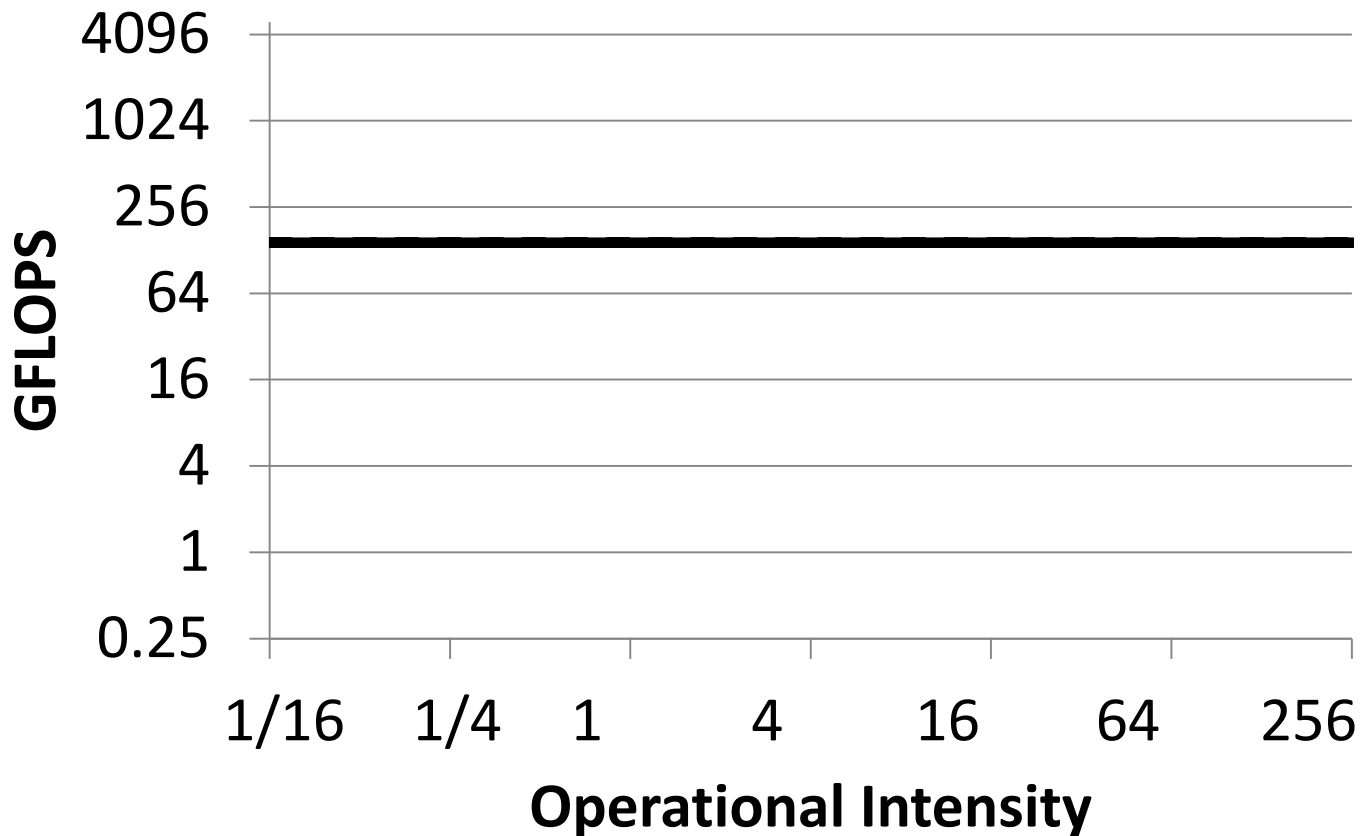
## When to stop? Roofline Model

- The Operational Intensity indicates if an application is bound by the peak performance or memory bandwidth.
- Operational Intensity =  $\frac{\text{Number of Floating Point Operations}}{\text{Number of Bytes loaded}}$



