intel®

# INTEL® VTUNE™ AMPLIFIER
## PERFORMANCE PROFILER
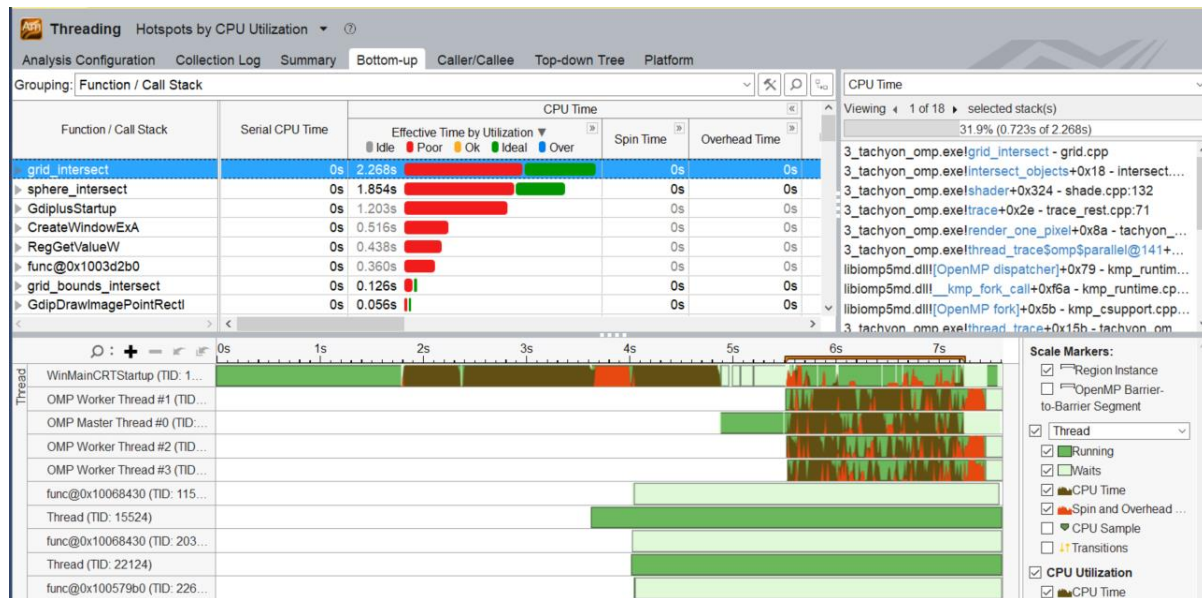
Kevin O'Leary – Intel Developer Products Division

# Intel® VTune™ Amplifier

## Tune Applications for Scalable Multicore Performance

**Agenda**

- Introduction
- Data Collection –
  Rich set of performance data
- Data Analysis -
  Find answers fast
- Flexible workflow –
  - User i/f and command line
  - Compare results
  - Remote collection
- Performance Analysis
  Details
- Summary

# Faster, Scalable Code, Faster

## Intel® VTune™ Amplifier Performance Profiler

### Accurate Data – Low Overhead
- CPU, GPU, FPU, threading, bandwidth...
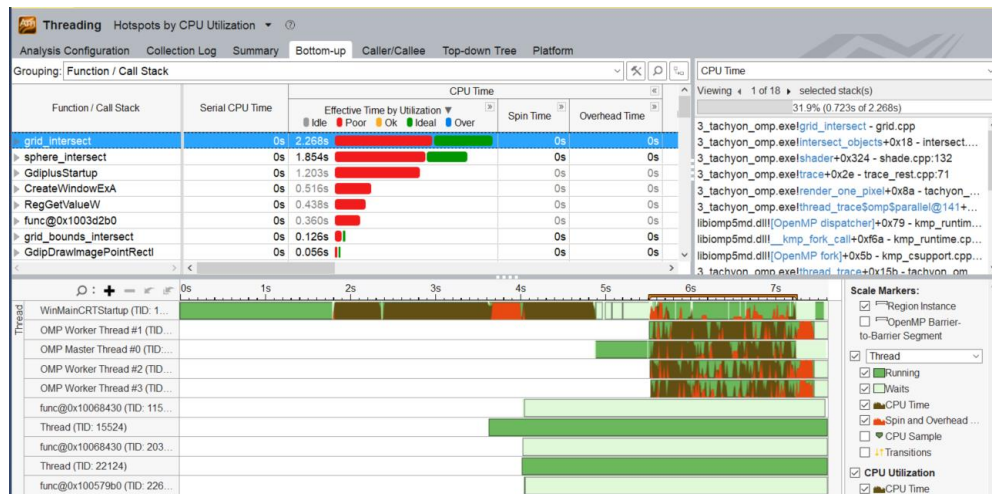
### Meaningful Analysis
- Threading, OpenMP region efficiency
- Memory access, storage device

### Easy
- Data displayed on the source code
- Easy set-up, no special compiles

"Last week, Intel® VTune™ Amplifier helped us find almost 3X performance improvement.  This week it helped us improve the performance another 3X."

Claire Cates
Principal Developer
**SAS Institute Inc.**

# Setting up a profile is easy



2. Choose Analysis Type

3. Collection options

**WHERE**

Local Host

**WHAT**

Launch Application

Specify and configure your analysis target: an application or a script to execute. Press F1 for more details.

Application:
/localdisk/temp/matrix/linux/matrix.gcc

Application parameters:

☑ Use application directory as working directory

Working directory:
/localdisk/jmarusar/temp/matrix/linux

Advanced ▶

1. What/where to profile

INTEL VTUNE AMPLIFIER 2019

*Find your analysis direction*

**Hotspots**
Want to find out where your app spends time and optimize your algorithms?

**Microarchitecture**
Want to see how efficiently your code is using the underlying hardware?

Basic Hotspots

General Exploration

Advanced Hotspots

Memory Access

Memory Consumption

**Parallelism**
Want to assess the compute efficiency of your multi-threaded app?

Concurrency    Locks and Waits

HPC Performance Characterization

4. Push Start

**HOW**

## Memory Access

Measure a set of metrics to identify memory access related issues (for example, specific for NUMA architectures). This analysis type is based on the hardware event-based sampling collection. Learn more (F1)

