

# Introduction to Performance Optimization and Tuning Tools

Steve Lantz, Cornell University

*CoDaS-HEP Summer School, July 19, 2023*



# Goals

- Give an overview of what is meant by performance optimization and tuning
- Provide basic guidance on how to understand the performance of a code using tools
- Provide a starting point for performance optimizations



# Performance Tuning: What Is It? Why Do It?

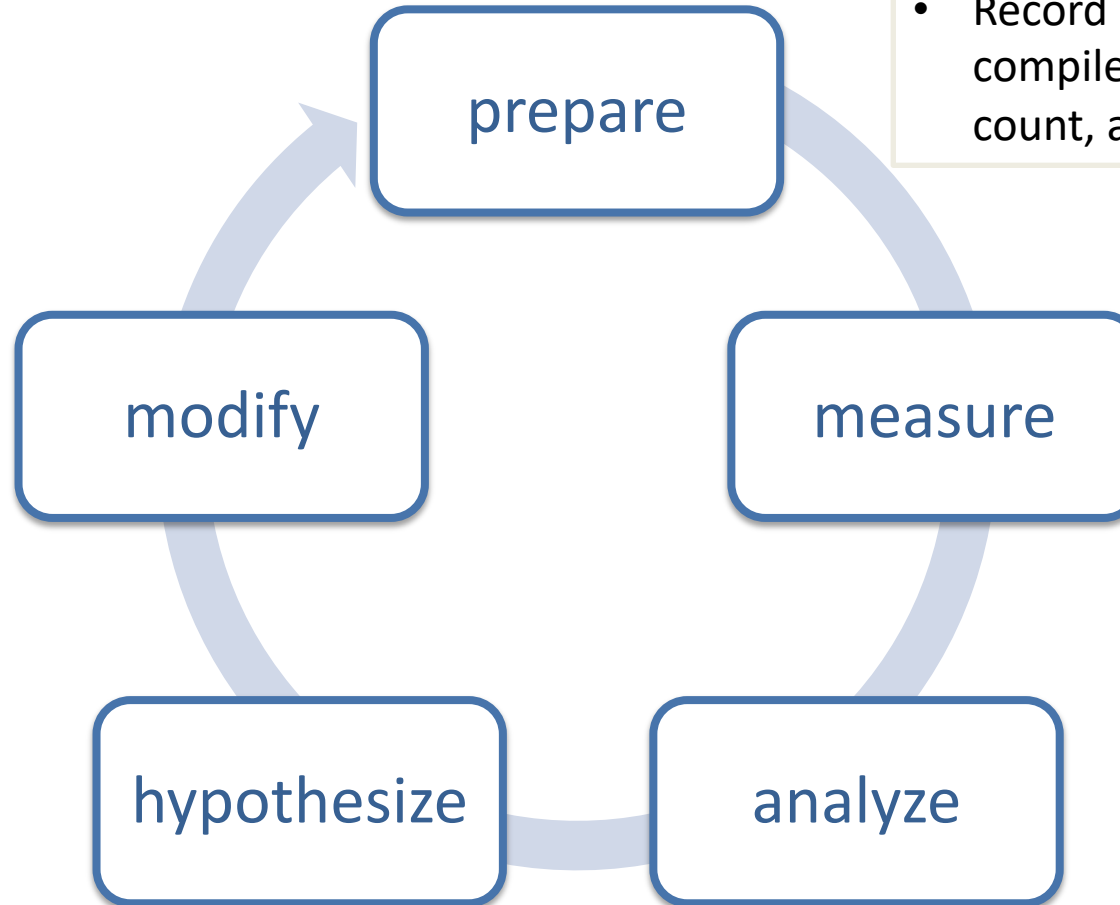
- What is performance tuning?
  - The process of improving the efficiency of an application to make better use of a given hardware resource
  - A cycle of identifying bottlenecks, eliminating these where possible, and rechecking efficiency – usually continued until performance objectives are satisfied
  - Writing code informed by one's understanding of the performance features of the given hardware (see previous presentations on “What Every Computational Physicist Should Know About Computer Architecture” and “Vector Parallelism on Multi-Core Processors”)
- Why does performance matter?
  - Energy efficiency is becoming increasingly important
  - Today's applications only use a fraction of the machine
  - Due to complex architectures, mapping applications onto architectures is hard



# The Performance Tuning Cycle

- Change only **one thing at a time**
- Consider the ease (difficulty) of implementation
- Keep **track** of all **changes**
- Apply regression test to **ensure correctness** after each change
- Remember: fast computing of a wrong result is completely irrelevant

- Choose a workload which is measurable, representative, static, reproducible, and quantifiable
- Record code generation, compiler version, compiler flags, input parameters, core count, affinity, etc.



# What Do I Measure?

- Choose metrics which quantify the performance of your code
  - **Time** spent at different levels: whole program, functions, lines of code
  - **Hardware counters** can help you figure out the reasons for slow spots
- What are some easy ways to make time measurements?
  - Wrap your executable command in the Linux “time” command
    - Get an idea of overall run time: `time ./my_exe` (or `/bin/time ./my_exe`)
    - No way to zero in on performance bottlenecks
  - Insert calls to timers around critical loops/functions
    - `gettimeofday()`, `MPI_Wtime()`, `omp_get_wtime()`
    - Available in common libraries (system, MPI, OpenMP respectively)
    - Good for checking known hotspots in a small code base
    - Hard to maintain, require significant a priori knowledge of the code



# Advantages of Performance Tools

- Performance tools (recommended)
  - Collect a lot data with varying granularity, cost and accuracy
  - Connect back to the source code (use -g compiler flag)
  - Analyze/visualize collected data using the tool
  - The learning curve is steep, but you can climb it gradually
- Tools generally work in one of two ways

## Sampling

- Records system state at periodic intervals
- Useful to get an overview
- Low and uniform overhead
- Ex. Profiling

## Instrumentation

- Records all events
- Provide detailed per event information
- High overhead for request events
- Ex. Tracing



# Performance Tools Overview

- Basic OS tools
  - /bin/time
  - **perf**, gprof, igprof (from HEP)
  - valgrind, callgrind
- Hardware counters
  - PAPI API & tool set
- Community open source
  - HPCToolkit (Rice Univ.)
  - TAU (Univ. of Oregon)
  - Open|SpeedShop (Krell)
- Commercial products
  - Linaro Forge (DDT, MAP)
- Vendor supplied (free)
  - **Intel Advisor, Intel VTune**
  - Intel Trace Analyzer and Collector (MPI)
  - AMD  $\mu$ Prof
  - CrayPat
  - NVIDIA Nsight Compute (CUDA)
  - NVIDIA pgprof (OpenACC)
  - AMD Omniprof (ROC)

No tool can do everything. Choose the right tool for the right task.