# When to stop? Roofline Model

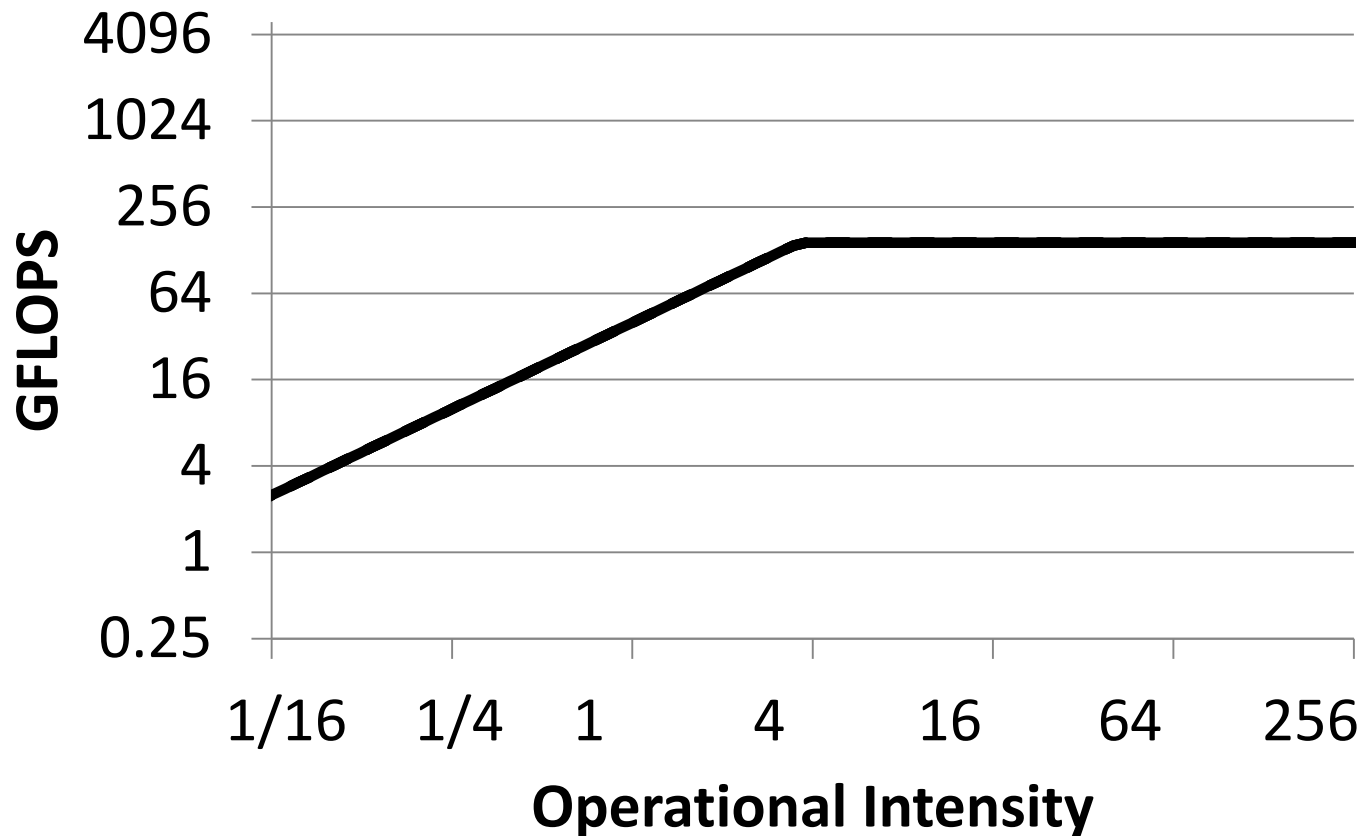
- Peak performance of a 2-socket Intel Westmere-EP (3 GHz) is:  
 $3 \text{ (GHz)} * 12 \text{ (cores)} * 2 \text{ (vector length)} * 2 \text{ (fused multiply-add)}$   
**= 144 GFLOPS**