**CPU sampling interval, ms**

1

☑ Analyze dynamic memory objects

**Minimal dynamic memory object size to track, in bytes**

1024

☑ Evaluate max DRAM bandwidth

☐ Analyze OpenMP regions

▼ **Details**

☐ Analyze I/O waits

**Collect I/O API data**

No

☐ Collect stacks

**Stack size, in bytes**
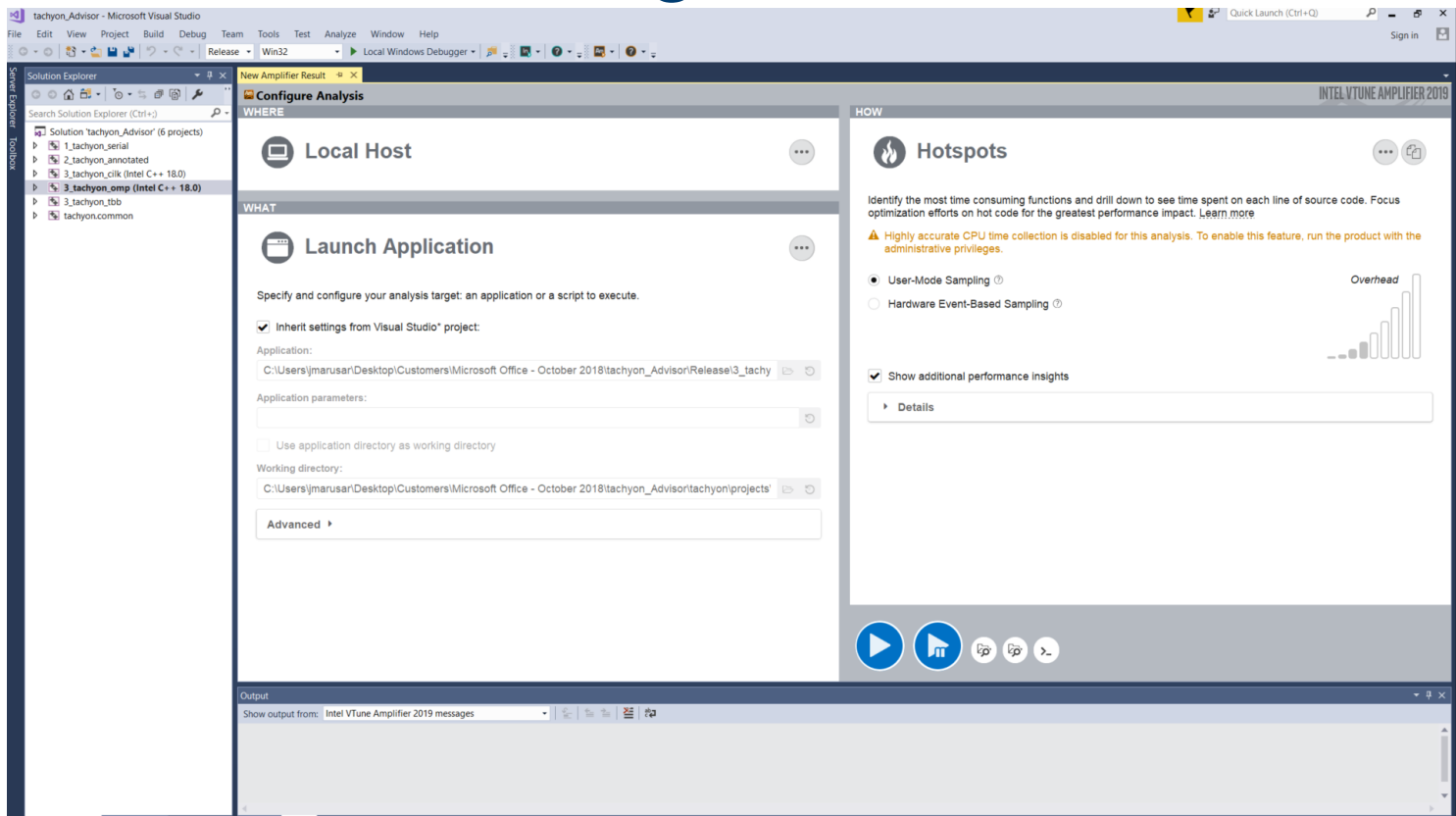
Stack type

* Full command-line also available

# Full Visual Studio* Integration

# Two Great Ways to Collect Data

Intel® VTune™ Amplifier

| Software Collector | Hardware Collector | |
|---|---|---|
| Uses OS interrupts | Uses the on chip Performance Monitoring Unit (PMU) | |
| Collects from a single process tree | Collect system wide or from a single process tree. | |
| ~10ms default resolution | ~1ms default resolution (finer granularity – finds small functions) | |
| Either an Intel® or a compatible processor | Requires a genuine Intel® processor for collection | |
| Call stacks show calling sequence | Optionally collect call stacks | |
| Works in virtual environments | Works in a VM only when supported by the VM (e.g., vSphere*, KVM) | |
| No driver required | Requires a driver | – Easy to install on Windows <br> – Linux requires root  (or use default perf driver) |

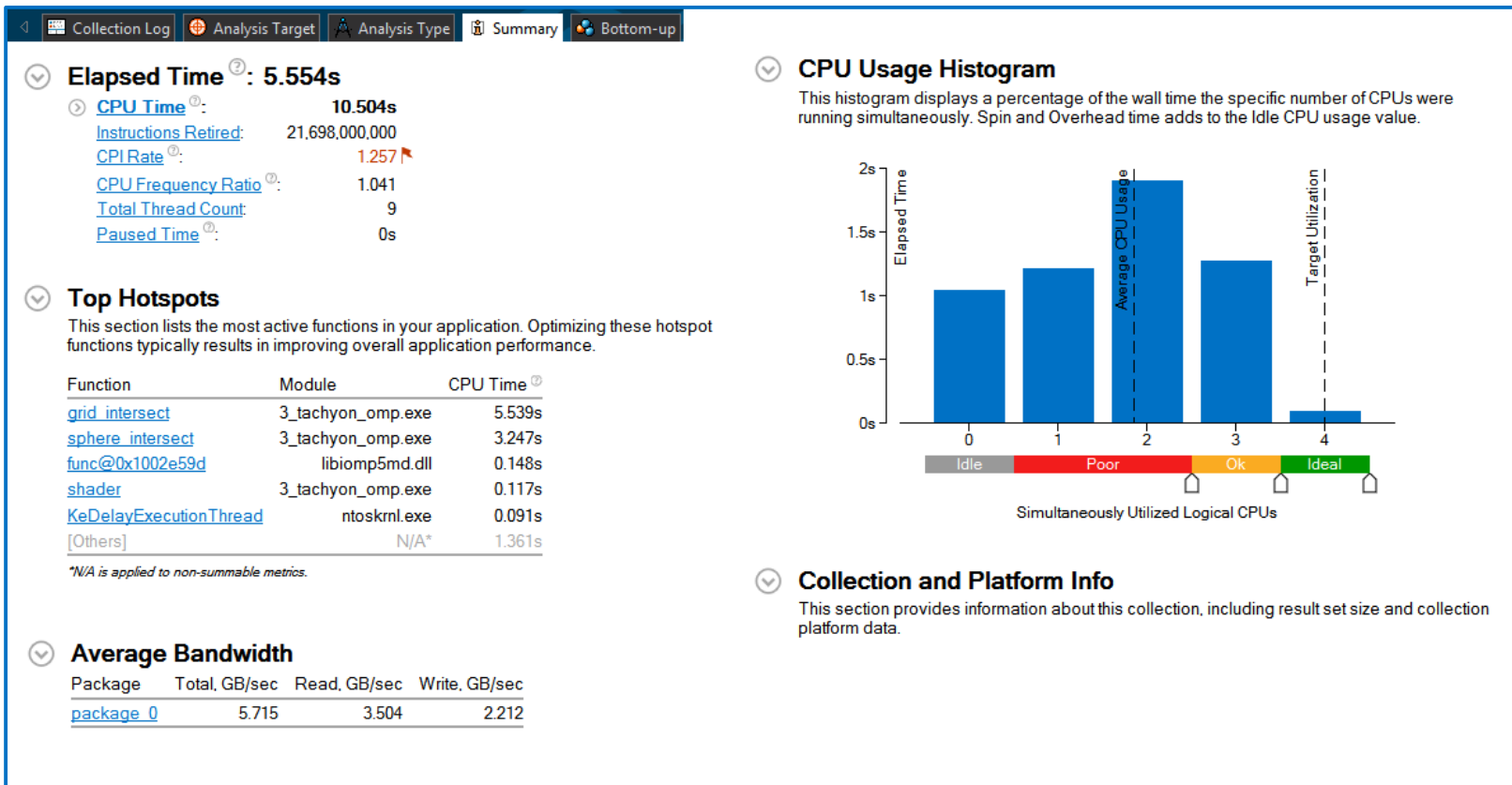**No special recompiles – C, C++, C#, Fortran, Java, Assembly**

# Example: Hotspots Analysis
## Summary View

**Elapsed Time** ⑦ : **5.554s**

| | |
|---|---|
| **CPU Time** ⑦: | **10.504s** |
| Instructions Retired: | 21,698,000,000 |
| CPI Rate ⑦: | 1.257 ⚑ |
| CPU Frequency Ratio ⑦: | 1.041 |
| Total Thread Count: | 9 |
| Paused Time ⑦: | 0s |

**Top Hotspots**

This section lists the most active functions in your application. Optimizing these hotspot functions typically results in improving overall application performance.

| Function | Module | CPU Time ⑦ |
|---|---|---|
| grid_intersect | 3_tachyon_omp.exe | 5.539s |
| sphere_intersect | 3_tachyon_omp.exe | 3.247s |
| func@0x1002e59d | libiomp5md.dll | 0.148s |
| shader | 3_tachyon_omp.exe | 0.117s |
| KeDelayExecutionThread | ntoskrnl.exe | 0.091s |
| [Others] | N/A* | 1.361s |

*N/A is applied to non-summable metrics.*

**Average Bandwidth**

| Package | Total, GB/sec | Read, GB/sec | Write, GB/sec |
|---|---|---|---|
| package_0 | 5.715 | 3.504 | 2.212 |

**CPU Usage Histogram**

This histogram displays a percentage of the wall time the specific number of CPUs were running simultaneously. Spin and Overhead time adds to the Idle CPU usage value.

Simultaneously Utilized Logical CPUs

Idle — Poor — Ok — Ideal

**Collection and Platform Info**

This section provides information about this collection, including result set size and collection platform data.

