



Perfilamento com Intel Advisor e Vtune

Silvio Stanzani , Raphael Cóbe , Rogério Iope,
Jefferson Fialho

Núcleo de Computação Científica – UNESP

ERAD-SP 2019

Agenda

- Parallel Architectures
- Roofline Model
- Optimization Workflow
- Intel Advisor
- Intel Vtune

Parallel Architectures

- Heterogeneous computational systems:
 - Scalar and Vector Instructions

Vector Instructions (SIMD)								Scalar Instructions	
A7	A6	A5	A4	A3	A2	A1	A0	A	
+								+	
B7	B6	B5	B4	B3	B2	B1	B0	B	
=								=	
A7+B7	A6+B6	A5+B5	A4+B4	A3+B3	A2+B2	A1+B1	A0+B0	A+B	

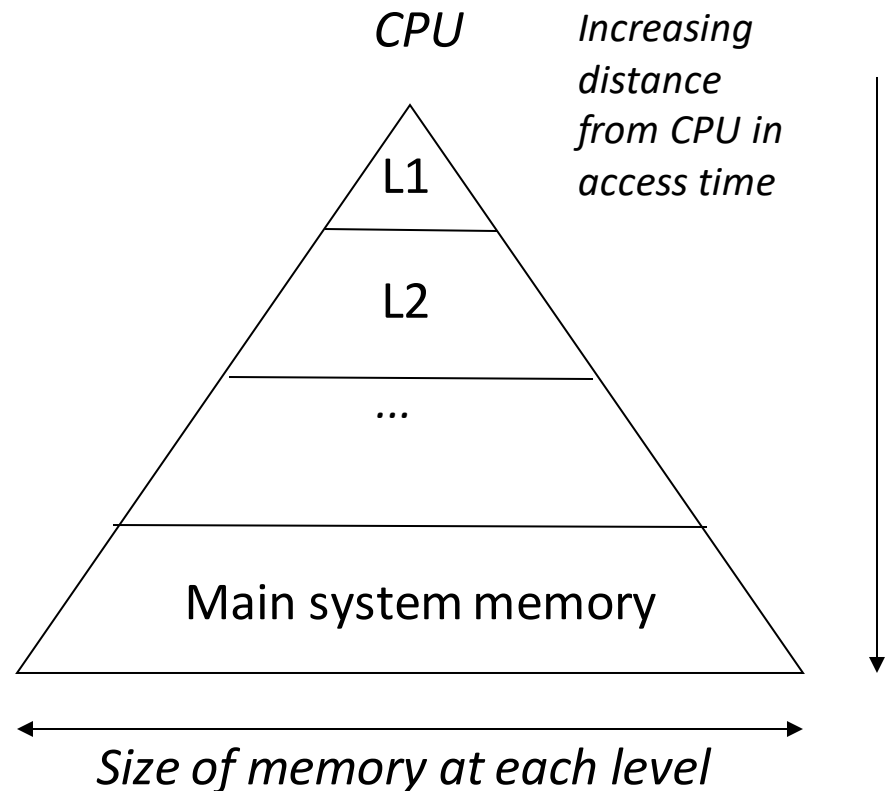
- Multi-level memory

- ❑ RAM Memory;
- ❑ Multi-level Cache.

Processor 1			Processor 2		
Core 1	Core 2	Core N	Core 1	Core 2	Core N
L1	L1	L1	L1	L1	L1
L2	L2	L2	L2	L2	L2
L3			L3		
RAM					

Parallel Architecture

- Temporal locality: if an item was referenced, it will be referenced again soon (e.g. cyclical execution in loops);
- Spatial locality: if an item was referenced, items close to it will be referenced too (the very nature of every program – serial stream of instructions)
- Stride: Step size between consecutive access of array elements;



Parallel Architecture

- AoS vs SoA (Array of Structures vs Structure of Arrays)

```
// Array of Structures (AoS)
struct coordinate {
    float x, y, z;
} crd[N];

...
for (int i = 0; i < N; i++)
    ... = ... f(crd[i].x, crd[i].y, crd[i].z);
```

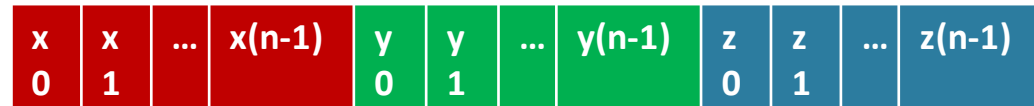
Non-Unit Stride



```
// Structure of Arrays (SoA)
struct coordinate {
    float x[N], y[N], z[N];
} crd;

...
for (int i = 0; i < N; i++)
    ... = ... f(crd.x[i], crd.y[i],
    crd.z[i]);
```

Unit Stride - Sequential Access



Parallel Architecture

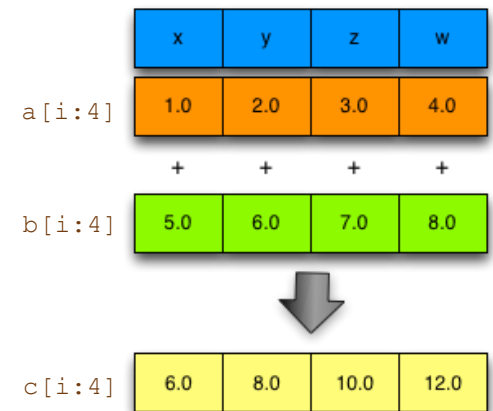
- Vectorization
 - Loading data into cache accordingly;
 - Store elements on SIMD registers or vectors;
 - Apply the same operation to a set of Data at the same time;
 - Iterations need to be independent;
 - Usually on inner loops.