# What Can I Learn From Performance Tools?

- Where am I spending my time?
  - Find the hotspots
- Is my code memory bound or compute bound?
  - Memory bound code has lots of events like these (tracked by hardware counters):
    - L1/L2/L3 cache misses
    - TLB misses
  - Compute bound code has lots of events like these:
    - Pipeline stalls not due to memory events
    - Type conversions
    - Time spent in unvectorized loops
- Is my I/O inefficient?





# Typical Performance Pitfalls on a Single Node

- Scattered memory accesses that constantly bring in new cache lines
  - Storing data as an array of structs (AoS) instead of a struct of arrays (SoA)
  - Looping through arrays with a large stride

More cache lines  $\Rightarrow$   
data must be fetched  
from more distant  
caches, or from RAM

	Registers	L1	L2	LLC	DRAM
Speed (cycle)	1	~4	~10	~30	~200
Size	< KB	~32KB	~256KB	~35MB	10-100GB

- Mismatched types in assignments

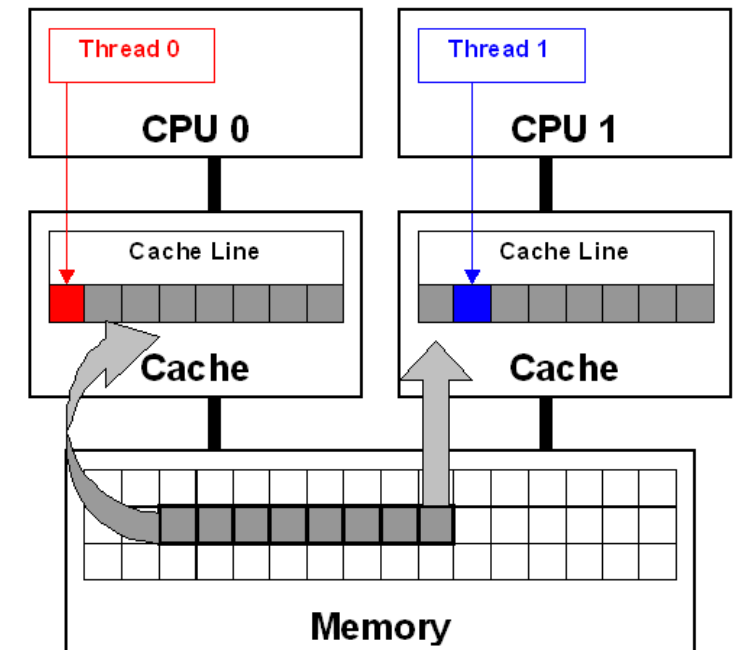
```
float x=3.14; //bad: 3.14 is a double
float s=sin(x); //bad: sin() is a double
precision function
long v=round(x); //bad: round() takes and
returns double
```

```
float x=3.14f; //good: 3.14f is a float
float s=sinf(x); //good: sin() is a single
precision function
long v=lroundf(x); //good: lroundf() takes
float and returns long
```



# Typical Performance Pitfalls: Multithreading

- Load imbalance
- False sharing: when CPUs alter different variables in the same cache line ↓
  - Data aren't really shared, but caches must stay coherent
  - Data always travel together in “cache lines” of 64 bytes
- Insufficient parallelism
- Synchronization
  - Use private thread storage to avoid synchronization
- Non-optimal memory placement
  - Memory is actually allocated on first touch
  - Thread that touches first has fastest access



<https://software.intel.com/en-us/articles/avoiding-and-identifying-false-sharing-among-threads>



# Linux Tool: *perf*

- Perf is a performance analyzing tool in Linux
  - *perf record*: measure and save sampling data for a single program
    - *-g*: enable call-graph (callers/callee information)
  - *perf report*: analyze the file generated by perf record, can be flat profile or graph
    - *-g*: enable call-graph (callers/callee information)
  - *perf stat*: measure total event count for a single program
    - *-e event-name-1,event-name-2*: choose from event names provided by *perf list*
  - *perf list*: list available hardware and software events for measurement
- When compiling the code, use the following flags for easier interpretation
  - *-g*: generate debug symbols needed to annotate source
  - *-fno-omit-frame-pointer*: provide stack chain/backtrace

<https://perf.wiki.kernel.org/index.php/Tutorial>  
<https://www.brendangregg.com/perf.html>



# Example: Finding Hot Spots with *perf*

- Compile the code: *g++ -g -fno-omit-frame-pointer -O3 -DNAIVE matmul\_2D.cpp -o mm\_naive.out*
- Collect profiling data: *perf record -g ./mm\_naive.out 500*
- Open the result: *perf report -g*

Samples: 7K of event 'cycles:uppp', Event count (approx.): 5629336320				
Children	Self	Command	Shared Object	Symbol
+ 99.95%	0.00%	mm_naive.out	libc-2.17.so	[.] __libc_start_main
+ 99.95%	0.00%	mm_naive.out	mm_naive.out	[.] main
- 99.69%	99.69%	mm_naive.out	mm_naive.out	[.] compute_naive
__libc_start_main				
main				
compute_naive				
0.09%	0.09%	mm_naive.out	mm_naive.out	[.] init_matrix_2D
0.06%	0.06%	mm_naive.out	libc-2.17.so	[.] __random
0.06%	0.06%	mm_naive.out	libc-2.17.so	[.] __memset_sse2
0.03%	0.03%	mm_naive.out	[unknown]	[.] 0xffffffff8196c4e7
0.03%	0.00%	mm_naive.out	[unknown]	[.] 0000000000000000
0.02%	0.02%	mm_naive.out	libc-2.17.so	[.] __random_r
0.01%	0.01%	mm_naive.out	mm_naive.out	[.] rand@plt
0.01%	0.01%	mm_naive.out	ld-2.17.so	[.] do_lookup_x
0.01%	0.01%	mm_naive.out	libc-2.17.so	[.] _int_malloc
0.01%	0.01%	mm_naive.out	libc-2.17.so	[.] intel_check_word
0.00%	0.00%	mm_naive.out	ld-2.17.so	[.] check_match.9523
0.00%	0.00%	mm_naive.out	[unknown]	[.] 0x00000000000c2698
0.00%	0.00%	mm_naive.out	ld-2.17.so	[.] _dl_sysdep_start
0.00%	0.00%	mm_naive.out	ld-2.17.so	[.] dl_main
0.00%	0.00%	mm_naive.out	ld-2.17.so	[.] _dl_load_cache_lookup
0.00%	0.00%	mm_naive.out	ld-2.17.so	[.] _etext
0.00%	0.00%	mm_naive.out	ld-2.17.so	[.] _dl_map_object
0.00%	0.00%	mm_naive.out	ld-2.17.so	[.] __libc_memalign@plt
0.00%	0.00%	mm_naive.out	ld-2.17.so	[.] _dl_start_user