# Example: Threading Analysis
## Bottom-up View

# Example: Memory Access Analysis
## Bottom-up View



Over-Time DRAM Bandwidth

Over-Time QPI/UPI Bandwidth

Grid Breakdown by Function (configurable)

# Find Answers Fast
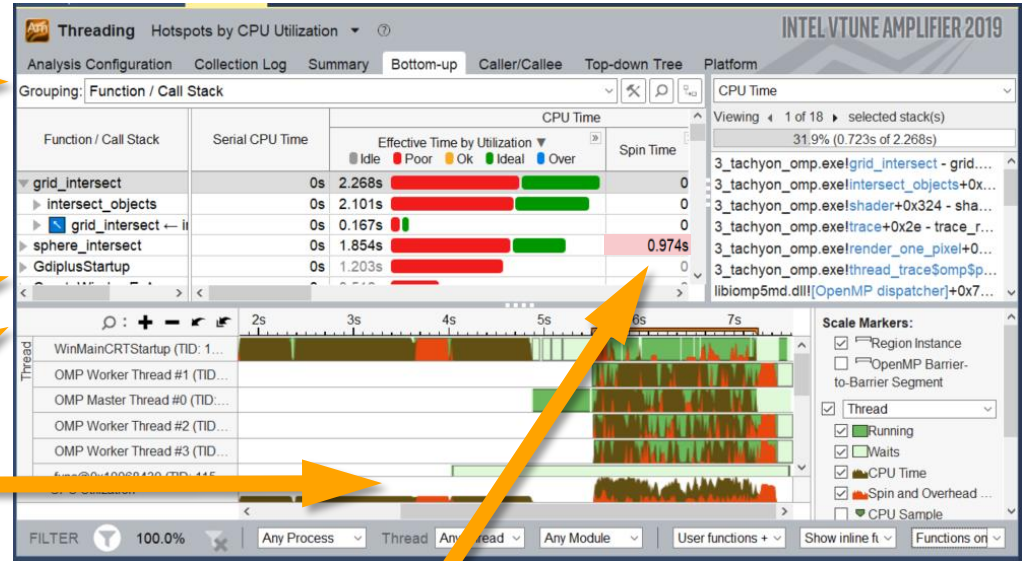
## Intel® VTune™ Amplifier



Adjust Data Grouping

- Function - Call Stack
- Module - Function - Call Stack
- Source File - Function - Call Stack
- Thread - Function - Call Stack

... (Partial list shown)

Double Click Function to View Source

Click > for Call Stack

Filter by Timeline Selection (or by Grid Selection)

Zoom In And Filter On Selection
Filter In by Selection
Remove All Filters

Filter by Process & Other Controls

Tuning Opportunities Shown in Pink. Hover for Tips

# See Profile Data On Source / Asm
## Double Click from Grid or Timeline



View Source / Asm or both

CPU Time

Right click for instruction reference manual

Quick Asm navigation:
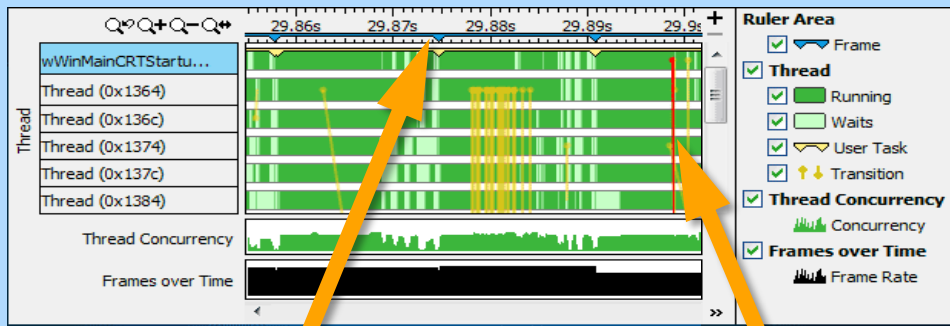Select source to highlight Asm

Scroll Bar "Heat Map" is an overview of hot spots

Click jump to scroll Asm

# Timeline Visualizes Thread Behavior

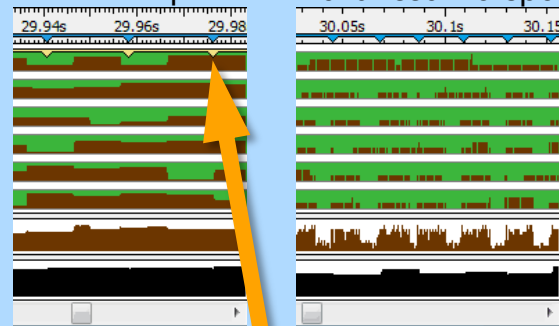Intel® VTune™ Amplifier



**Transitions**
Locks & Waits

**CPU Time**
Basic Hotspots      Advanced Hotspots

Ruler Area
- ☑ Frame
- ☑ **Thread**
  - ☑ Running
  - ☑ Waits
  - ☑ User Task
  - ☑ Transition
- ☑ **Thread Concurrency**
  - Concurrency
- ☑ **Frames over Time**
  - Frame Rate

Hovers:

Frame
Frame
Start: 29.858s Duration: 0.017s
Frame: 72
Frame Domain: Smoke::Framework::execute()
Frame Type: Good
Frame Rate: 59.8242179

Transition
Transition
wWinMainCRTStartup (0x12d4) to Thread (0x138c) (29.899s to 29.899s)
Sync Object: TBB Scheduler
Object Creation File: taskmanagertbb.cpp
Object Creation Line: 318

User Task
User Task
Start: 29.958s Duration: 0.018s
Task Type: Smoke::FrameWork::execute()::Other
Task End Call Stack: Framework::Execute

CPU Time
94.233472%

Optional: Use API to mark frames and user tasks   Frame   User Task

Optional: Add a mark during collection   Mark Timeline

# Tune OpenMP for Efficiency and Scalability

## Fast Answers:  Is My OpenMP Scalable?  How Much Faster Could It Be?



The summary view shown above gives fast answers to four important OpenMP tuning questions:

1) Is the serial time of my application significant enough to prevent scaling?
2) How much performance can be gained by tuning OpenMP?
3) Which OpenMP regions / loops / barriers will benefit most from tuning?
4) What are the inefficiencies with each region? (click the link to see details)

# Command Line Interface

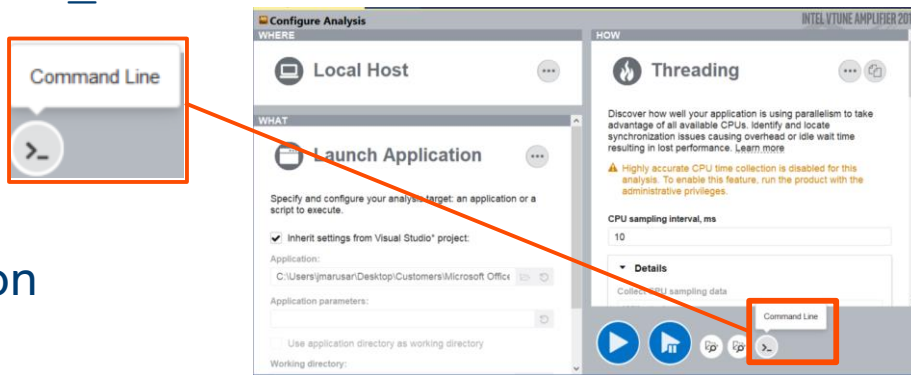Automate analysis

## amplxe-cl is the command line:

- **Windows:** `C:\Program Files (x86)\IntelSWTools\VTune Amplifier\bin[32|64]\amplxe-cl.exe`
- **Linux:** `/opt/intel/vtune_amplifier/bin[32|64]/amplxe-cl`

Help: `amplxe-cl –help`

## Use UI to setup
1) Configure analysis in UI
2) Press "Command Line..." button
3) Copy & paste command



**Great for regression analysis – send results file to developer
Command line results can also be opened in the UI**

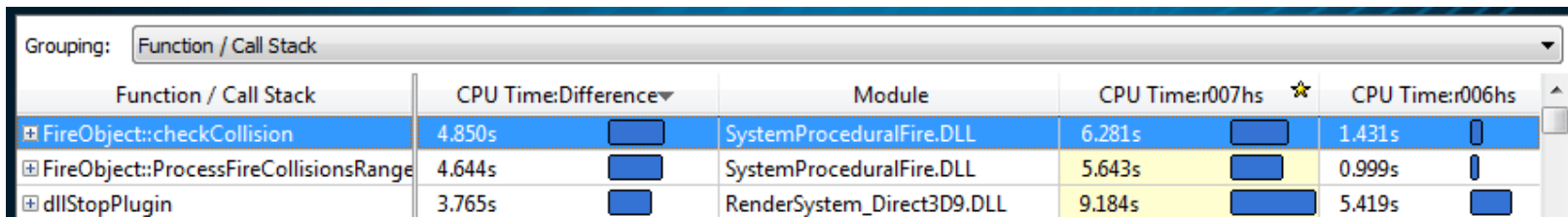# Compare Results Quickly – Sort By Difference

Intel® VTune™ Amplifier

Quickly identify cause of regressions.

- Run a command line analysis daily

- Identify the function responsible so you know who to alert

Compare 2 optimizations – What improved?