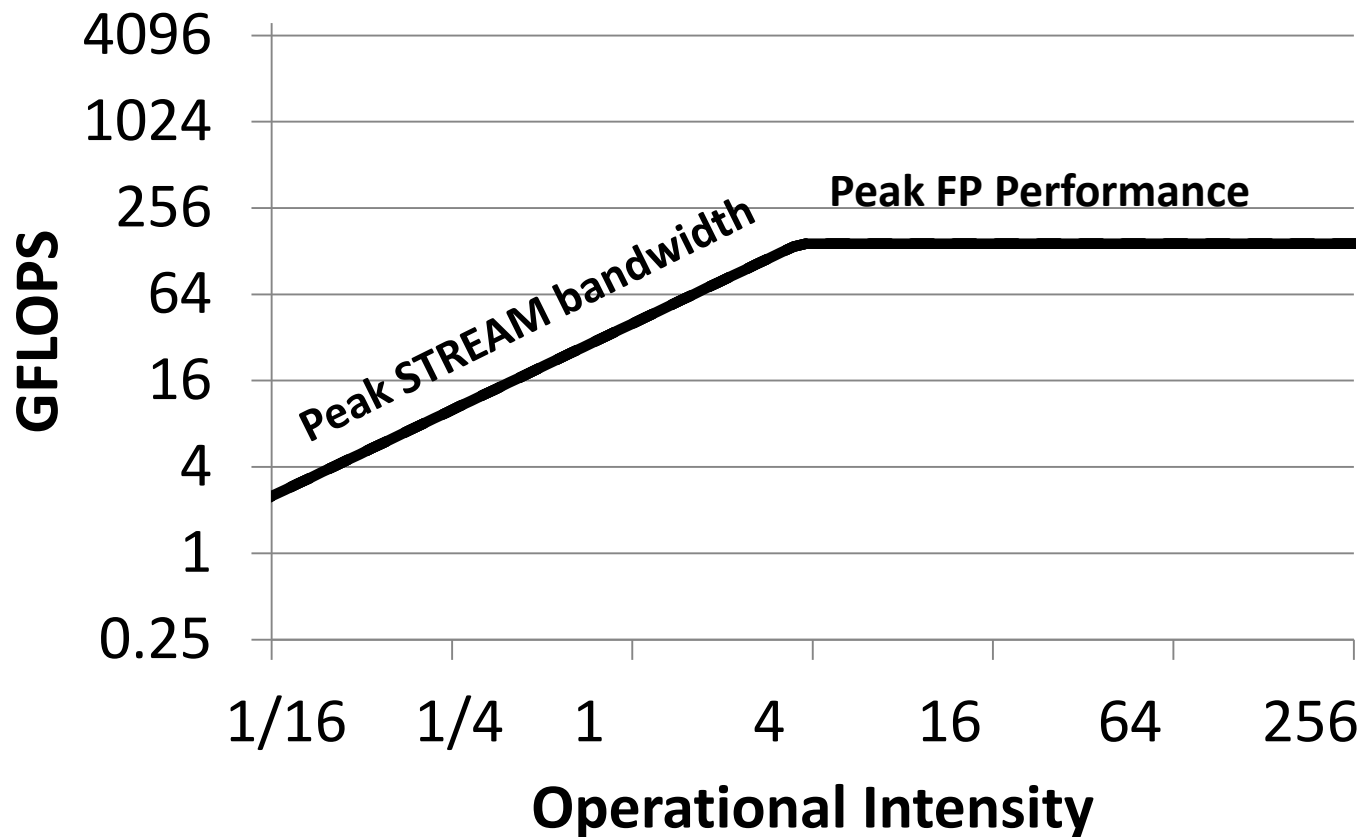
# When to stop? Roofline Model

- Memory bandwidth measured with the STREAM benchmark is about 40 GB/s (Triad:  $\vec{a} = \vec{b} + \alpha * \vec{c}$ )