Press "A" →

```
init_matrix_2D /home/beiwang/codas_perftools/mm_naive.out
Percent
__attribute__((noinline)) void init_matrix_2D(double **A, double **B, double **C, int matrix_size){
#pragma omp parallel for
for (int i=0; i<matrix_size; i++) {
test %ecx,%ecx
↓ jle 401558 <init_matrix_2D(double**, double**, b8
__attribute__((noinline)) void init_matrix_2D(double **A, double **B, double **C, int matrix_size){
push %rbp
lea -0x1(%rcx),%eax
mov %rsp,%rbp
push %r15
push %r14
mov %rsi,%r14
push %r13
lea 0x8(%rdi,%rax,8),%rsi
push %r12
push %rbx
lea 0x8(%rax,8),%r13
mov %rdi,%r12
mov %rdx,%r15
sub $0x18,%rsp
mov %rsi,-0x38(%rbp)
nop
xor %ebx,%ebx
nop
for (int j = 0 ; j < matrix_size; j++) {
A[i][j]=((double) rand() / (RAND_MAX));
48: → callq rand@plt
pxor %xmm0,%xmm0
mov (%r12),%rdx
55.56 cvtsi2sd %eax,%xmm0
divsd 0x5e7(%rip),%xmm0 # 401ae8 <__dso_handle+0x60>
22.22 movsd %xmm0, (%rdx,%rbx,1)
B[i][j]=((double) rand() / (RAND_MAX));
→ callq rand@plt
pxor %xmm0,%xmm0
mov (%r14),%rdx
cvtsi2sd %eax,%xmm0
C[i][j]=0.0;
mov (%r15),%rax
B[i][j]=((double) rand() / (RAND_MAX));
divsd 0x5c7(%rip),%xmm0 # 401ae8 <__dso_handle+0x60>
11.11 movsd %xmm0, (%rdx,%rbx,1)
C[i][j]=0.0;
11.11 movq $0x0, (%rax,%rbx,1)
add $0x8,%rbx
for (int j = 0 ; j < matrix_size; j++) {
cmp %rbx,%r13
4014e8 <init_matrix_2D(double**, double**, 48
↑ jne $0x8,%r12
add $0x8,%r14
add $0x8,%r15
for (int i=0; i<matrix_size; i++) {
cmp -0x38(%rbp),%r12
↑ jne 4014e0 <init_matrix_2D(double**, double**, 40
}
}
```



# Example: Counting Cache Misses with *perf stat*

List of pre-defined events (to be used in `-e`):

branch-instructions OR branches	[Hardware event]
branch-misses	[Hardware event]
bus-cycles	[Hardware event]
cache-misses	[Hardware event]
cache-references	[Hardware event]
cpu-cycles OR cycles	[Hardware event]
instructions	[Hardware event]
ref-cycles	[Hardware event]
alignment-faults	[Software event]
bpf-output	[Software event]
context-switches OR cs	[Software event]
cpu-clock	[Software event]
cpu-migrations OR migrations	[Software event]
dummy	[Software event]
emulation-faults	[Software event]
major-faults	[Software event]
minor-faults	[Software event]
page-faults OR faults	[Software event]
task-clock	[Software event]
L1-dcache-load-misses	[Hardware cache event]
L1-dcache-loads	[Hardware cache event]
L1-dcache-stores	[Hardware cache event]
L1-icache-load-misses	[Hardware cache event]
LLC-load-misses	[Hardware cache event]
LLC-loads	[Hardware cache event]
LLC-store-misses	[Hardware cache event]
LLC-stores	[Hardware cache event]
branch-load-misses	[Hardware cache event]
branch-loads	[Hardware cache event]
dTLB-load-misses	[Hardware cache event]
dTLB-loads	[Hardware cache event]
dTLB-store-misses	[Hardware cache event]
dTLB-stores	[Hardware cache event]
iTLB-load-misses	[Hardware cache event]
iTLB-loads	[Hardware cache event]
node-load-misses	[Hardware cache event]
node-loads	[Hardware cache event]
node-store-misses	[Hardware cache event]
node-stores	[Hardware cache event]

- The *perf list* command lists all available CPU counters
  - Check *man perf\_event\_open* to see what each event measures
- The *perf stat* command instruments and summarizes selected CPU counters

*perf stat -e cpu-cycles,instructions,L1-dcache-loads,L1-dcache-load-misses ./mm\_naive.out 500*

Performance counter stats for './mm\_naive.out 500':

5,564,503,540	cpu-cycles		
10,063,662,841	instructions	#	1.81 insn per cycle
3,767,490,743	L1-dcache-loads		
1,475,374,174	L1-dcache-load-misses	#	39.16% of all L1-dcache hits

1.691104619 seconds time elapsed

- Make changes, see if L1 load misses improve, e.g.