Scalar loop

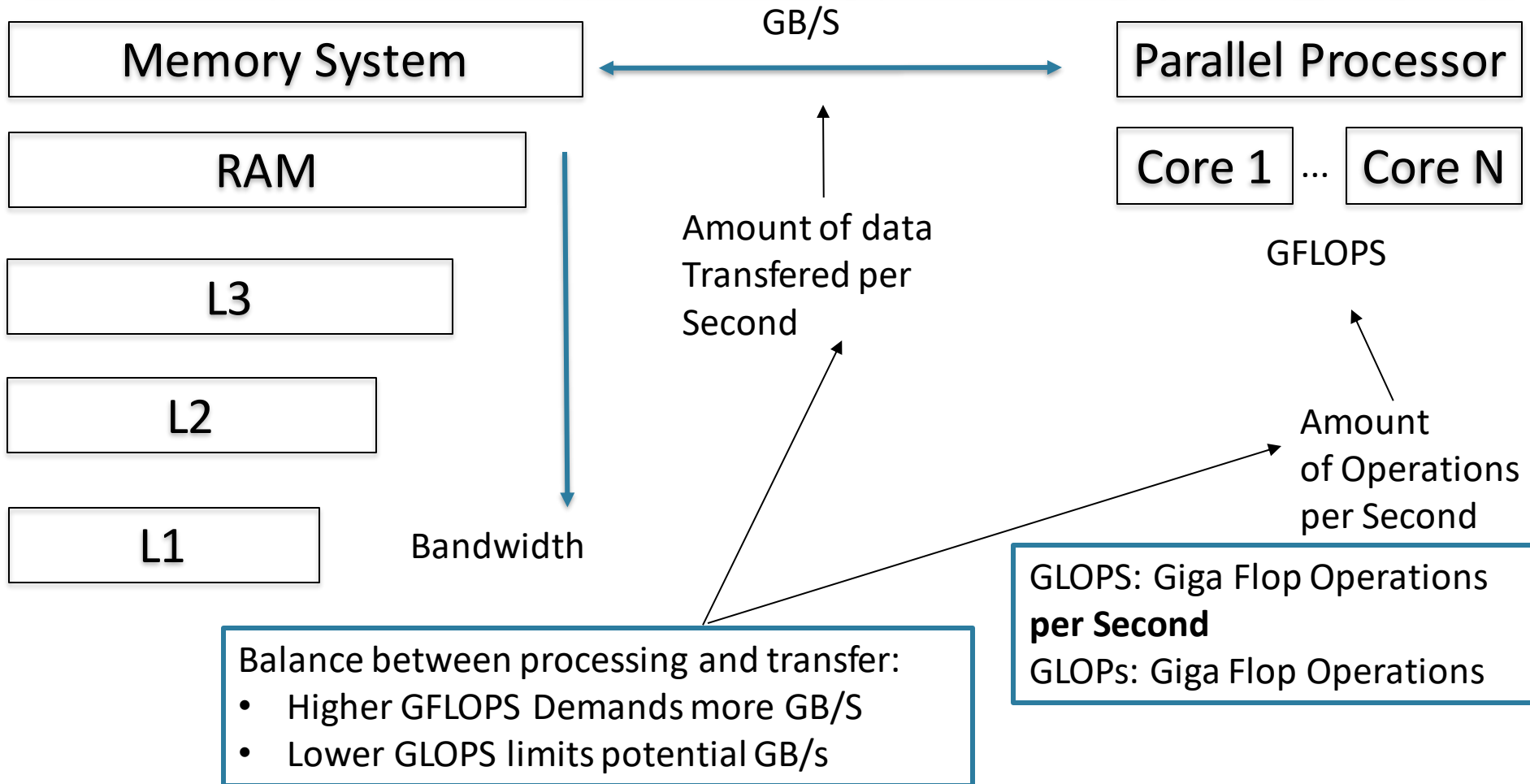
```
for (int i = 0; i < N; i++)  
    c[i] = a[i] + b[i];
```

SIMD loop (4 elements)

```
for (int i = 0; i < N; i += 4)  
    c[i:4] = a[i:4] + b[i:4];
```