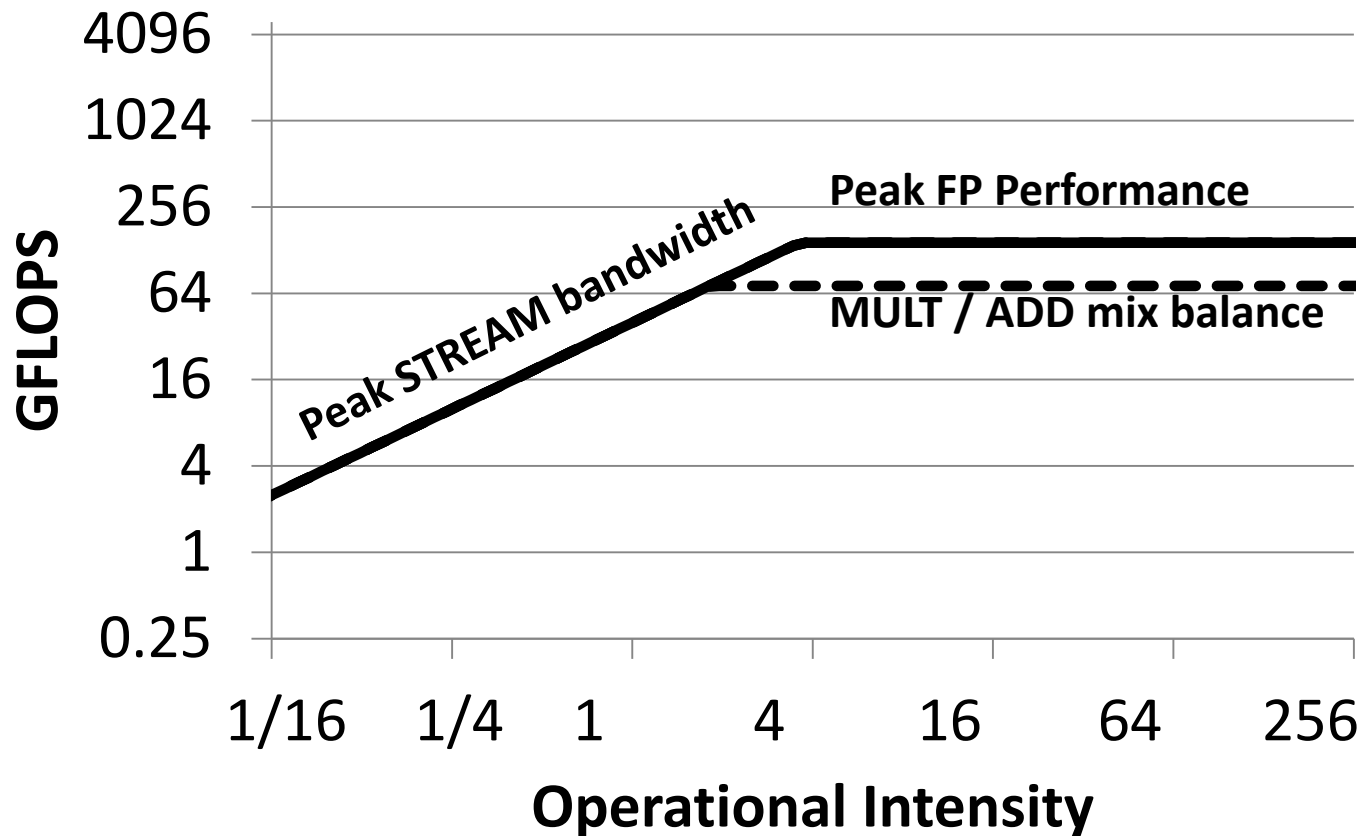


## When to stop? Roofline Model

- The “Roofline” is the peak performance depending on the algorithms' “operational intensity”.