# Intel Advisor

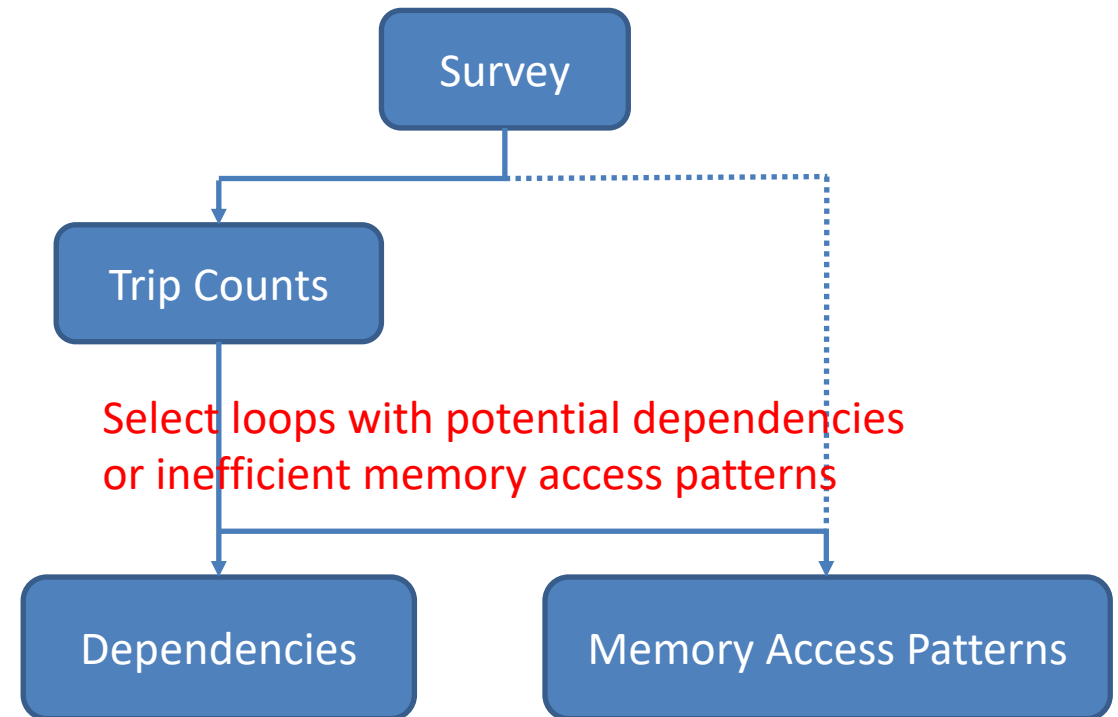
Two very useful analyses in Intel Advisor will be highlighted:

- Vectorization advisor
  - Provide vectorization information from vectorization report
  - Identify the hotspots where your efforts pay off the most
  - Provide call graph information
  - Identify the performance and vectorization issues
  - Check memory access pattern, dependencies, more
- Roofline
  - How much performance is being left on the table
  - Where are the bottlenecks
  - Which can be improved
  - Which are worth improving



# Workflow of Vectorization Advisor

- **Survey:** find the vectorization information for loops and provide suggestions for improvement
- **Trip Counts:** generate a **Roofline** Chart
- **Memory Access Patterns (MAP):** see how you access the data
- **Dependencies:** determine if it is safe to force vectorization





# Advisor Advises You About Performance Issues

The screenshot displays the Intel Advisor 2019 interface. At the top, there are filters for 'Elapsed time: 4.12s', 'Vectorized', 'Not Vectorized', and a 'FILTER' section with dropdowns for 'All Modules', 'All Sources', 'Loops And Functions', and 'All Threads'. Below this is a 'Summary' tab with sub-tabs for 'Survey & Roofline' and 'Refinement Reports'. The main table shows a list of function call sites and loops, with columns for 'CPU Time' (Self Time and Total Time), 'Type', 'Why No Vectorization?', and 'Vectorized Loops' (Vecto..., Efficiency, Gain, VL (V...)). A red circle highlights the 'Performance Issues' icon in the first row. Below the table, there are tabs for 'Source', 'Top Down', 'Code Analytics', 'Assembly', 'Recommendations', and 'Why No Vectorization?'. The 'Recommendations' tab is active, showing a list of issues. A red circle highlights the first issue: 'Possible inefficient memory access patterns present'. Below this, there is a section for 'Confirm inefficient memory access patterns' with a link to 'Memory Access Patterns analysis'. Another red circle highlights the second issue: 'Data type conversions present'. Below this, there is a section for 'Use the smallest data type' with a link to 'Use the smallest data type'. On the right side, there are two additional sections: 'Possible inefficient memory access patterns present' and 'Data type conversions present'. At the bottom, there is a footer with the text: 'Intel, and the Intel logo are trademarks of Intel Corporation in the U.S. and/or other countries. \*Other names and brands may be claimed as the property of others. © Intel Corporation'.

Elapsed time: 4.12s Vectorized Not Vectorized FILTER: All Modules All Sources Loops And Functions All Threads Customize View OFF

Summary Survey & Roofline Refinement Reports

INTEL ADVISOR 2019

Function Call Sites and Loops	Performance Issues	CPU Time		Type	Why No Vectorization?	Vectorized Loops				Instruction Traits
		Self Time	Total Time			Vecto...	Efficiency	Gain ...	VL (V...	
[loop in MoveParticles at nbody.cpp:76]	2 Possible ine...	4.090s	4.090s	Vectorized (Body)		AVX512	~53%	8.47x	16	Appr. Reci
MoveParticles		0.010s	4.100s	Inlined Function						2-Source P
_start		0.000s	4.100s	Function						
[loop in main at nbody.cpp:55]	1 Data type con...	0.000s	4.100s	Scalar	inner loop was already ...					Appr. Reci
main		0.000s	4.100s	Function						2-Source P
[loop in main at nbody.cpp:204]	1 Data type con...	0.000s	4.100s	Scalar	compile time constraint...					Divisions; F

Source Top Down Code Analytics Assembly Recommendations Why No Vectorization?

All Advisor-detectable issues: C++ / Fortran

**! Possible inefficient memory access patterns present**

Inefficient memory access patterns may result in significant vector code execution slowdown or block automatic vectorization by the compiler. Improve performance by investigating.

**Confirm inefficient memory access patterns**

There is no confirmation inefficient memory access patterns are present. To fix: Run a [Memory Access Patterns analysis](#).

**! Data type conversions present**

There are multiple data types within loops. Utilize hardware vectorization support more effectively by avoiding data type conversion.

**Use the smallest data type**

The source loop contains data types of different widths. To fix: Use the smallest data type that gives the needed precision to use the entire vector register width.