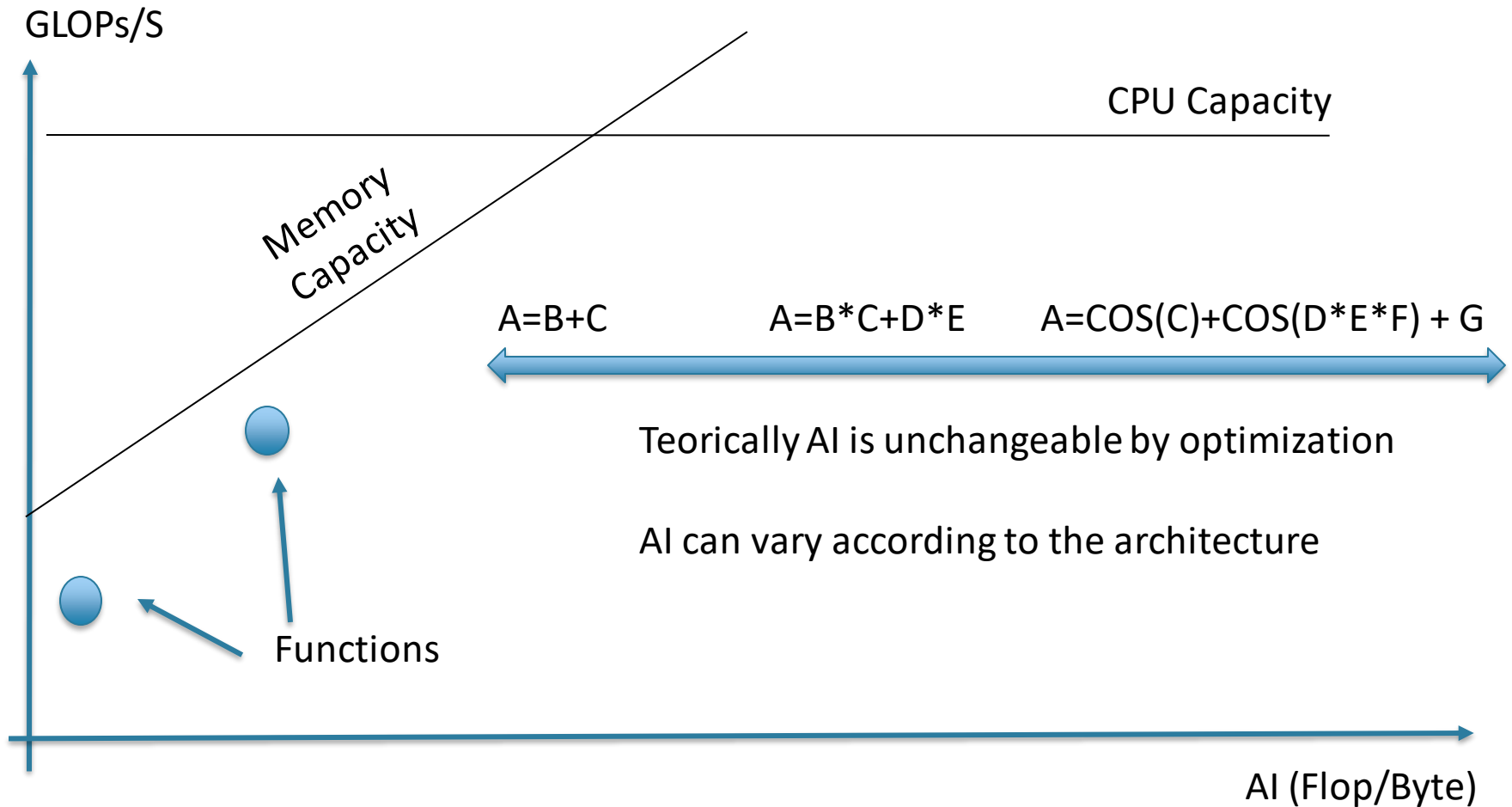
Performance Model



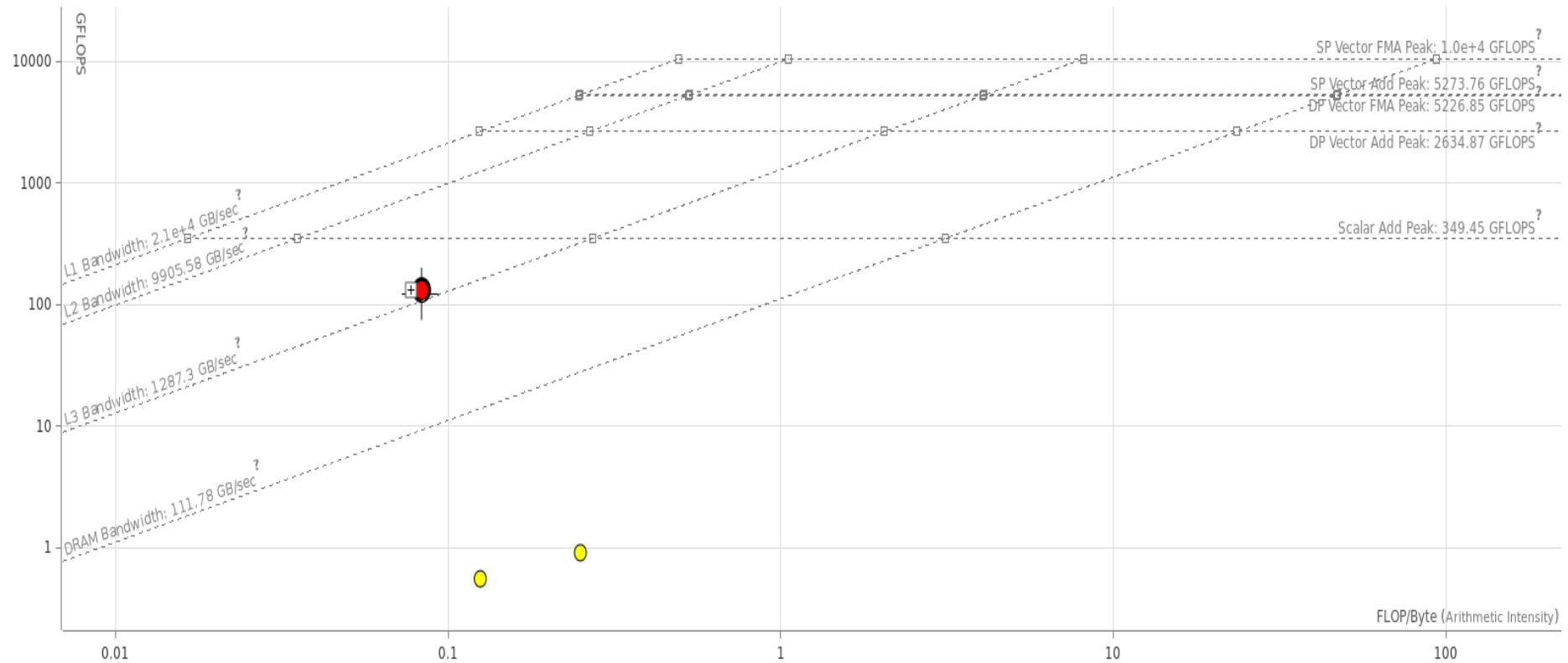
How to Measure such Balance?

Arithmetic Intensity (AI): Ratio between work performed and data transferred: Flop/Byte.

Roofline Model



Roofline Model

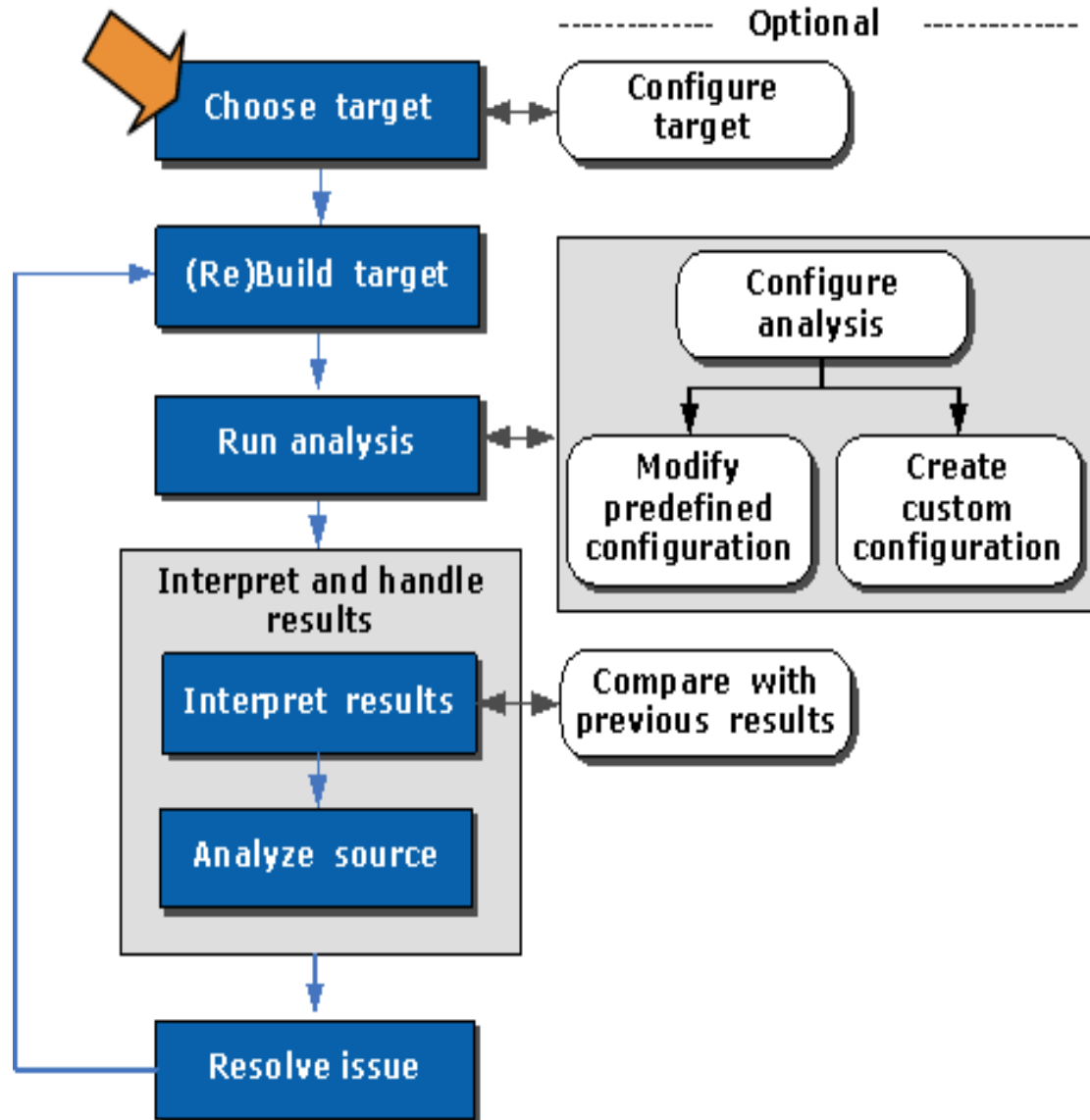


Optimization Workflow

Profiling Techniques include:

- Instrument Code
- Measure Hardware Events
- Performance Monitoring Units

Iterative process to improve Performance



Intel Advisor

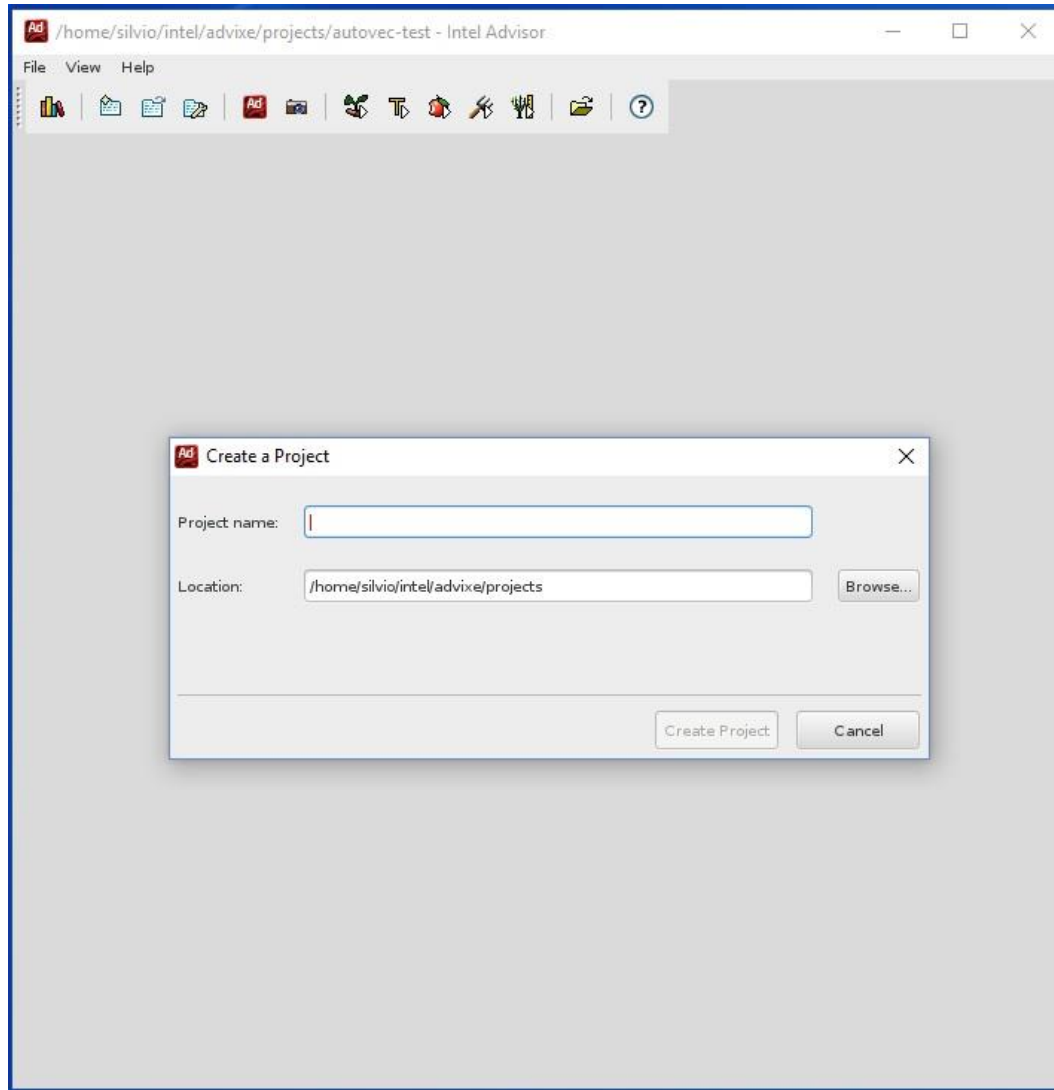
- Evaluate multi-threading parallelization
- Intel® Advisor XE
 - ❑ Performance modeling using several frameworks for multi-threading in processors and co-processors:
 - OpenMP, Intel® Cilk™ Plus, Intel® Threading Building Blocks
 - C, C++, Fortran (OpenMP only) and C# (Microsoft TPL)
 - ❑ Identify parallel opportunities
 - Detailed information about vectorization;
 - Check loop dependencies;
 - Memory Access Patterns
 - ❑ Scalability prediction: amount of threads/performance gains
 - ❑ Correctness (deadlocks, race conditions)



Intel Advisor

- Roofline Model
 - Obtain the Roofline of an application
- Survey Target;
 - Vectorization of loops: detailed information about vectorization;
 - Total Time: elapsed time on each loop considering the time involved in internal loops;
 - Self Time: elapsed time on each loop not considering the time involved in internal loops;
- Find Trip Counts;
 - Analysis to identify how many time particular loops run;
- Check Dependencies;
 - Analysis it there are loop-carried dependencies;
- Check Memory Access Patterns.
 - Analysis to identify how your code is iterating with memory.

Advisor – New Project



Advisor – Application and Parameters

teste - Project Properties (on skl01.ncc.unesp.br)

Analysis Target | Binary/Symbol Search | Source Search

Survey Analysis Types

- Survey Hotspots Analysis
- Trip Counts and Flow
- Suitability Analysis

Refinement Analysis Types

- Memory Access Pattern
- Dependencies Analysis

Launch Application

Specify and configure the application executable (target) to analyze. Press F1 for more details.

No application executable (target) file specified.

Application: Browse...

Application parameters: Modify...

☒ Use application directory as working directory

Working directory: Browse...

User-defined environment variables: Modify...

Modules:

☐ Include only the following module(s)

☒ Exclude the following module(s)

Modify...

Managed code profiling mode: Auto


OK Cancel

Advisor – Source

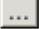
teste - Project Properties (on skl01.ncc.unesp.br)




Analysis Target | Binary/Symbol Search | **Source Search**

Additional Source File Locations
Specify local directories to include in the search. Press F1 for more details.

 Specify at least one directory in the Source Search tab before using the Suitability or Dependencies tools.


Search Directories

Add new search location 

☐ Search recursively

Exclude the following files:

Mask	File
Add new file mask. Example: b	Add new file to exclude 

OK Cancel

Advisor - Analysis

The screenshot shows the Vectorization Advisor interface. On the left, the 'Vectorization Workflow' is visible with steps: 'Run Roofline', '1. Survey Target', 'Mark Loops for Deeper Analysis' (with sub-steps 1.1, 2.1, and 2.2), and 'Batch mode' (OFF). On the right, the 'Threading Workflow' is visible. The top bar shows 'Elapsed time: 10,46s', 'Vectorized' and 'Not Vectorized' buttons, and filters for 'All Modules' and 'All Sources'. The 'Summary' tab is active, showing the 'Vectorization Advisor' section with a description. Arrows point from the workflow steps to descriptive text on the right.

Vectorization Workflow

OFF ☐ Batch mode

Run Roofline

1. Survey Target

Mark Loops for Deeper Analysis
Select checkboxes in the **Survey & Roofline** tab to mark loops for other Advisor analyses.
-- There are no marked loops --

1.1 Find Trip Counts and FLOP

-- Analyze all loops --

2.1 Check Memory Access Patterns

-- No loops selected --

2.2 Check Dependencies

-- No loops selected --

Threading Workflow

Elapsed time: 10,46s **Vectorized** **Not Vectorized** FILTER: All Modules All Sources

Summary Survey & Roofline Refinement Reports

Vectorization Advisor

Vectorization Advisor is a vectorization analysis toolset that lets you identify loops that will benefit most from vector parallelism, discover performance issues preventing from effective vectorization and characterize your memory vs. vectorization bottlenecks with Advisor Roofline model automation.