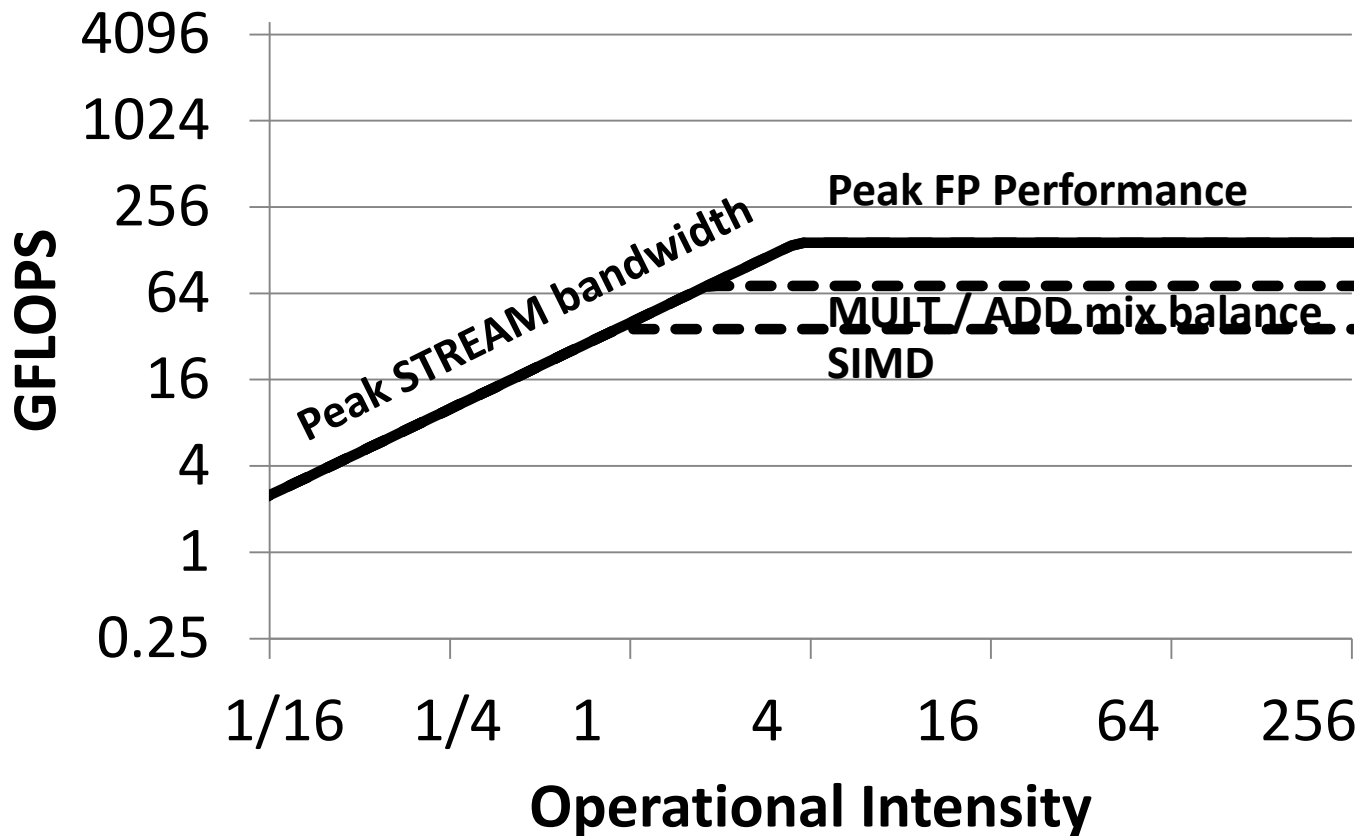


- To reach the peak performance an even mix of multiply and add operations is need (“fused multiply add”)