**Possible inefficient memory access patterns present**

Confirm inefficient memory access patterns

**Data type conversions present**

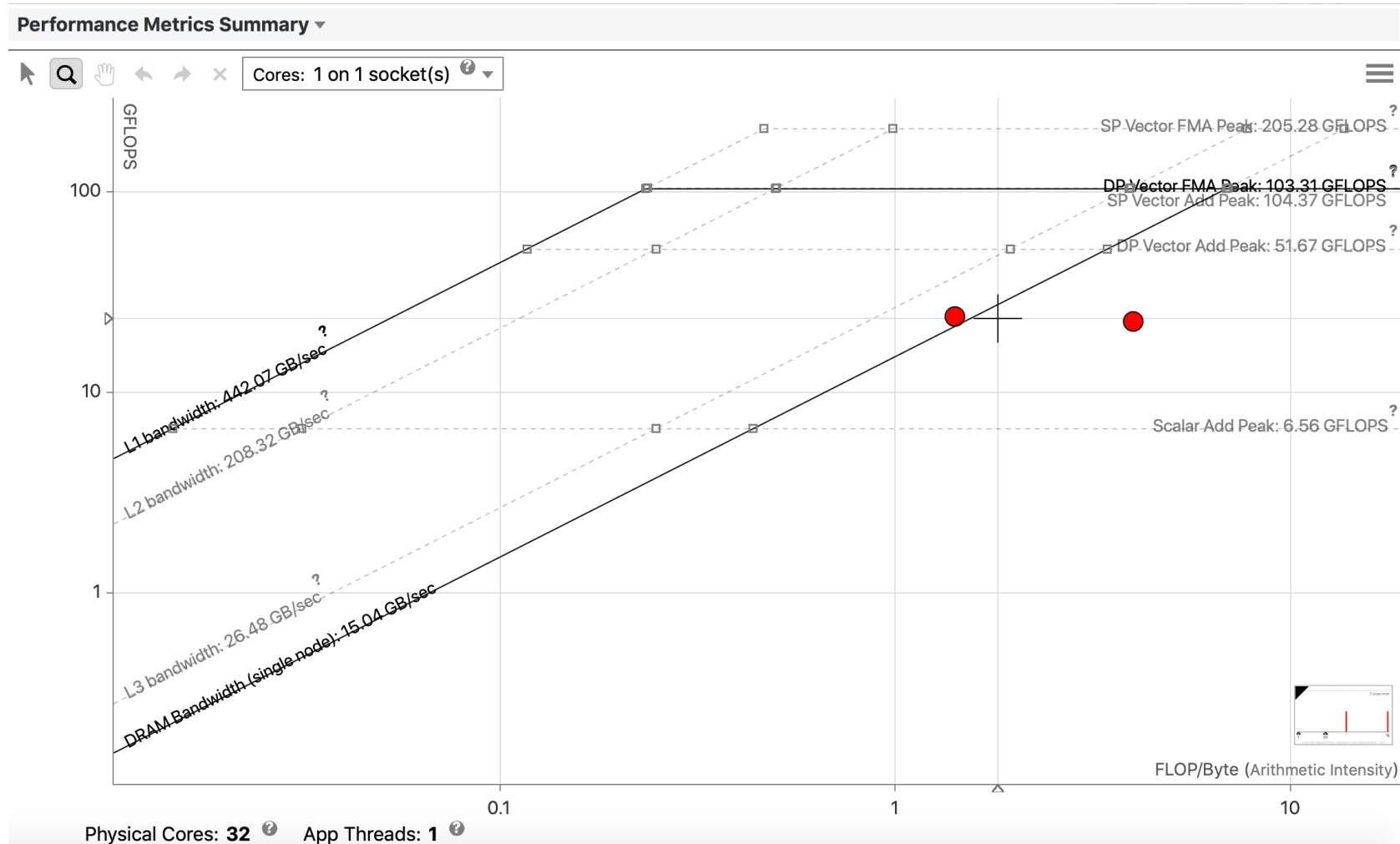
Use the smallest data type

Intel, and the Intel logo are trademarks of Intel Corporation in the U.S. and/or other countries.  
\*Other names and brands may be claimed as the property of others.  
© Intel Corporation





# Roofline Analysis: What Is It?



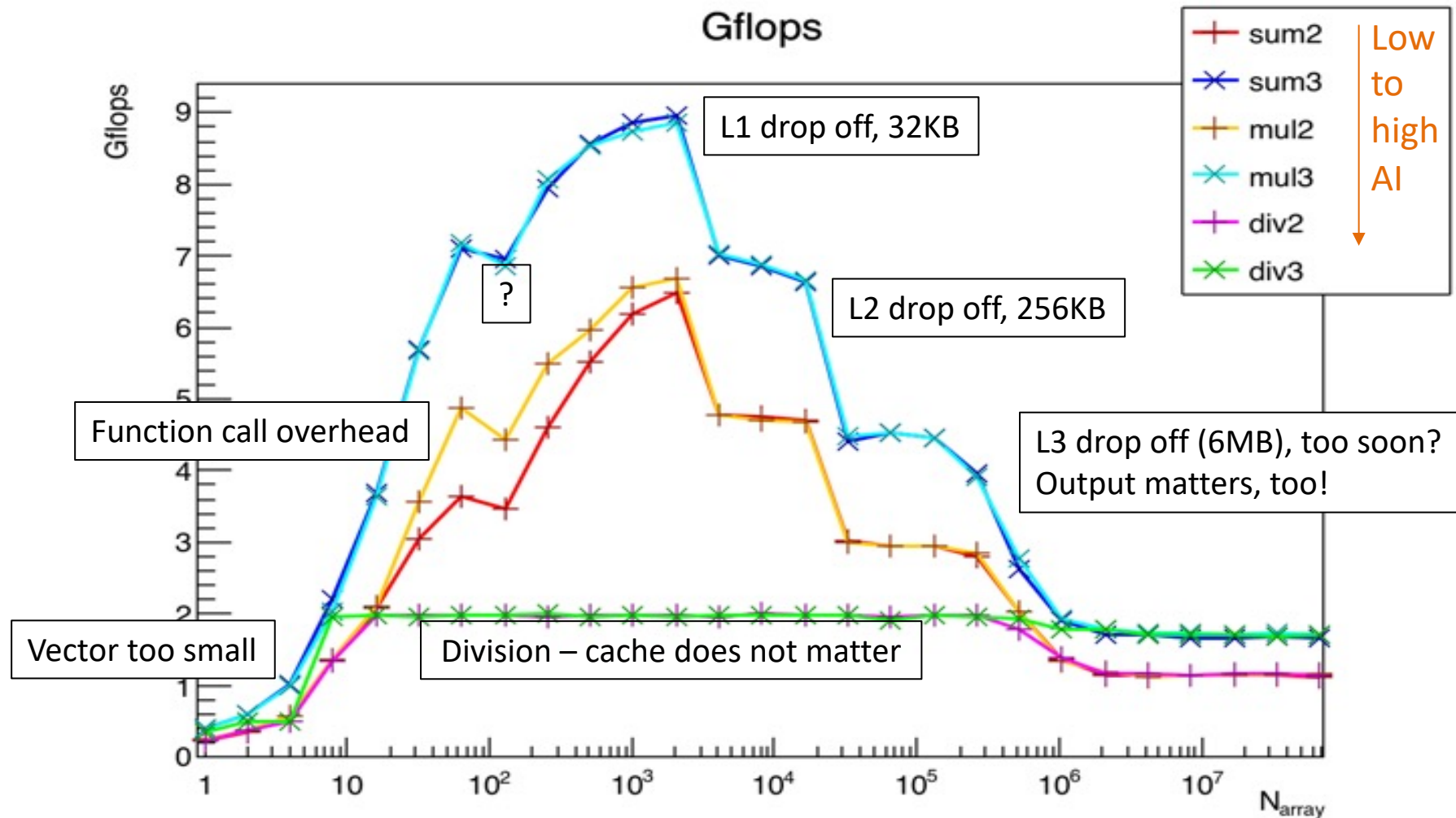
# Towards Peak Flop/s: Arithmetic Intensity