Run the Roofline Model

Collect Profiling Information

Obtain the amount of times a marked loop is executed

Obtain the stride distribution of marked loops

Verify if a marked loop has dependencies between iterations

Advisor – Survey Target

Summary

Survey Report

Refinement Reports

Annotation Report

Function Call Sites and Loops

Vector Issues

Self Time

Total Time

Type

Vectorized Loops

Instruction Set Analysis

					Vec...	Efficiency	Gai...	VL (Vector Length)	Compiler Estimated Gain	Traits
[loop in main at vec.c:50]	2 Inefficient memory access patterns present	0.730s	0.730s	Vectorized (Body)	AVX2		2.20x	16	2.20x	Extracts; Type Conversions
[loop in main at vec.c:63]	2 Data type conversions present	0.680s	0.680s	Vectorized Versions	AVX		5.71x	2; 8	<5.71x	Extracts; Inserts; Type Conversions
[loop in main at vec.c:55]	1 Data type conversions present	0.490s	0.490s	Vectorized (Body)	AVX2		5.79x	8	5.79x	Extracts; Inserts; Type Conversions

Source

Top Down

Loop Analytics

Loop Assembly

Recommendations

Compiler Diagnostic Details

File: vec.c:50 main

Line	Source	Total Time	%	Loop Time	%	Traits
36	randV=0;					
37						
38	a = (double*) mm_malloc(msize * msize * sizeof(double), 64);					
39	b = (double*) mm_malloc(msize * msize * sizeof(double), 64);					
40	c = (double*) mm_malloc(msize * msize * sizeof(double), 64);					
41						
42	srand(time(NULL));					
43	randV=rand();					
44	randV=randV*0.11;					
45	printf("randV %f\n", randV);					
46						
47						
48	for(j=0; j<10000; j++) {			1,899.980ms		
49						
50	for(i=0; i<40000; i++) {			729.970ms		
51	aosobj[i].x=i + randV;	229.961ms				Extracts; Type Conversions
52	aosobj[i].y=i*i+ randV;	290.000ms				Extracts; Type Conversions
53	aosobj[i].z=i*i+ randV;	210.009ms				Extracts; Type Conversions
54	}					
55	for(i=0; i<40000; i++) {	30.003ms		490.029ms		
56	soaobj.x[i]=i+i+ randV;	129.997ms				Extracts; Inserts; Type Conversions
57	soaobj.y[i]=i-i+ randV;	40.001ms				Type Conversions
58	soaobj.z[i]=i*i+ randV;	290.027ms				Extracts; Inserts; Type Conversions
59	}					
60	randV=rand();					Type Conversions
61	randV=randV*0.11;					
62						
63	for(i=0; i<40000; i++) {	9.999ms		679.985ms		
64	aosobj[i].x= aosobj[i].y+ aosobj[i].z + randV;	489.987ms				Extracts; Inserts; Type Conversions
65	}					
66						
67	for(i=0; i<40000; i++) {					
68						
69	soaobj.x[i]= soaobj.y[i]+ soaobj.z[i] + randV;	179.999ms				Extracts; Inserts; Type Conversions
70	}					
71	}					

Advisor – Survey Target

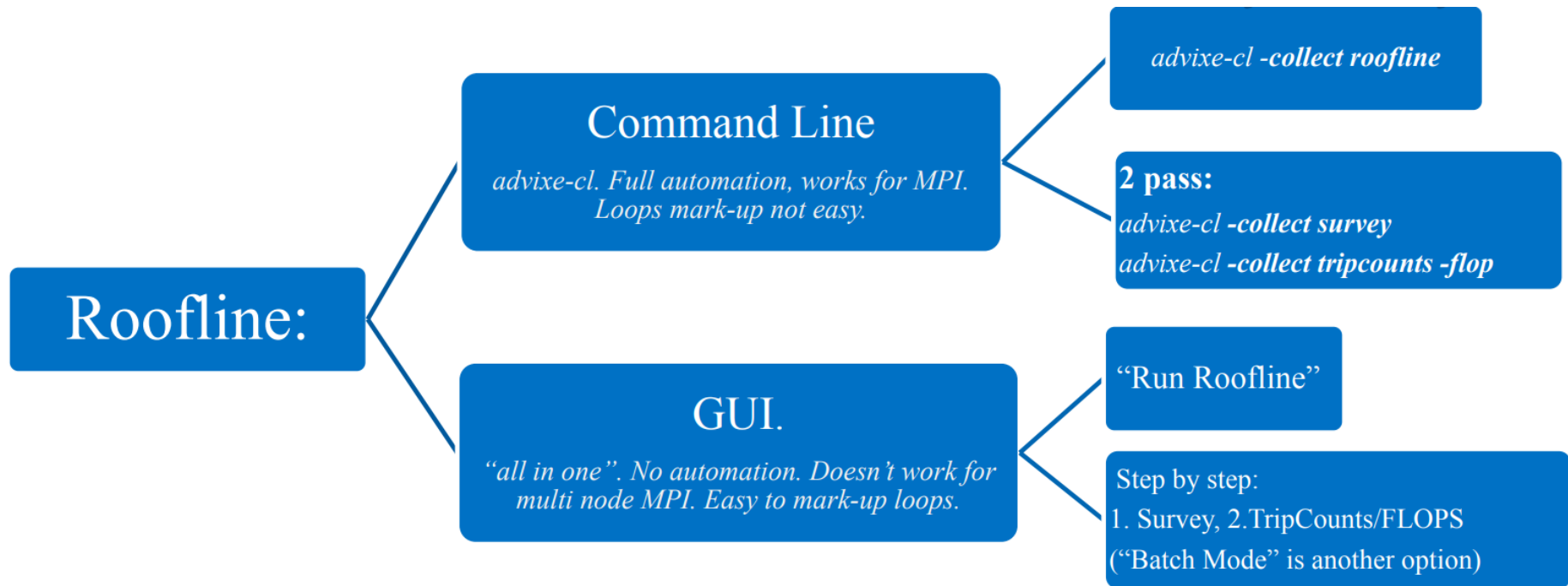
Instruction Set Analysis

Traits	Data Types	Number of Vector Registers	Vector Widths	Instruction Sets
Extracts; Type Conversions	Float32; Float64; Int32; UIn...	15	128/256	AVX; AVX2
Extracts; Inserts; Type Conversions	Float32; Float64	14; 15	128; 256	AVX
Extracts; Inserts; Type Conversions	Float32; Float64; Int32; UIn...	16	256	AVX; AVX2