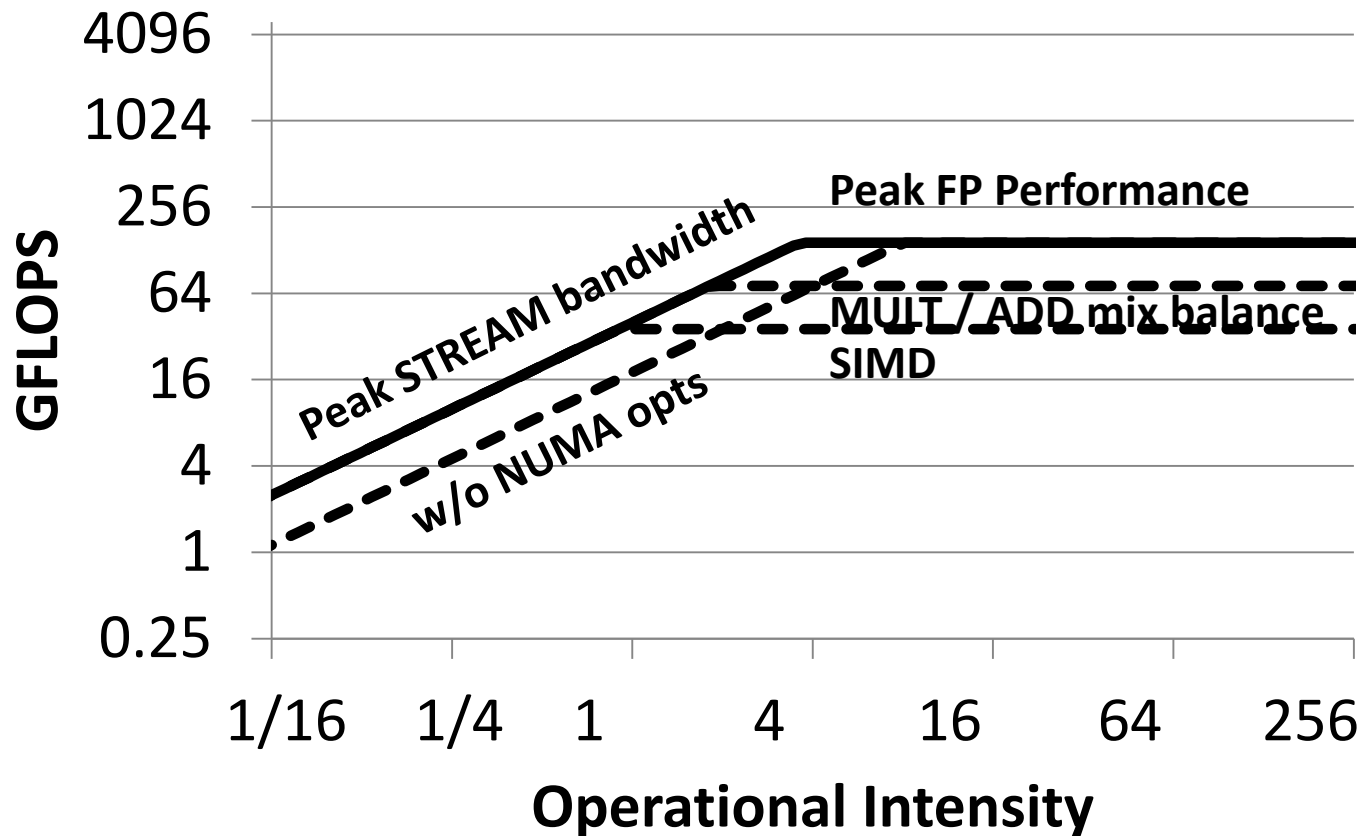
# When to stop? Roofline Model

- Without SIMD vectorization only  $\frac{1}{4}$  of the peak performance is achievable