- Arithmetic intensity or AI is the number of flops executed by a code divided by the bytes of memory that are required to perform the computations
  - AI is an intrinsic property of the code
- Even a simple stride-1 loop may not get the peak flop/s rate, if its AI is low
  - VPU becomes stalled waiting for loads and stores to complete
  - Delays become longer as the memory request goes further out in the hierarchy from L1 to L2 (to L3?) to RAM
  - Even if the right vectors are in L1 cache, there is limited bandwidth from L1 to registers!
- If the goal is to maximize flop/s, you'll want to try to improve AI
- Also want threads to work on independent, cache-size chunks of data
  - Watch out for false sharing, where 2 threads fight needlessly over a cache line

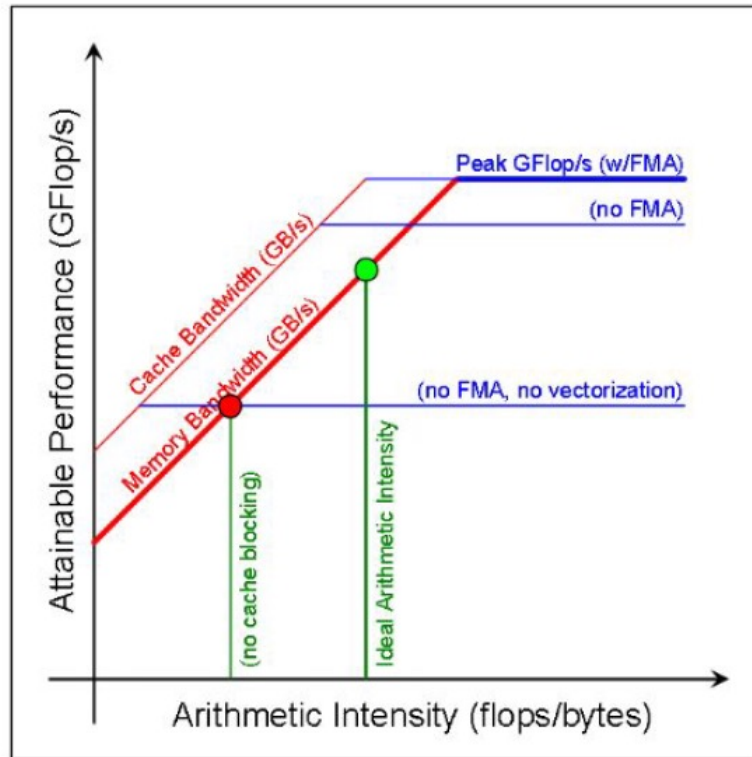


# Effect of AI and Caches on GFLOP/s

Data taken on a laptop (2.6 GHz, vector width 8):

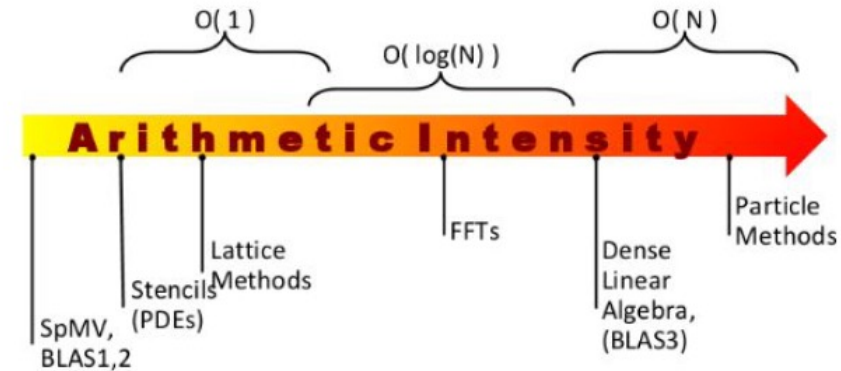


# Roofline Analysis Explained



$$\text{Attainable FLOPs/sec} = \min \left\{ \begin{array}{l} \text{Peak FLOPs/sec,} \\ \text{Peak Memory Bandwidth} \times \text{Arithmetic Intensity} \end{array} \right.$$

$$\text{Arithmetic Intensity} = \frac{\text{Total FLOPs}}{\text{Total Bytes}}$$



Deslippe et al., "Guiding Optimization Using the Roofline Model," tutorial presentation at IXPUG2016, Argonne, IL, Sept. 21, 2016.