Compare 2 systems – What didn't speed up as much?

| Grouping: | Function / Call Stack | | | | | |
|---|---|---|---|---|---|---|
| Function / Call Stack | CPU Time:Difference▼ | | Module | CPU Time:r007hs ⭐ | | CPU Time:r006hs |
| ⊞ FireObject::checkCollision | 4.850s | ▭ | SystemProceduralFire.DLL | 6.281s | ▭ | 1.431s ▯ |
| ⊞ FireObject::ProcessFireCollisionsRange | 4.644s | ▭ | SystemProceduralFire.DLL | 5.643s | ▭ | 0.999s ▯ |
| ⊞ dllStopPlugin | 3.765s | ▭ | RenderSystem_Direct3D9.DLL | 9.184s | ▭ | 5.419s ▭ |

# Optimize Memory Access
## Memory Access Analysis - Intel® VTune™ Amplifier 2017
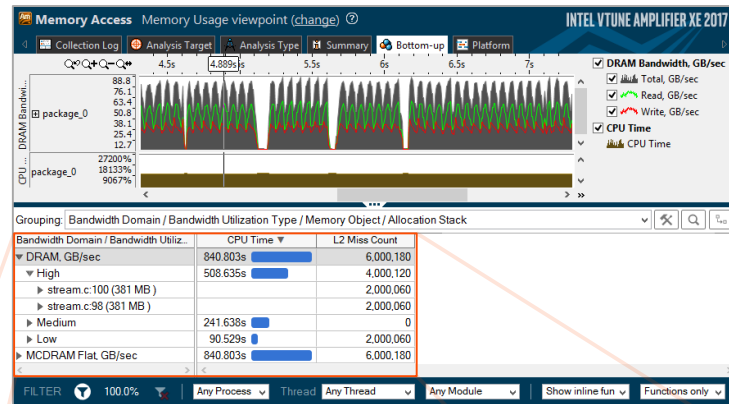
## Tune data structures for performance

- Attribute cache misses to data structures (not just the code causing the miss)
- Support for custom memory allocators

## Optimize NUMA latency & scalability

- True & false sharing optimization
- Auto detect max system bandwidth
- Easier tuning of inter-socket bandwidth

## Easier install, Latest processors

- No special drivers required on Linux*
- Intel® Xeon Phi™ processor MCDRAM (high bandwidth memory) analysis



| Bandwidth Domain / Bandwidth Utiliz... | CPU Time ▼ | L2 Miss Count |
|---|---|---|
| ▼ DRAM, GB/sec | 840.803s | 6,000,180 |
| ▼ High | 508.635s | 4,000,120 |
| ▶ stream.c:100 (381 MB ) | | 2,000,060 |
| ▶ stream.c:98 (381 MB ) | | 2,000,060 |
| ▶ Medium | 241.638s | 0 |
| ▶ Low | 90.529s | 2,000,060 |
| ▶ MCDRAM Flat, GB/sec | 840.803s | 6,000,180 |

# Memory Object Identification

# Memory Object Identification



Assembly view also available

Double-click to see allocation site in source view

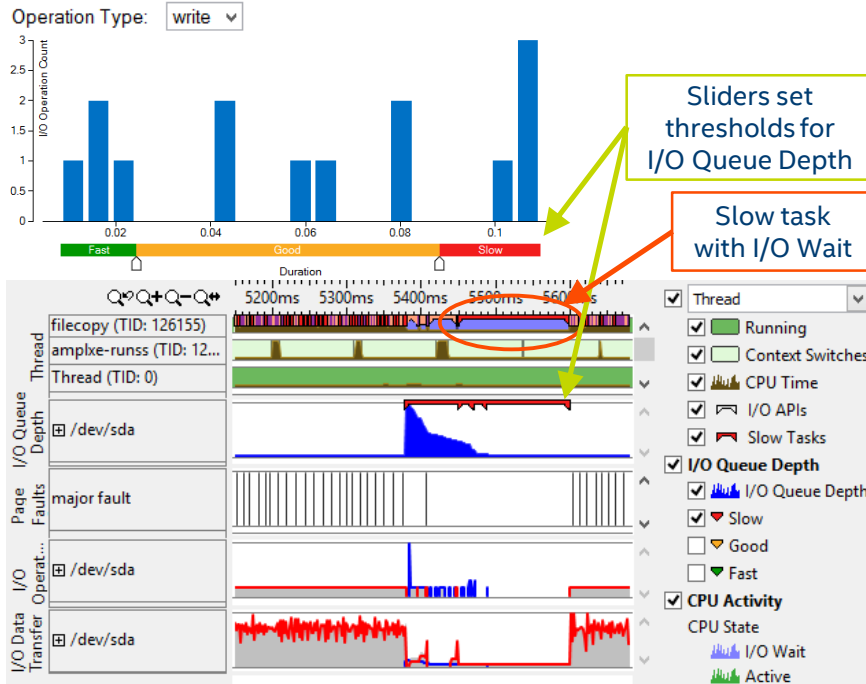# Storage Device Analysis (HDD, SATA or NVMe SSD)
## Intel® VTune™ Amplifier

### Are You I/O Bound or CPU Bound?

- Explore imbalance between I/O operations (async & sync) and compute
- Storage accesses mapped to the source code
- See when CPU is waiting for I/O
- Measure bus bandwidth to storage

### Latency analysis

- Tune storage accesses with latency histogram
- Distribution of I/O over multiple devices



Disk Input and Output Histogram

Sliders set thresholds for I/O Queue Depth

Slow task with I/O Wait

# A Quick Question for the Audience

# TANGENT ON A COUPLE OTHER TOOLS

# INTEL® ADVISOR 2019

## VECTORIZATION OPTIMIZATION AND THREAD PROTOTYPING

- Vectorization Advisor
- Threading Advisor
- Flow Graph Analyzer

# Get Faster Code Faster!  Intel® Advisor
## Thread Prototyping

## Have you:

- Threaded an app, but seen little benefit?
- Hit a "scalability barrier"?
- Delayed release due to sync. errors?

## Data Driven Threading Design:

- Quickly prototype multiple options
- Project scaling on larger systems
- Find synchronization errors before implementing threading
- Design without disrupting development

**Add Parallelism with Less Effort, Less Risk and More Impact**



Scalability of Maximum Site Gain

"**Intel® Advisor** has allowed us to quickly prototype ideas for parallelism, saving developer time and effort"

*Simon Hammond*
*Senior Technical Staff*
***Sandia National Laboratories***

http://intel.ly/advisor-xe

# Get Faster Code Faster!  Intel® Advisor
## Vectorization Optimization

**Have you:**

- Recompiled for AVX2 with little gain
- Wondered where to vectorize?
- Recoded intrinsics for new arch.?
- Struggled with compiler reports?

**Data Driven Vectorization:**

- What vectorization will pay off most?
- What's blocking vectorization?  Why?
- Are my loops vector friendly?
- Will reorganizing data increase performance?
- Is it safe to just use pragma simd?

| | | | Perfor... Issues | Self Time▼ | Total Time | Type | Why No Vectorization? | Vect... | Efficiency | Gain... | VL.. | Com... | Traits | Da |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ⊞ ⟳ [loop in main at roofline.cpp:295] | ☐ | | | 18.538s | 18.538s | Vectorized (B... | | AVX | ~100% | 5.34x | 4 | 5.34x | | Flo |
| ⊞ ⟳ [loop in main at roofline.cpp:310] | ☐ | | | 18.394s | 18.394s | Vectorized (Bo... | | AVX | ~100% | 5.34x | 4 | 5.34x | | Flo |
| ⊞ ⟳ [loop in main at roofline.cpp:221] | ☑ | | | 14.741s | 14.741s | Scalar | ⊟ novector dire... | | | | | | | Flo |
| ⊞ ⟳ [loop in main at roofline.cpp:234] | ☐ | | | 11.117s | 11.117s | Scalar | ⊟ inner loop w... | | | | | | | Flo |
| ⊞ ⟳ [loop in main at roofline.cpp:247] | ☐ | | | 6.967s | 6.967s | Vectorized (Bo... | | AVX | ~31% | 1.22x | 4 | 1.22x | Inserts; U... | Flo |
| ⊞ ⟳ [loop in main at roofline.cpp:138] | ☐ | | | 6.949s | 6.949s | Scalar | ⊟ novector dire... | | | | | | | Flo |
| ⊞ ⟳ [loop in main at roofline.cpp:260] | ☐ | | | 3.285s | 3.285s | Vectorized (Bo... | | AVX | ~100% | 5.09x | 4 | 5.09x | | Flo |
| ⊞ ⟳ [loop in main at roofline.cpp:199] | ☐ | | | 2.454s | 2.454s | Vectorized (Bo... | | AVX | ~100% | 5.14x | 4 | 5.14x | | Flo |
| ⊞ ⟳ [loop in main at roofline.cpp:273] | ☐ | | | 2.258s | 2.258s | Vectorized (Bo... | | AVX2 | ~100% | 4.73x | 4 | 4.73x | FMA | Flo |
| ⊞ ⟳ [loop in main at roofline.cpp:151] | ☐ | | | 1.899s | 1.899s | Vectorized (Bo... | | AVX | ~100% | 4.80x | 4 | 4.80x | | Flo |
| ⊞ ⟳ [loop in main at roofline.cpp:256] | ☐ | | ⚠ 1 Oppo... | 0.042s | 3.327s | Scalar | ⊟ inner loop w... | | | | | | | Flo |
| ⊞ ⟳ [loop in main at roofline.cpp:304] | | | | 0.040s | 18.434s | Scalar | ⊟ inner loop w... | | | | | | | |