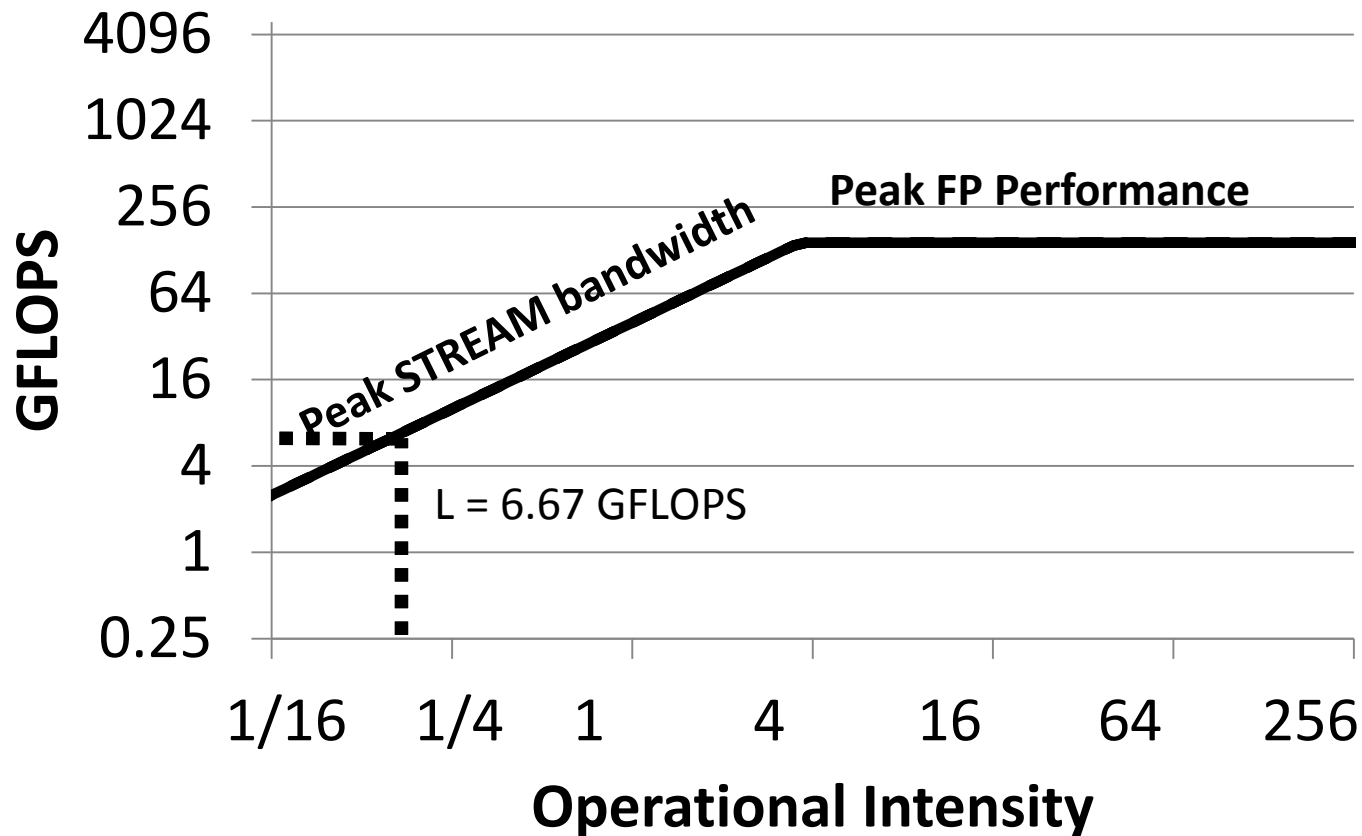
# When to stop? Roofline Model

- Peak STREAM Performance is only achievable if data is distributed optimally across NUMA nodes.