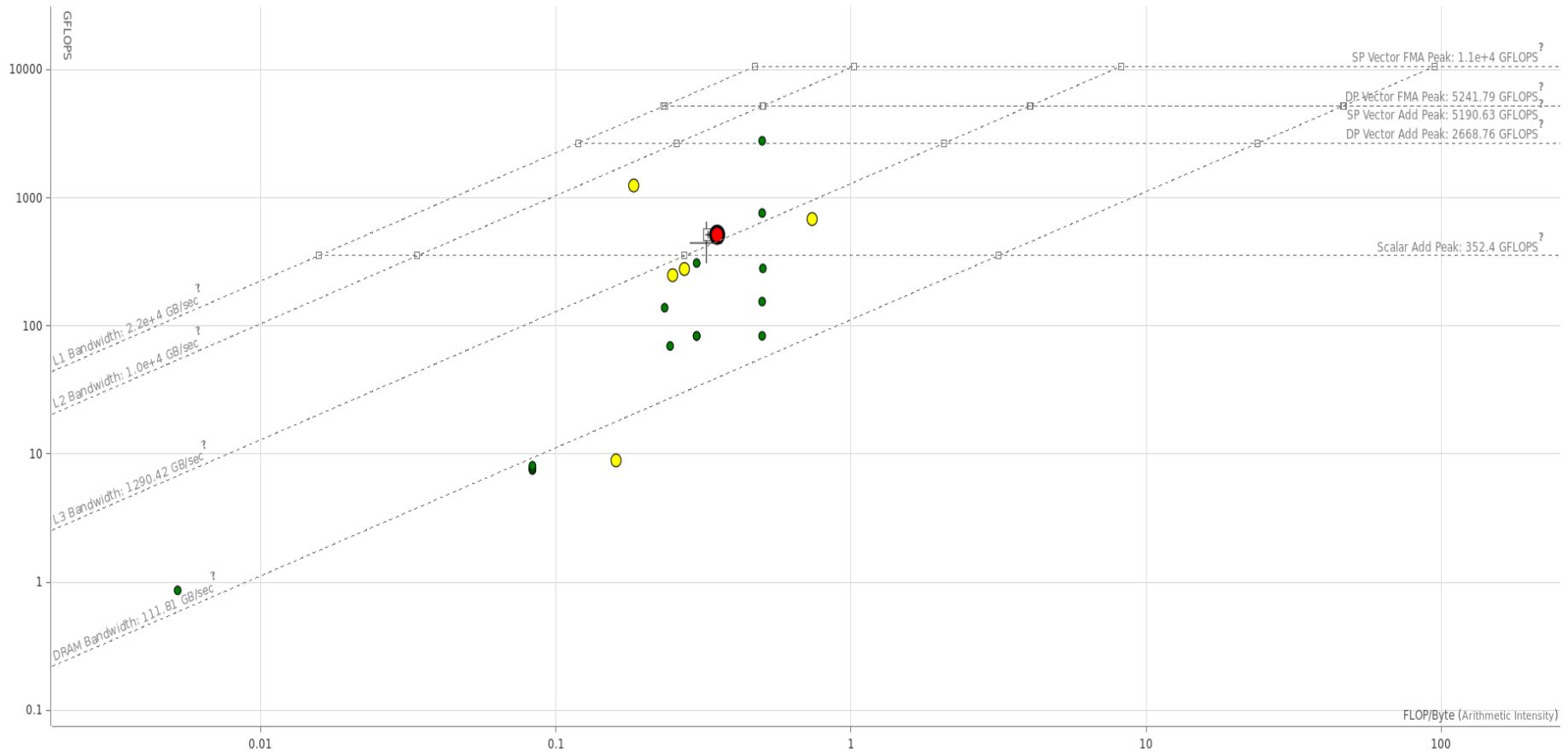
Advanced

Transformations	Unroll Factor	Vectorization Details	Optimization Details	Location
				vec.c:50
Fused; Unrolled	4		LOOP WAS DISTRIBUTED, CHUNK 1; LOOP WAS DISTRIBUTED, C...	vec.c:63
				vec.c:55

Roofline no Advisor



Roofline Example



Advisor – Memory Access Patterns

Vectorization Workflow | **Threading Workflow**

OFF | Batch mode

1. Survey Target

Collect

1.1 Find Trip Counts

Collect

Mark Loops for Deeper Analysis

Select loops in the Survey Report for Dependencies and/or Memory Access Patterns analysis.

4 loops are marked

2.1 Check Dependencies

Collect

2.2 Check Memory Access Patterns

Collect

Check memory access patterns in your application

Elapsed time: 1.90s | Vectorized | Not Vectorized | FILTER: All Modules

Summary | Survey Report | Refinement Reports | Annotation Report

Site Location	Loop-Carried Dependencies	Strides Distribution	Access Pattern	Site Name
[loop in main at vec.c:50]	No information available	0% / 100% / 0%	All const strides	loop_site_7
[loop in main at vec.c:55]	No information available	100% / 0% / 0%	All unit strides	loop_site_4
[loop in main at vec.c:63]	No information available	0% / 100% / 0%	All const strides	loop_site_8
[loop in main at vec.c:63]	No information available	0% / 100% / 0%	All const strides	loop_site_9

Memory Access Patterns Report | Dependencies Report | Recommendations

ID	Stride	Type	Source	Site Name	Nested Func
P1	48	Constant stride	vec.c:51	loop_site_7	
<pre>49 50 for(i=0; i<40000; i++) { 51 aosobj[i].x=i + randV; 52 aosobj[i].y=i*i+ randV; 53 aosobj[i].z=i+i+ randV; 54 }</pre>					
P2	48	Constant stride	vec.c:52	loop_site_7	
<pre>50 for(i=0; i<40000; i++) { 51 aosobj[i].x=i + randV; 52 aosobj[i].y=i*i+ randV; 53 aosobj[i].z=i+i+ randV; 54 }</pre>					
P3	48	Constant stride	vec.c:53	loop_site_7	
<pre>51 aosobj[i].x=i + randV; 52 aosobj[i].y=i*i+ randV; 53 aosobj[i].z=i+i+ randV; 54 } 55 for(i=0; i<40000; i++) {</pre>					
P6		Parallel site information	vec.c:50	loop_site_7	
<pre>48 for(j=0; j<10000; j++) { 49 50 for(i=0; i<40000; i++) { 51 aosobj[i].x=i + randV; 52 aosobj[i].y=i*i+ randV;</pre>					

Random-Stride

Constant-Stride

Unit-Stride

“Top-Down” Methodology Definition

- hierarchical organization of event-based metrics that identifies the dominant performance bottlenecks in an application;
- Such methodology aims at show, on average, how well the CPU's pipeline(s) were being utilized while running an application;
- Several hardware events are monitored according to the micro architecture;
- The results of this monitoring guides performance optimization .

Intel® VTune™ Amplifier XE

- **Compute-Intensive Application Analysis**
 - HPC Performance Characterization
- **Algorithm Analysis:**
 - Basic Hotspots
 - Advanced Hotspots
 - Concurrency
 - Locks & Waits
 - Memory Consumption

Views:

- Summary
- Bottom-up
- Top-down tree
- Tasks and frames
- Caller/Callee
- Source code
- Assembly
- Thread timeline
- Filters

View points:

- Hardware Event Counts
- Hardware Event Sample
- Hardware Issues
- Hotspots
- General Exploration
- Task Time
- Outros...

Intel® VTune™ Amplifier XE

- **Platform Analysis**

- CPU/GPU Concurrency
- System Overview
- GPU Hotspots
- GPU In-kernel Profiling
- Disk Input and Output
- Graphics Rendering (preview)
- CPU/FPGA Interaction (preview)

- **Microarchitecture Analysis**

- General Exploration
- Memory Access
- TSX Exploration
- TSX Hotspots
- SGX Hotspots