Elapsed time: 125.72s  Vectorized  Not Vectorized   OFF  Smart Mode

INTEL ADVISOR 2019

FILTER: All Modules  All Sources  Loops And Functions  All Threads

Summary  Survey & Roofline  Refinement Reports

Function Call Sites and Loops    Vectorized Loops    Instruction Set

"Intel® Advisor's Vectorization Advisor permitted me to focus my work where it really mattered.  When you have only a limited amount of time to spend on optimization, it is invaluable."

*Gilles Civario*
*Senior Software Architect*
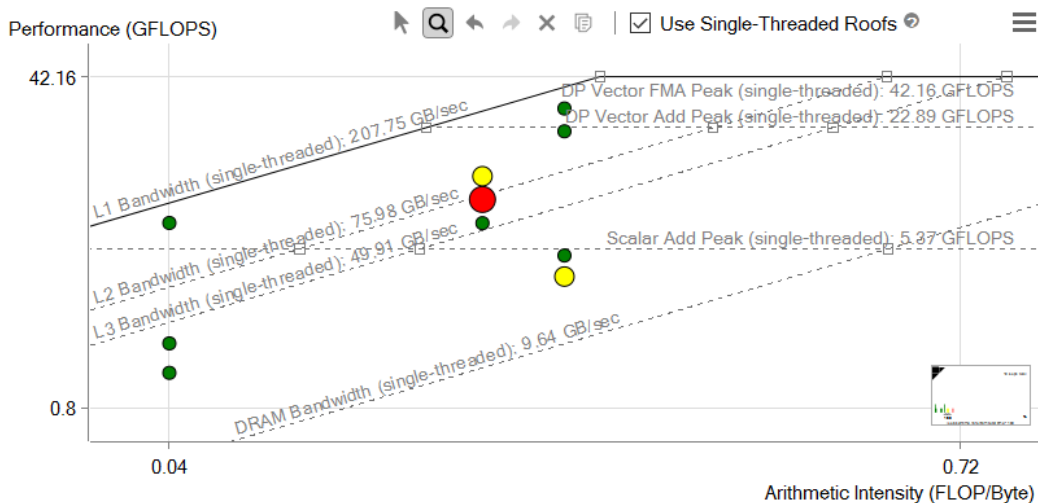***Irish Centre for High-End Computing***

# What is a Roofline Chart?

A Roofline Chart plots application performance against hardware limitations.

- Where are the bottlenecks?

- How much performance is being left on the table?

- Which bottlenecks can be addressed, and which *should* be addressed?

- What's the most likely cause?

- What are the next steps?



Roofline first proposed by University of California at Berkeley:
*Roofline: An Insightful Visual Performance Model for Multicore Architectures*, 2009
Cache-aware variant proposed by University of Lisbon:
*Cache-Aware Roofline Model: Upgrading the Loft*, 2013

# Debug Memory & Threading with Intel® Inspector
## Find & Debug Memory Leaks, Corruption, Data Races, Deadlocks



## Correctness Tools Increase ROI by 12%-21%[1]
- Errors found earlier are less expensive to fix
- Races & deadlocks not easily reproduced
- Memory errors are hard to find without a tool

## Debugger Integration Speeds Diagnosis
- Breakpoint set just before the problem
- Examine variables and threads with the debugger

## What's New in 2019 Release
Find Persistent Memory Errors
- Missing / redundant cache flushes
- Missing store fences
- Out-of-order persistent memory stores
- PMDK transaction redo logging errors

Learn More: bit.ly/intel-inspector

# Debug Memory & Threading Errors

Intel® Inspector

## Find and eliminate errors

- Memory leaks, invalid access…
- Races & deadlocks
- C, C++ and Fortran  (or a mix)

## Simple, Reliable, Accurate

- No special recompiles
  Use any build, any compiler[1]
- Analyzes dynamically generated or linked code
- Inspects 3rd party libraries without source
- Productive user interface + debugger integration
- Command line for automated regression analysis

Clicking an error instantly displays source
code snippets and the call stack

**Fits your existing process**

# DIVING DEEPER INTO ANALYSIS
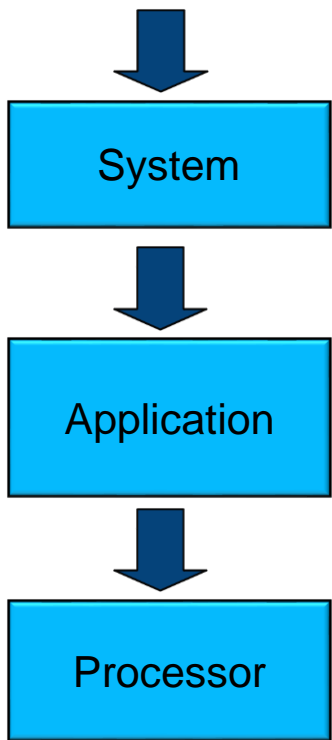
# Introduction to Performance Tuning

**System**
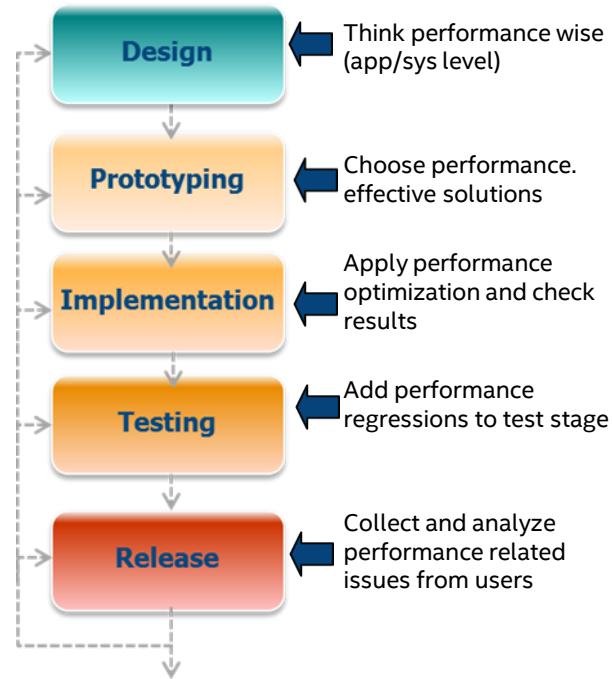
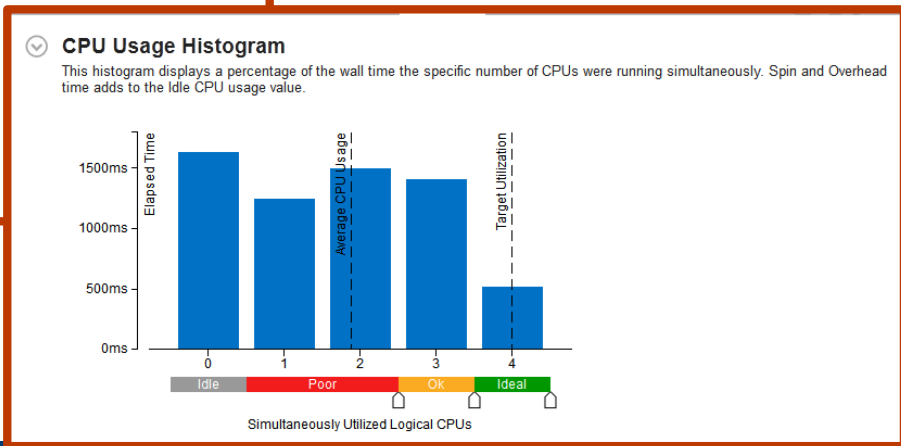H/W tuning:
BIOS (TB, HT)
Memory
Network I/O
Disk I/O