<https://anl.app.box.com/v/IXPUG2016-presentation-29>



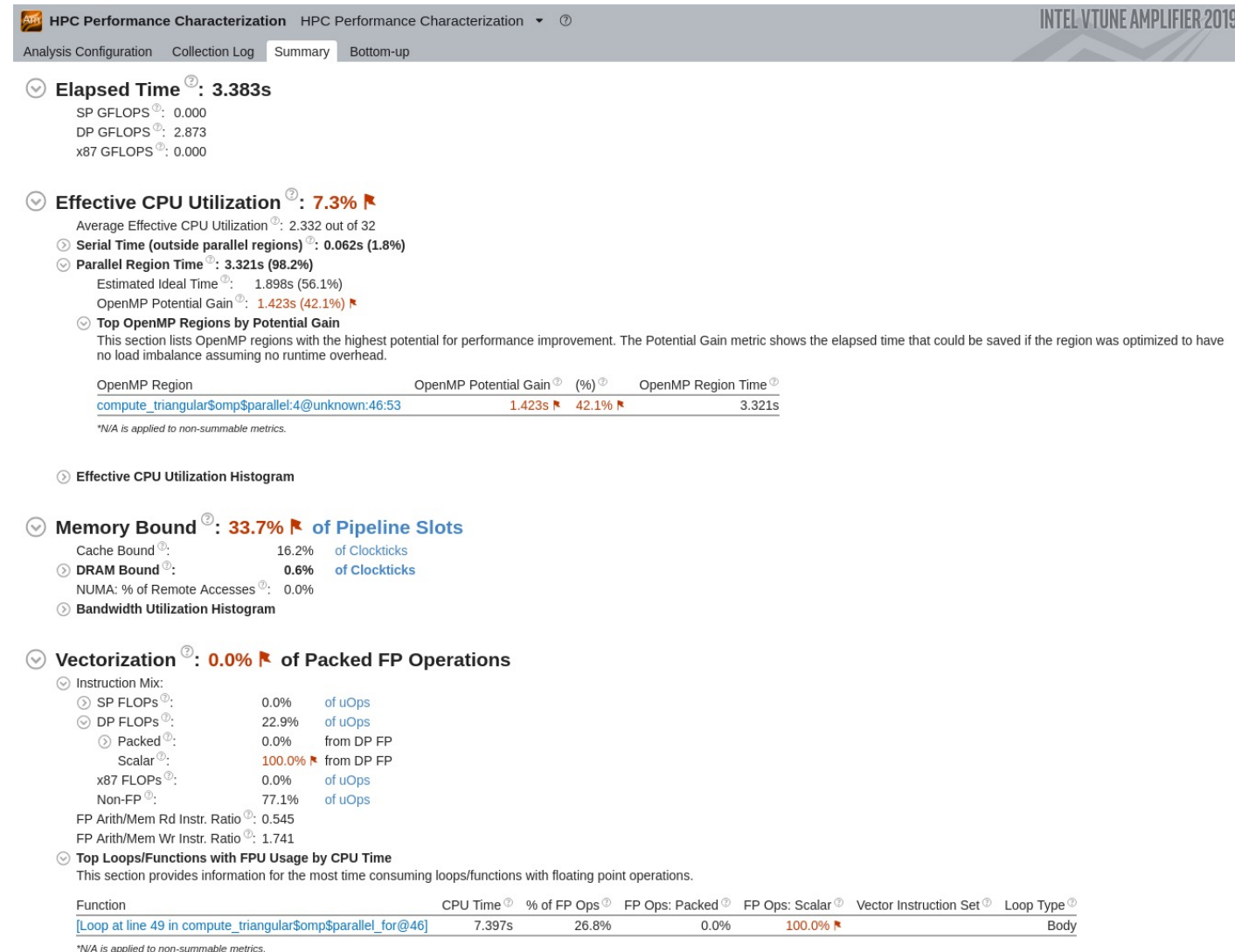
# What Does Roofline Analysis Tell You?

- Roofline analysis is a way of telling whether a piece of code is *compute bound* or *memory bound*
  - The “roofline” is a performance ceiling related to *hardware characteristics*
- The *arithmetic intensity* or AI (flop/byte) of a code tells you what part of the roof the code is under
  - AI is a *software characteristic* telling you the extent to which the code is limited by its need to load and store data from/to memory
- The roofline sets the highest flop/s rate possible for a given piece of code
  - If some of your functions fall way below that rate, you may need to investigate why
  - It’s possible to show that the AI needed for reaching *theoretical peak* flop/s (the highest flat roof) implies that 50% of operands are vector constants, i.e., they are loaded just once and never leave registers!