SGX Intel® Software Guard Extensions

TSX Intel® Transactional Synchronization Extensions

Intel® VTune™ Amplifier XE

Projects

Analysis

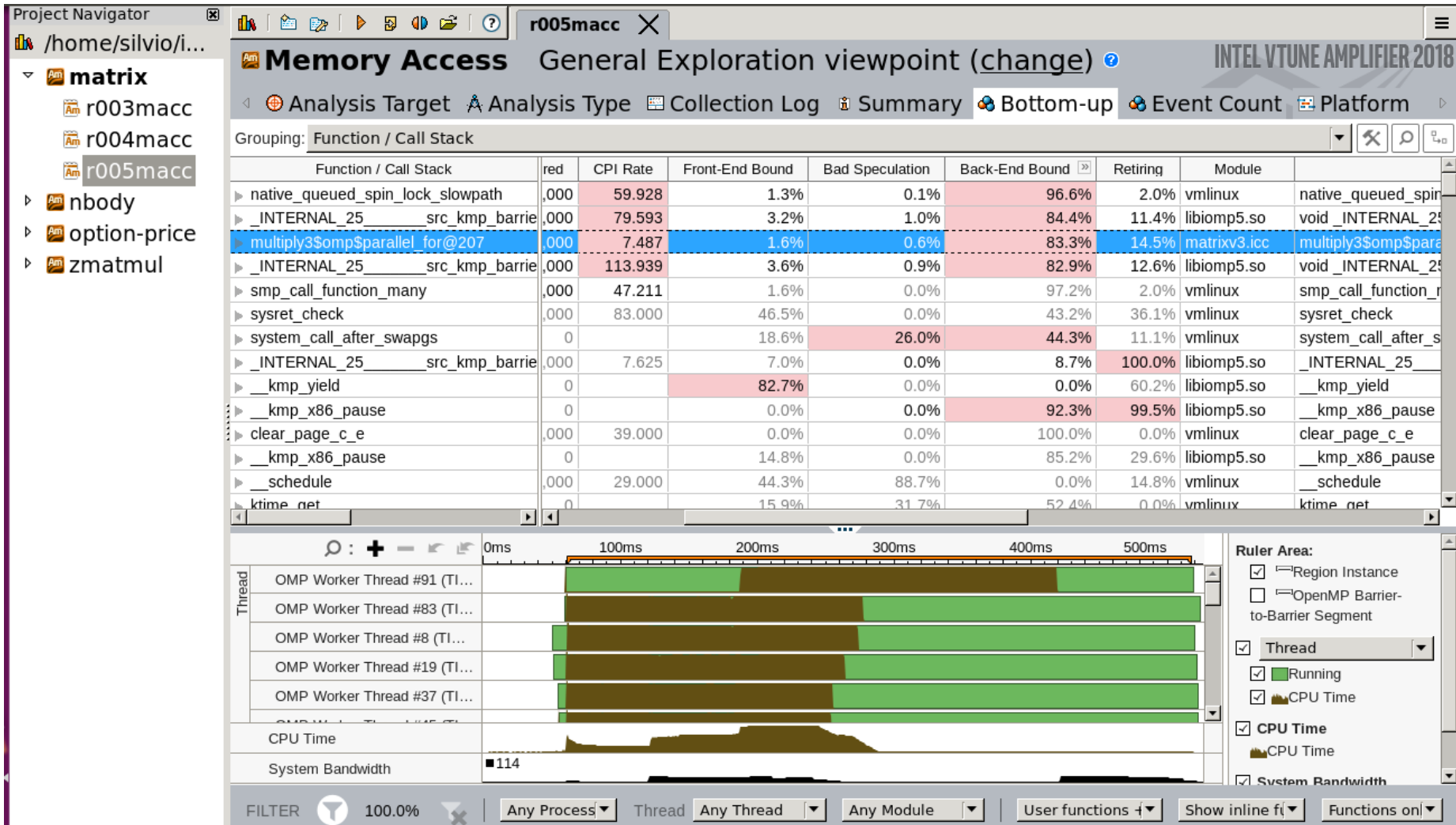
Results

The screenshot displays the Intel VTune Amplifier XE interface. On the left, the 'Project Navigator' shows a project named 'matrix' with sub-projects 'r003macc' and 'r004macc'. The main window shows the 'Memory Access' analysis for the 'r003macc' target. The analysis is in the 'General Exploration viewpoint'. The results section shows an elapsed time of 2.205s and various performance metrics.

Metric	Value	Unit
Elapsed Time	2.205s	
Clockticks	492,714,600,000	
Instructions Retired	5,294,100,000	
CPI Rate	93.069	
Front-End Bound	2.5%	of Pipeline Slots
Bad Speculation	0.0%	of Pipeline Slots
Back-End Bound	94.9%	of Pipeline Slots
Memory Bound	92.9%	of Pipeline Slots
L1 Bound	3.0%	of Clockticks
L2 Bound	N/A*	of Clockticks
L3 Bound	32.6%	of Clockticks
DRAM Bound	N/A*	of Clockticks
Store Bound	0.0%	of Clockticks
Core Bound	2.0%	of Pipeline Slots
Retiring	2.6%	of Pipeline Slots
Total Thread Count	98	
Paused Time	0s	

*N/A is applied to metrics with undefined value. There is no data to calculate the metric.

Intel® VTune™ Amplifier XE



General Exploration Analysis – Intel Xeon

- Results:
 - Intel Xeon: high level of LLC hit
- Suggestion: improve data locality in order to make better use of cache
 - loop interchange;

Loop Interchange

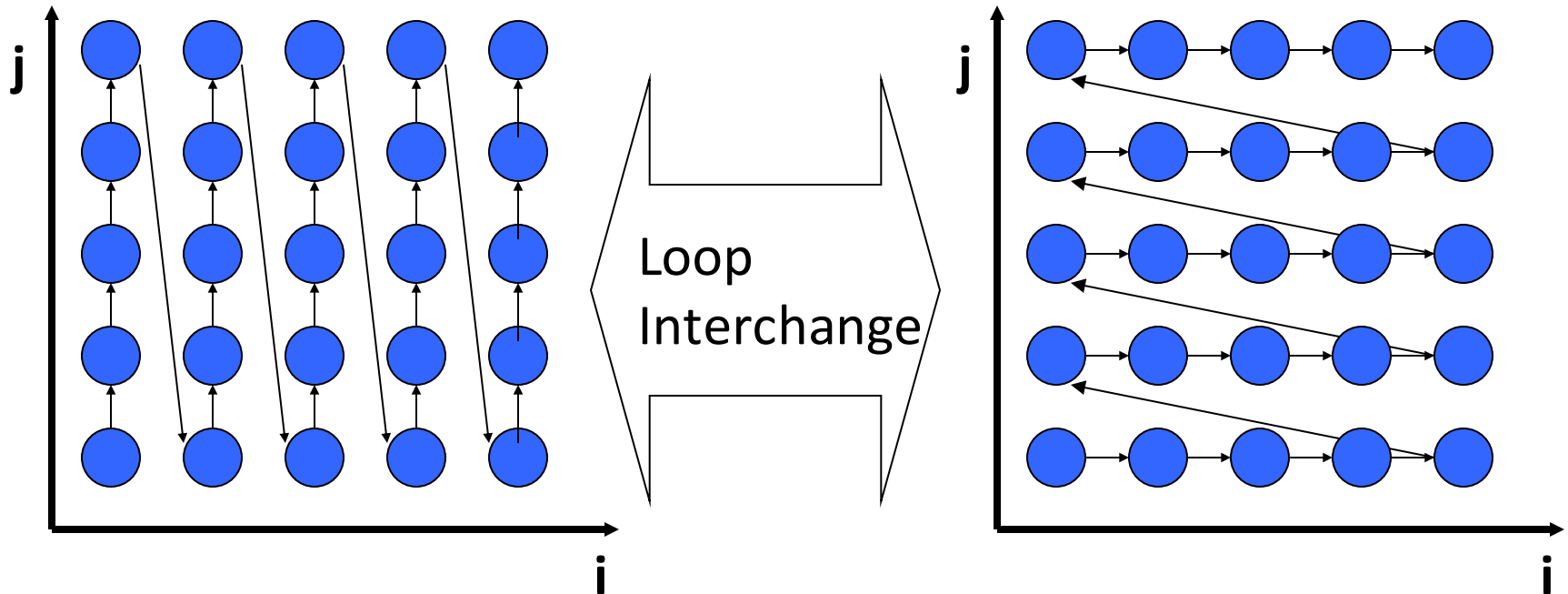
- Row-major order storage.