OS tuning:
Page size
Swap file
RAM Disk
Power settings
Network protocols

**Application**

Better application design:
Parallelization
Fast algorithms / data bases
Programming language and RT libs
Performance libraries
Driver tuning

**Processor**

Tuning for Microarchitecture:
Compiler settings/Vectorization
Memory/Cache usage
CPU pitfalls

**Design** — Think performance wise (app/sys level)

**Prototyping** — Choose performance. effective solutions

**Implementation** — Apply performance optimization and check results

**Testing** — Add performance regressions to test stage

**Release** — Collect and analyze performance related issues from users

(intel)

# System-Level Profiling – High-level Overviews

# System-Level Profiling – Process/Module Breakdowns

# System-Level Profiling – Disk I/O Analysis

**Are You I/O Bound or CPU Bound?**

- Explore imbalance between I/O opera (async & sync) and compute
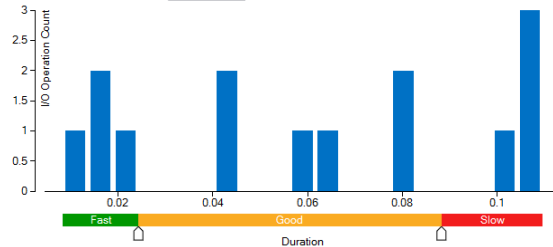- Storage accesses mapped to the source code

**See when CPU is waiting for I/O**

- Measure bus bandwidth to storage
- Latency analysis
- Tune storage accesses with latency histogram
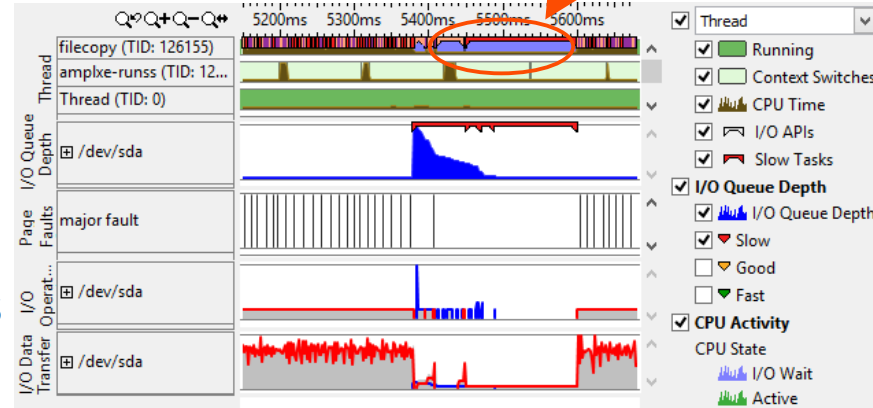- Distribution of I/O over multiple devices

```
> amplxe-cl –collect disk-io –d 10
```



**Disk Input and Output Histogram**

Slow task with I/O Wait

# System-Level Profiling – HPC Characterizaton

## Three Metric Classes

- **CPU Utilization**
  - Logical core % usage
  - Includes parallelism and OpenMP information
- **Memory Bound**
  - Break down each level of the memory hierarchy
- **FPU Utilization**
  - Floating point GFLOPS and density

```
> amplxe-cl -collect hpc-performance -d 10
```

**CPU Utilization** : 60.9%

| | | |
|---|---|---|
| Average CPU Usage | 14.611 Out of 24 logical CPUs | |
| Serial Time | 0.013s (0.1%) | |
| **Parallel Region Time** | **11.986s (99.9%)** | |
| Estimated Ideal Time | 8.205s (68.4%) | |
| OpenMP Potential Gain | 3.781s (31.5%) | |

The time wasted on load imbalance or parallel work arrangement is significant and negatively impacts the application performance and scalability. Explore OpenMP regions with the highest metric values. Make sure the workload of the regions is enough and the loop schedule is optimal.

**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...

OpenMP...
conj_grad...
MAIN__$o...
MAIN__$o...
MAIN__$o...
MAIN__$o...
[Others]

*N/A is app...

**Memory Bound** : 91.8%

Cache Bound :
DRAM Latency Bound :
DRAM Bandwidth Bound :

This metric represents a fraction of ...
main memory (DRAM). This metric d...
Consider improving data locality in N...

NUMA: % of Remote Accesses :

A significant amount of DRAM loads ...
same core, or at least the same pack...

**FPU Utilization** : 1.3%

| | |
|---|---|
| SP FLOPs per Cycle | 0.211 Out of 16 |
| Vector Capacity Usage | 48.3% |
| FP Instruction Mix: | |
| % of Packed FP Instr. | 93.1% |
| % of 128-bit | 93.1% |
| % of 256-bit | 0.0% |
| % of Scalar FP Instr. | 6.9% |
| FP Arith/Mem Rd Instr. Ratio | 0.264 |
| FP Arith/Mem Wr Instr. Ratio | 6.298 |

**Top 5 hotspot loops (functions) by FPU usage**

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 575 in conj_grad_$omp$parallel@517] | 126.149s | 1.6% | SSE2(128) | Body |
| [Loop at line 678 in conj_grad_$omp$parallel@517] | 5.004s | 1.7% | SSE2(128) | Body |
| [Loop at line 575 in conj_grad_$omp$parallel@517] | 2.678s | 2.1% | [Unknown] | Remainder |
| [Loop at line 573 in conj_grad_$omp$parallel@517] | 0.995s | 4.0% | SSE2(128) | Body |
| [Loop at line 661 in conj_grad_$omp$parallel@517] | 0.952s | 1.3% | SSE(128); SSE2(128) | Body |
| [Others] | 2.437s | N/A* | N/A* | N/A* |

*N/A is applied to non-summable metrics.

intel

# System-Level Profiling – Memory Bandwidth



Find areas of high and low bandwidth usage. Compare to max system bandwidth based on Stream benchmarks.

```
-knob collect-memory-bandwidth=true
```
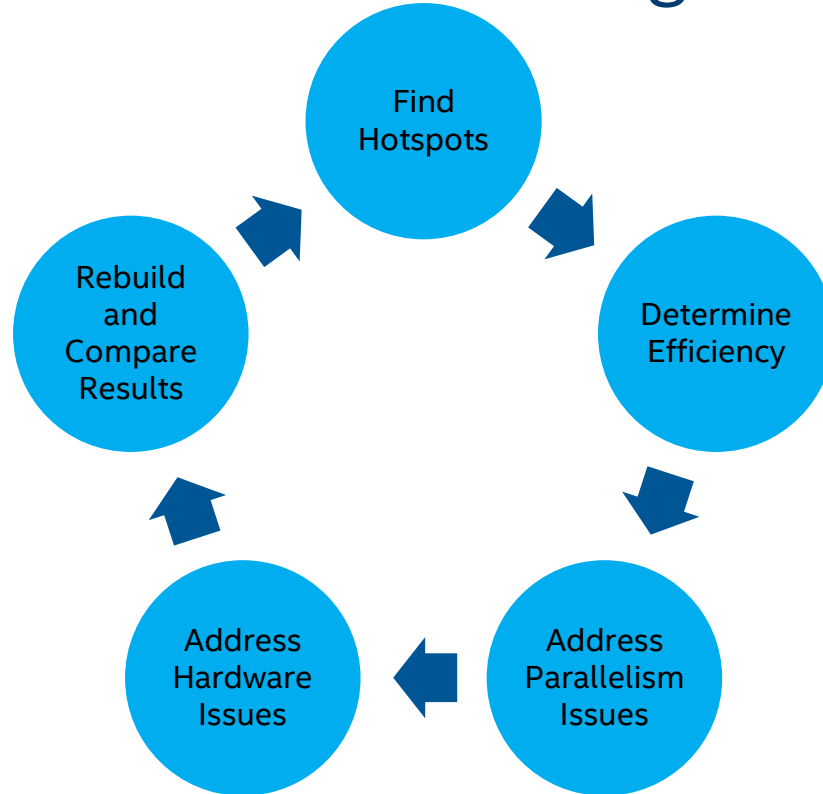
# Application Performance Tuning Process

# Find Hotspots



```
> amplxe-cl -collect hotspots -- ./myapp.out
```

# Find Hotspots

- Drill to source or assembly
- Hottest areas easy to ID
- Is this the expected behavior
- Pay special attention to loops and memory accesses

- Learn how your code behaves
- What did the compiler generate
- What are the expensive statements

# Determine Efficiency



Look for Parallelism, Cycles-per-Instruction (CPI), and Retiring %

# Address Parallelism Issues

- Use Concurrency Analysis to ensure you're using all your threads as often as possible.

- Common concurrency problems can often be diagnosed in the timeline.

- Switch to the Locks And Waits viewpoint or run a Locks and Waits analysis to investigate contention.