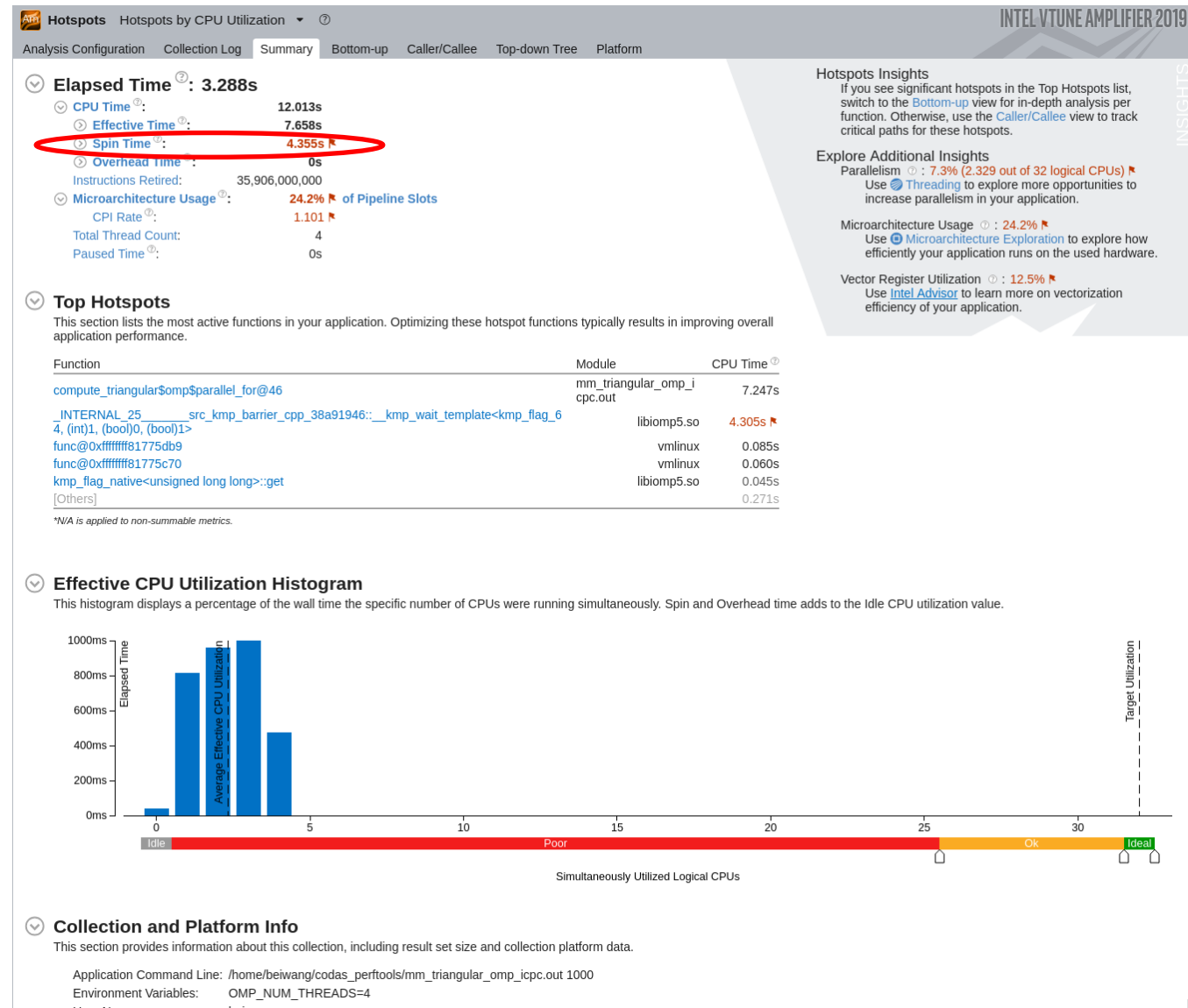


# Intel VTune

- Covers all aspects of execution
  - Hotspots
  - Processor microarchitecture
  - Memory accesses
  - Threading
  - I/O
- Flexible
  - GUI in Linux, Windows and macOS
  - Drills down to source code, assembly
  - Easy setup, no special compiling
- Shared memory only
  - Serial or OpenMP
  - MPI, but only within a single node

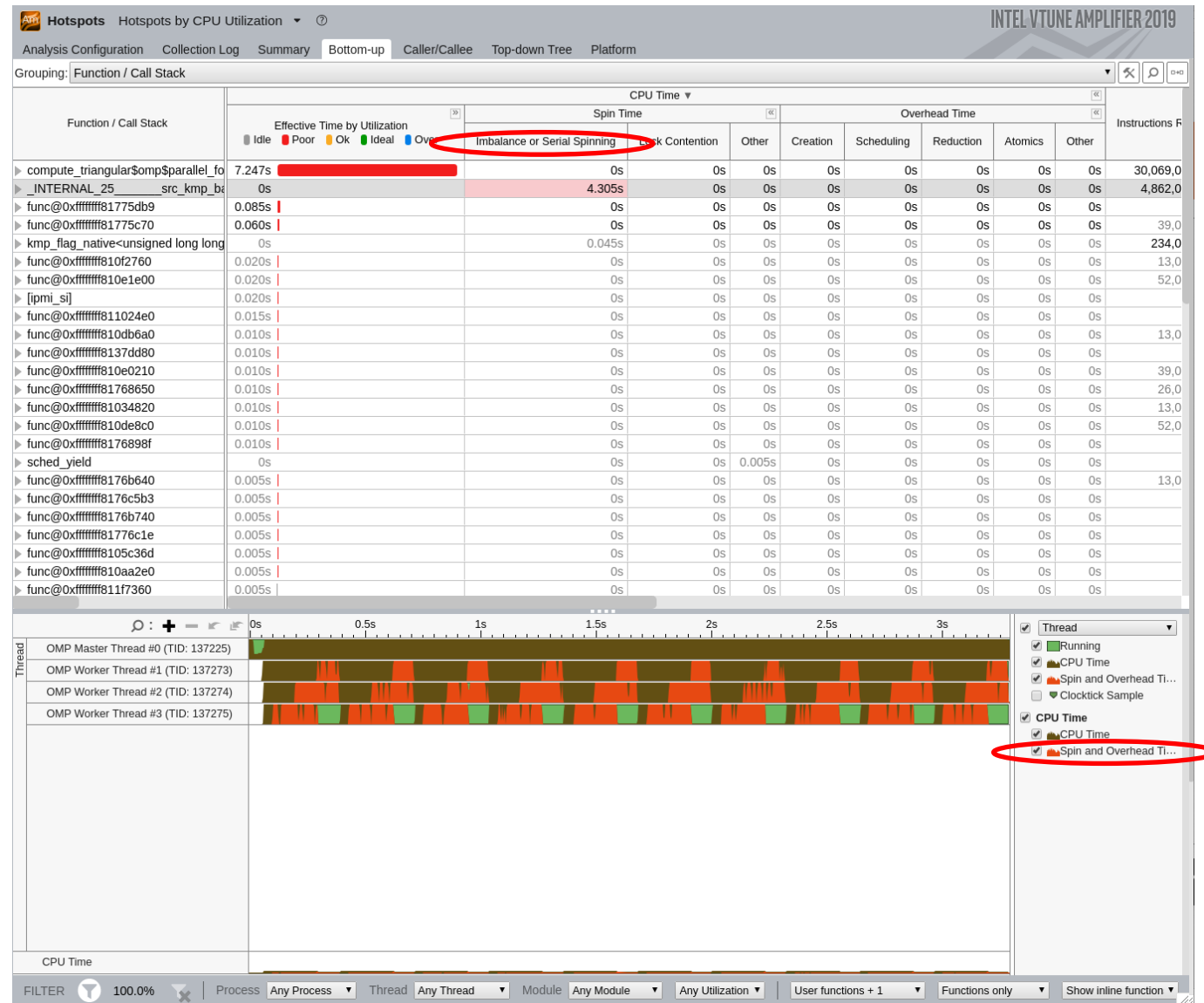


# Hotspots Analysis





# Thread Timelines Showing “Spin and Overhead”





# CPU Utilization by Threads