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----

Loop Interchange

Reading in Row –Major order improves locality




```
for(i=0; i<W; i++)  
  for(j=0; j<H; j++)  
    A[i][j] = ...;
```

```
for(j=0; j<H; j++)  
  for(i=0; i<W; i++)  
    A[i][j] = ...;
```

Matrix Multiplication Loop Interchange

```
void multiply1(int msize, int tid, int numt, TYPE a[][NUM], TYPE
b[][NUM], TYPE c[][NUM], TYPE t[][NUM])
{
    int i,j,k;
    #pragma omp parallel for
    for(i=0; i<msize; i++) {
        for(k=0; k<msize; k++) {
            for(j=0; j<msize; j++) {
                c[i][j] = c[i][j] + a[i][k] * b[k][j];
            }
        }
    }
}
```



Matrix Multiplication Loop Interchange

Interchange

Clockticks:	86,272,200,000	
Instructions Retired:	3,641,400,000	
CPI Rate ^⑦ :	23.692	🚩
Front-End Bound ^⑦ :	2.8%	of Pipeline Slots
Bad Speculation ^⑦ :	0.0%	of Pipeline Slots
⌵ Back-End Bound ^⑦ :	88.2%	🚩 of Pipeline Slots
⌵ Memory Bound ^⑦ :	66.3%	🚩 of Pipeline Slots
⌵ L1 Bound ^⑦ :	3.5%	of Clockticks
FB Full ^⑦ :	100.0%	of Clockticks
L2 Bound ^⑦ :	0.1%	of Clockticks
⌵ L3 Bound ^⑦ :	4.5%	of Clockticks
Contested Accesses ^⑦ :	0.9%	of Clockticks
Data Sharing ^⑦ :	0.0%	of Clockticks
L3 Latency ^⑦ :	6.1%	of Clockticks
⌵ DRAM Bound ^⑦ :	2.1%	of Clockticks
Memory Bandwidth ^⑦ :	6.5%	of Clockticks
⌵ Memory Latency ^⑦ :	1.7%	of Clockticks
Local DRAM ^⑦ :	1.5%	of Clockticks
Remote DRAM ^⑦ :	0.0%	of Clockticks
Remote Cache ^⑦ :	0.2%	of Clockticks
Store Bound ^⑦ :	0.0%	of Clockticks
⌵ Core Bound ^⑦ :	21.8%	🚩 of Pipeline Slots
Port Utilization ^⑦ :	3.3%	of Clockticks
Retiring ^⑦ :	9.0%	of Pipeline Slots
Total Thread Count:	98	
Paused Time ^⑦ :	0s	

Without Interchange

Clockticks:	492,714,600,000	
Instructions Retired:	5,294,100,000	
CPI Rate ^⑦ :	93.069	🚩
Front-End Bound ^⑦ :	2.5%	of Pipeline Slots
Bad Speculation ^⑦ :	0.0%	of Pipeline Slots
⌵ Back-End Bound ^⑦ :	94.9%	🚩 of Pipeline Slots
⌵ Memory Bound ^⑦ :	92.9%	🚩 of Pipeline Slots
⌵ L1 Bound ^⑦ :	3.0%	of Clockticks
FB Full ^⑦ :	N/A*	of Clockticks
L2 Bound ^⑦ :	N/A*	of Clockticks
⌵ L3 Bound ^⑦ :	32.6%	🚩 of Clockticks
Contested Accesses ^⑦ :	0.7%	of Clockticks
Data Sharing ^⑦ :	0.0%	of Clockticks
L3 Latency ^⑦ :	3.5%	of Clockticks
⌵ DRAM Bound ^⑦ :	N/A*	of Clockticks
Memory Bandwidth ^⑦ :	38.6%	of Clockticks
⌵ Memory Latency ^⑦ :	22.7%	of Clockticks
Local DRAM ^⑦ :	4.2%	of Clockticks
Remote DRAM ^⑦ :	0.2%	of Clockticks
Remote Cache ^⑦ :	0.6%	of Clockticks
Store Bound ^⑦ :	0.0%	of Clockticks
⌵ Core Bound ^⑦ :	2.0%	of Pipeline Slots
Port Utilization ^⑦ :	1.4%	of Clockticks
Retiring ^⑦ :	2.6%	of Pipeline Slots
Total Thread Count:	98	
Paused Time ^⑦ :	0s	

HPC Performance Characterization

SP GFLOPS[?]: 28.970

Effective Physical Core Utilization[?]: 35.0% (16.788 out of 48) 🚩 📄

Effective Logical Core Utilization[?]: 32.1% (30.819 out of 96) 🚩

Serial Time (outside parallel regions)[?]: 0.062s (29.9%) 🚩

Parallel Region Time[?]: 0.145s (70.1%) 📄

Estimated Ideal Time[?]: 0.085s (41.1%)

OpenMP Potential Gain[?]: 0.060s (29.0%) 🚩

Top OpenMP Regions by Potential Gain

This section lists OpenMP regions with the highest potential for performance improvement. The Potential Gain metric shows the elapsed time that could be saved if the region was optimized to have no load imbalance assuming no runtime overhead.

OpenMP Region	OpenMP Potential Gain [?]	(%) [?]	OpenMP Region Time [?]
multiply3\$omp\$parallel:96@unknown:207:215	0.060s 🚩	29.0% 🚩	0.145s

Effective CPU Utilization Histogram

Memory Bound[?]: 60.2% 🚩 of Pipeline Slots

Cache Bound[?]: 10.5% of Clockticks

DRAM Bound[?]: 9.6% of Clockticks

NUMA: % of Remote Accesses[?]: 0.0%

Bandwidth Utilization Histogram

HPC Performance Characterization

⌵ FPU Utilization [?]: 0.3% 🚩

SP FLOPs per Cycle [?]: 0.221 Out of 64 🚩

Vector Capacity Usage [?]: 100.0%

⌵ FP Instruction Mix:

⌵ % of Packed FP Instr. [?]: 100.0%

% of 128-bit [?]: 0.0%

% of 256-bit [?]: 0.0%

% of 512-bit [?]: 100.0%

% of Scalar FP Instr. [?]: 0.0%

FP Arith/Mem Rd Instr. Ratio [?]: 0.132 🚩

FP Arith/Mem Wr Instr. Ratio [?]: 0.363 🚩

⌵ Top Loops/Functions with FPU Usage by CPU Time

This section provides information for the most time consuming loops/functions with floating point operations.

Function	CPU Time [?]	FPU Utilization [?]	Vector Instruction Set [?]	Loop Type [?]
[Loop at line 211 in multiply3\$omp\$parallel_for@207]	5.701s	0.3%	AVX512F_512(512)	Body
__kmp_fork_barrier	4.700s	0.4%		
osq_lock	0.112s	1.7%		
[Loop at line 208 in multiply3\$omp\$parallel_for@207]	0.073s	3.9%	AVX(256); AVX2(256); AVX512F_256(256); AVX512F_512(512)	Body