Coarse-Grain Locks



Thread Imbalance



High Lock Contention

(intel)

# Address Hardware Issues



The X86 Processor Pipeline (simplified)

# Pipeline Slot Categorization

- Pipeline slots can be sorted into one of four categories on a given cycle by what happens to the uop in that slot.
  - Retiring
  - Bad Speculation
  - Back End Bound
  - Front End Bound

- Each category has an expected range of values in a well tuned application.



| App. Type: Category | | Client/Desktop | Server/Database/ Distributed | High Performance Computing |
|---|---|---|---|---|
| Retiring | ▲ | 20-50% | 10-30% | 30-70% |
| Bad Speculation | ▼ | 5-10% | 5-10% | 1-5% |
| Front End Bound | ▼ | 5-10% | 10-25% | 5-10% |
| Back End Bound | ▼ | 20-40% | 20-60% | 20-40% |

# The uop Pipeline
## Categorizing the hotspots

- Modern CPUs "pipeline" instructions. This pipeline can be generally divided into two sections.

  - The Front End fetches instructions, decodes them into uops, and allocates them to…

  - The Back End, which is responsible for executing the uops. Once successfully completed, a uop is considered "retired".

- A Pipeline Slot is an abstract representation of the hardware resources needed to process a uop.

- The front end can only allocate so many uops per cycle, and the same is true of the back end and retiring them. This determines the number of Pipeline Slots. As a general rule, this number is four.

# Pipeline Slot Categorization
## Retiring

This is the good category! You want as many of your slots in this category as possible. However, even here there may be room for optimization.



FRONT END

BACK END

EXECUTION UNIT

RETIREMENT

uop

uop

uop

uop

Fetch & Decode Instructions, Predict Branches

Re-order and Execute Instructions

Commit Results to Memory

# Pipeline Slot Categorization
## Bad Speculation

This occurs when a uop is removed from the back end without retiring; effectively, it's cancelled, most often because a branch was mispredicted.



FRONT END

BACK END

EXECUTION UNIT

RETIREMENT

uop
uop
uop
uop

Fetch & Decode Instructions, Predict Branches

Re-order and Execute Instructions

Commit Results to Memory

# Pipeline Slot Categorization
## Back End Bound

This is when the back end can't accept uops, even if the front end can send them, because it already contains uops waiting on data or long execution.

# Pipeline Slot Categorization
## Front End Bound

This is when the front end can't deliver uops even though the back end can take them, usually due to delays in fetching code or decoding instructions.



FRONT END

BACK END

EXECUTION UNIT

RETIREMENT

uop

uop

Fetch & Decode Instructions, Predict Branches

Re-order and Execute Instructions

Commit Results to Memory

# Identifying and Diagnosing Inefficiency
## *Microarchitecture Analysis*

```
> amplxe-cl -collect uarch-exploration -- ./myapp.out
```

- Microarchitecture Exploration (previously General Exploration) is a hardware events analysis. It is preconfigured to sample the appropriate events on your architecture and calculates the proper metrics from them.

- Potential tuning opportunities are highlighted in pink.

- To check the efficiency of a hotspot, look at the Retiring metric. If it's less than the expected number for your application type, it's probably inefficient.
  - Hotspots with high retiring values may still have room for improvement.

| App Type | Expected |
|----------|----------|
| Client/ Desktop | **20-50%** |
| Server/ Database/ Distributed | **10-30%** |
| HPC | **30-70%** |

| ◄ | Collection Log | Analysis Target | Analysis Type | Summary | Bottom-up | Event Count | Platform | ▷ |

Grouping: Function / Call Stack

| Function / Call Stack | Instructions Retired | CPI Rate | Front-End Bound » | Bad Speculation » | Back-End Bound » | Retiring » |
|---|---|---|---|---|---|---|
| ▶ initialize_2D_buffer | 85,219,200,000 | 0.266 | 0.5% | 0.0% | 0.0% | 100.0% |
| ▶ grid_intersect | 10,963,200,000 | 0.706 | 4.6% | 15.1% | 46.4% | 33.9% |
| ▶ sphere_intersect | 10,946,400,000 | 0.601 | 2.1% | 1.6% | 47.5% | 48.8% |
| ▶ grid_bounds_intersect | 480,000,000 | 1.105 | 13.0% | 2.2% | 52.3% | 32.5% |
| ▶ tri_intersect | 216,000,000 | 0.789 | 0.0% | 20.2% | 39.3% | 40.5% |

# Categorizing and Correcting Inefficiencies
## *Microarchitecture Exploration Analysis*

- Intel® VTune™ Amplifier has hierarchical expanding metrics categorized by the four slot types.

- You can expand your way down, following the hotspot, to identify the root cause of the inefficiency.

  - Sub-metrics highlight pink on their own merits, just like top level metrics.

- Hovering over a metric produces a helpful, detailed tooltip (not shown).

  - There are tooltips on Summary tabs too: hover over any ⑦ icon.

# Categorizing and Correcting Inefficiencies
## Retiring: *Microarchitecture Exploration Analysis, Intel® Advisor*

- High Retiring percentage is generally good, but may be inefficient if you're doing work that doesn't need to be done at all, or could be done faster.

- Retiring can be split based on whether the uops being retired came from the microcode sequencer or not.

  - **Yes?** Try reworking code to avoid microcode assists.

  - **No?** Make sure the code is well vectorized.

| Retiring | | |
|---|---|---|
| General Retirement | | Microc... |
| FP Arithmetic | Other | Assists |
| 0.0% | 100.0% | 0.0% |
| 24.6% | 75.4% | 0.0% |
| 19.0% | 81.0% | 0.0% |
| 0.0% | 100.0% | 0.0% |
| 16.7% | 83.3% | 0.0% |
| 20.0% | 80.0% | 0.0% |
| 10.0% | 90.0% | 0.0% |

**Tip:**

Use Vectorization Advisor to fine-tune your vectorization.

# HPC Characterization: FPU Utilization

## FPU utilization

% of FPU load (100% - FPU is fully loaded, threshold 50%)

Calculation based on PMU events representing scalar and packed single and double precision SIMD instructions

## Metrics in FPU utilization section

FLOPs broken down by scalar and packed

Instruction Mix

Top 5 loops/functions by FPU usage

- Detected with static binary analysis

Vectorized vs. Non-vectorized, ISA, and characterization detected by static analysis and Intel Compiler diagnostics





# HARDWARE IS BECOMING MORE VECTORIZED, SO SHOULD YOU!

# Categorizing and Correcting Inefficiencies
Bad Speculation: *Microarchitecture Exploration Analysis*

- Bad Speculation is caused by either Machine Clears or Branch Mispredicts.

  - **Machine Clears** can be caused by self-modifying code, etc.

  - **Branch mispredicts** are more common. These occur when the paths taken by `if`, `switch`, `for`, `do-while`, and other conditional branches are incorrectly predicted and the uops have to be thrown out.

| Bad Speculation | «|
|---|---|
| Branch Mispredict | Machine Clears |
| 0.0% | 0.0% |
| 15.1% | 0.0% |
| 0.0% | 1.6% |
| 2.2% | 0.0% |
| 0.0% | 20.2% |

- Use Intel® VTune™ Amplifier's Source Viewer to identify problematic branches.

- Avoid unnecessary branching:

  - Remove branches entirely if possible

  - Move branches outside of loops if possible.

| Source | Assembly | | | | | | | | | | Assembly grouping: | Address |

| S. Li. | Source | Bad Spec... |
|---|---|---|
| 580 | `while (cur != NULL) {` | 0.1% |
| 581 | `if (ry->mbox[cur->obj->id] != ry->serial) {` | 4.3% |
| 582 | `ry->mbox[cur->obj->id] = ry->serial;` | 1.3% |

# Categorizing and Correcting Inefficiencies
## Front End: *Microarchitecture Exploration Analysis*

| Front-End Bound | | | | | | | | | «|
|---|---|---|---|---|---|---|---|---|---|
| Front-End Latency | | | | | | Front-End Bandwidth | | | |
| ICache Misses | ITLB Overhead | Branch Resteers | DSB Switches | Length ... | MS Switches | Front-End Bandwidth MITE | Front-End Bandwidth DSB | Front-End Bandwidth LSD | |
| 0.0% | 0.0% | 0.2% | 0.6% | 0.0% | 0.0% | 0.6% | 6.5% | 0.0% | |
| 0.0% | 0.1% | 1.8% | 1.2% | 0.0% | 0.0% | 1.2% | 10.1% | 1.8% | |

- Front End Bound pipeline slots are common in JIT or interpreted code.

- Front End Bound can be bandwidth or latency:

  - **Bandwidth** issues are caused by inefficient instruction decoding, or restrictions in caching decoded instructions, etc.