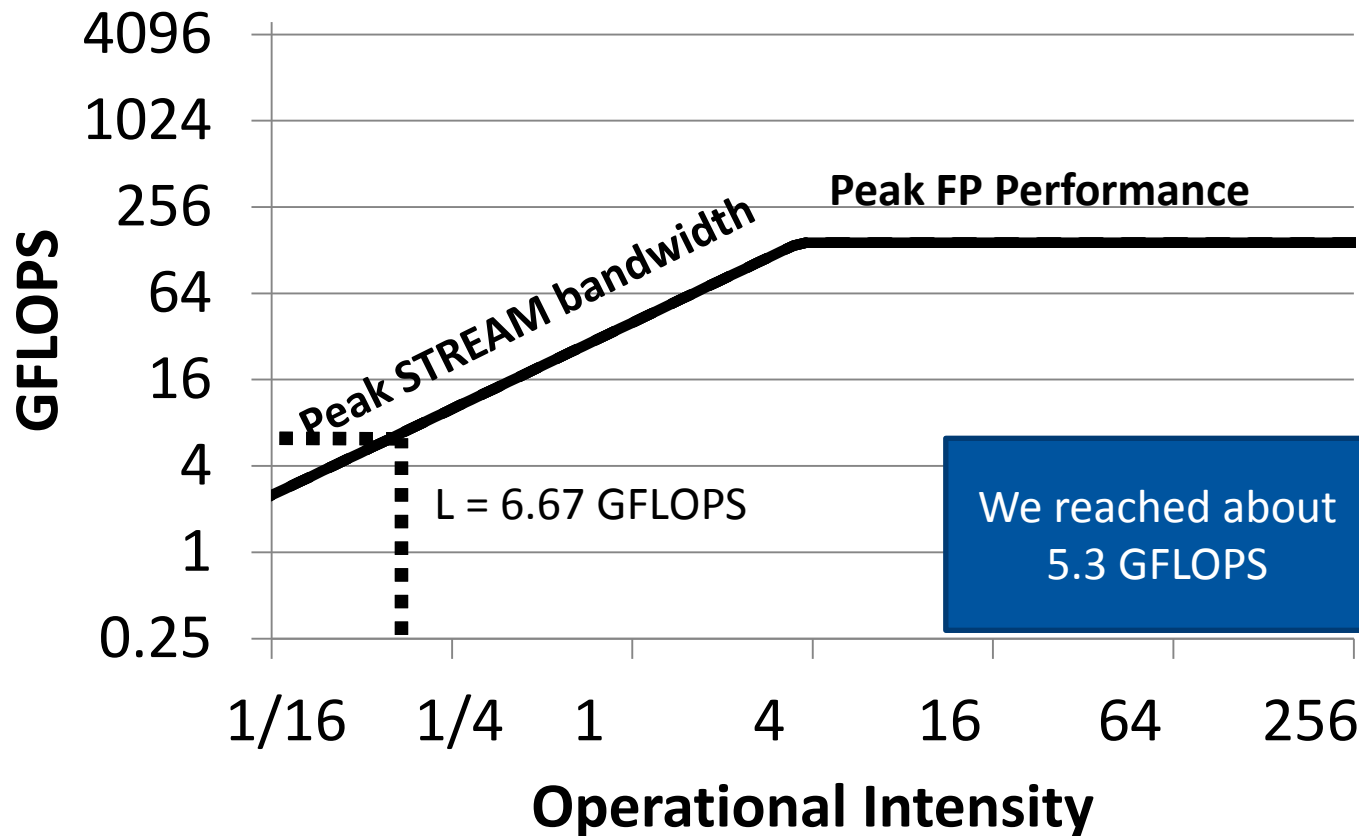
# When to stop? Roofline Model

- vectors can be kept in the cache
- every matrix element must be loaded for one multiply and add operation
- one element is a value (double) and an index (int)
- Operational intensity  $O = \frac{2 \text{ FLOPS}}{12 \text{ Byte}} = \frac{1 \text{ FLOPS}}{6 \text{ Byte}}$



# When to stop? Roofline Model

- vectors can be kept in the cache
- every matrix element must be loaded for one multiply and add operation
- one element is a value (double) and an index (int)
- Operational intensity  $O = \frac{2 \text{ FLOPS}}{12 \text{ Byte}} = \frac{1 \text{ FLOPS}}{6 \text{ Byte}}$



## Correctness:

- **Data Races are very hard to find, since they do not show up every program run.**
- **Intel Inspector XE helps a lot in finding these errors.**
- **Use really small datasets, since the runtime explores.**
- **If possible use non optimized code.**

## Performance:

- **Start with simple performance measurements like hotspots analyses and then focus on these hot spots.**
- **In OpenMP applications analyze the waiting time of threads. Is the waiting time balanced?**
- **Hardware counters might help for a better understanding of an application, but they are hard to interpret.**