  - **Latency** is caused by instruction cache misses, delays in instruction fetching after branch mispredicts, switching to the microcode sequencer too often, etc.

# Categorizing and Correcting Inefficiencies
## Back End: *Microarchitecture Exploration Analysis, Memory Bandwidth*

| | | Back-End Bound | | | | | | | | | | | | « |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Memory Bound | | | | | | | | | | « | Core Bound | « |
| | | L3 Bound | | | | « | DRAM Bound | | « | Store Bound | | | « | | |
| L1 Bound » | L2 Bound | Contested Acc... | Data Sharing | L3 Latency | SQ Full | Memory Band... | Memory Lat... « / LLC Miss | Store Latency | False Shari... | Split Sto... | DTLB Store ... | Divider | Port Utilization » |
| 3.2% | | 0.0% | 0.0% | 0.0% | 0.0% | 0.2% | 0.0% | 3.3% | 0.0% | 0.0% | 0.2% | 0.0% | 26.6% |
| 11.3% | 4.8% | 0.0% | 0.0% | 100.0% | 0.0% | 9.5% | 0.0% | 1.1% | 0.0% | 0.2% | 0.2% | 4.8% | 17.2% |

- Back End bound is the most common bottleneck type for most applications.

- It can be split into Core Bound and Memory Bound
  - **Core Bound** includes issues like not using execution units effectively and performing too many divides.
  - **Memory Bound** involves cache misses, inefficient memory accesses, etc.
    - Store Bound is when load-store dependencies are slowing things down.
    - The other sub-categories involve caching issues and the like. Memory Access Analysis may provide additional information for resolving this performance bottleneck.

# Rebuild and Compare Results

# Example: Poor NUMA Utilization

Memory Access  Memory Usage viewpoint (change) ?

◄ ☐ Collection Log  ⊕ Analysis Target  A Analysis Type  🗎 Summary  ⬤ Bottom-up  🖿 Platform

## ⊙ Elapsed Time ? : 32.626s

| | |
|---|---|
| CPU Time ? : | 508.508s |
| ⊙ Memory Bound ? : | 73.3% 🚩 of Pipeline Slots |
| L1 Bound ? : | 8.4% 🚩 of Clockticks |
| L2 Bound ? : | 0.0% of Clockticks |
| L3 Bound ? : | 7.4% 🚩 of Clockticks |
| ⊙ DRAM Bound ? : | 8.0% of Clockticks |
| DRAM Bandwidth Bound ? : | 0.0% of Elapsed Time |
| ⊙ Memory Latency: | |
| Remote / Local DRAM Ratio ? : | 0.000 |
| Local DRAM ? : | 1.2% of Clockticks |
| Remote DRAM ? : | 0.0% of Clockticks |
| Remote Cache ? : | 11.4% of Clockticks |
| Loads: | 163,640,209,059 |
| Stores: | 34,303,629,078 |
| ⊙ LLC Miss Count ? : | 331,669,899 |
| Average Latency (cycles) ? : | 12 |
| Total Thread Count ? : | 17 |
| Paused Time ? : | 0s |

If Memory Bound is high and local caches are not the problem

Focus on "Remote" metrics

# Example: Poor NUMA Utilization



Look for areas of high QPI/UPI bandwidth

## QPI/UPI BANDWIDTH IS COMMUNICATION BETWEEN THE SOCKETS. THIS MAY INDICATE SOME SORT OF NUMA ISSUE.

# EXAMPLE: POOR NUMA UTILIZATION

Common causes of poor NUMA utilization

- Allocation vs. first touch memory location

- False sharing of cache lines
  - Use padding when necessary

- Arbitrary array accesses

- Poor thread affinity

## WHERE IS YOUR MEMORY ALLOCATED AND WHERE ARE YOUR THREADS RUNNING?

# Tuning Guides Available Online

- http://intel.com/vtune-tuning-guides

Intel® VTune™ Amplifier Tuning Guides

Our tuning guides explain how to identify common software performance issues using VTune Amplifier and give suggestions for optimization.

| Microarchitecture Code Name | Processors Covered | Tuning Guide |
|---|---|---|
| Apollo Lake | Intel Atom® Processor E3900 Series, and Intel® Pentium® and Celeron® Processor N- and J-Series | Download PDF |
| Skylake-X | Intel® Xeon Processor Scalable Family 1st Gen | Download New PDF Download Old PDF |
| Knights Landing | Intel® Xeon Phi™ Processor | Download PDF |
| Broadwell-E* (Server) | Intel® Xeon Processor E5 v4 Family | Download PDF |
| Skylake | 6th Generation Intel® Core™ Processor Family | Download PDF |
| Broadwell | 5th Generation Intel® Core™ Processor Family | Download PDF |
| Haswell-E* (Server) | Intel® Xeon® Processor E5 v3 Family | Download PDF |
| Ivy Bridge-E* (Server) | Intel® Xeon® Processor E5/E7 v2 Family | Download PDF |
| Haswell | 4th Generation Intel® Core™ Processor Family | Download PDF |
| Sandy Bridge-EP/EX/EN (Server) | Intel® Xeon® Processor E5 Family | Download PDF |
| Ivy Bridge | 3rd Generation Intel® Core™ Processor Families | Download PDF |
| Sandy Bridge | 2nd Generation Intel® Core™ Processor Families | Download PDF |
| Many Integrated Core Architecture | Intel® Xeon Phi™ coprocessor | Read the Article |

# Example 1 – Matrix Multiply



Bound by the L3 Cache while reading arrays

**General Exploration** General Exploration viewpoint (change)

Analysis Target | Analysis Type | Collection Log | Summary | Bottom-up

**Elapsed Time: 6.241s**

| | |
|---|---|
| Clockticks: | 71,522,000,000 |
| Instructions Retired: | 77,472,000,000 |
| CPI Rate: | 0.923 |
| MUX Reliability: | 0.984 |
| Front-End Bound: | 3.0% |
| Bad Speculation: | 0.1% |
| Back-End Bound: | 69.5% |
| Memory Bound: | 43.6% |
| L1 Bound: | 0.0% |
| L2 Bound: | 0.0% |
| L3 Bound: | 33.8% |
| DRAM Bound: | 5.5% |
| Store Bound: | 0.0% |
| Core Bound: | 25.9% |
| Retiring: | 27.4% |
| Total Thread Count: | 4 |
| Paused Time: | 0s |

| S. Li. | Source | Clo. | Ins... Re... | CPI Rate | Fr. B. | B. S. | L. B. | L2 Bou. | L3 Bound | D B. | S. B. | Cor. Bo. | R. | Sou... File |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 176 | int i,j,k; | | | | | | | | | | | | | |
| 177 | | | | | | | | | | | | | | |
| 178 | // Basic parallel implementation | | | | | | | | | | | | | |
| 179 | #pragma omp parallel for | | | | | | | | | | | | | |
| 180 | for(i=0; i<msize; i++) { | | | | | | | | | | | | | |
| 181 | for(j=0; j<msize; j++) { | 26,.. | 7,0.. | 3.714 | 0.1% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | mu... |
| 182 | for(k=0; k<msize; k++) { | 42,.. | 46,.. | 0.907 | 2.9% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | 15.0% | 16... | mu... |
| 183 | c[i][j] = c[i][j] + a[i][k] * b[k][j]; | 29,.. | 31,.. | 0.936 | 0.1% | 0.1% | 0.0% | 0.0% | 33.9% | 5.5% | 0.0% | 11.0% | 11... | mu... |
| 184 | } | | | | | | | | | | | | | |

# Example 1 – Matrix Multiply

Interchange loop indices and collapse loops

| S. Li. | Source | Cl. | Ins... Re... | CPI Rate | Fr. B. | B. S. |
|---|---|---|---|---|---|---|
| 203 | void multiply3(int msize, int tidx, int numt, TYPE a[][ | | | | | |
| 204 | { | | | | | |
| 205 | int i,j,k; | | | | | |
| 206 | | | | | | |
| 207 | #pragma omp parallel for collapse (2) | | | | | |
| 208 | for(i=0; i<msize; i++) { | | | | | |
| 209 | for(k=0; k<msize; k++) { | | | | | |
| 210 | #pragma ivdep | | | | | |
| 211 | for(j=0; j<msize; j++) { | 64,.. | 24.. | 0.267 | 0.0% | 0.0 |
| 212 | c[i][j] = c[i][j] + a[i][k] * b[k][j]; | 18,.. | 60,.. | 0.303 | 6.3% | 0.5 |
| 213 | } | | | | | |
| 214 | } | | | | | |

## General Exploration    General Exploration viewpoint (change) ?

◁ ⊕ Analysis Target    A Analysis Type    Collection Log    Summary    Bottom-up
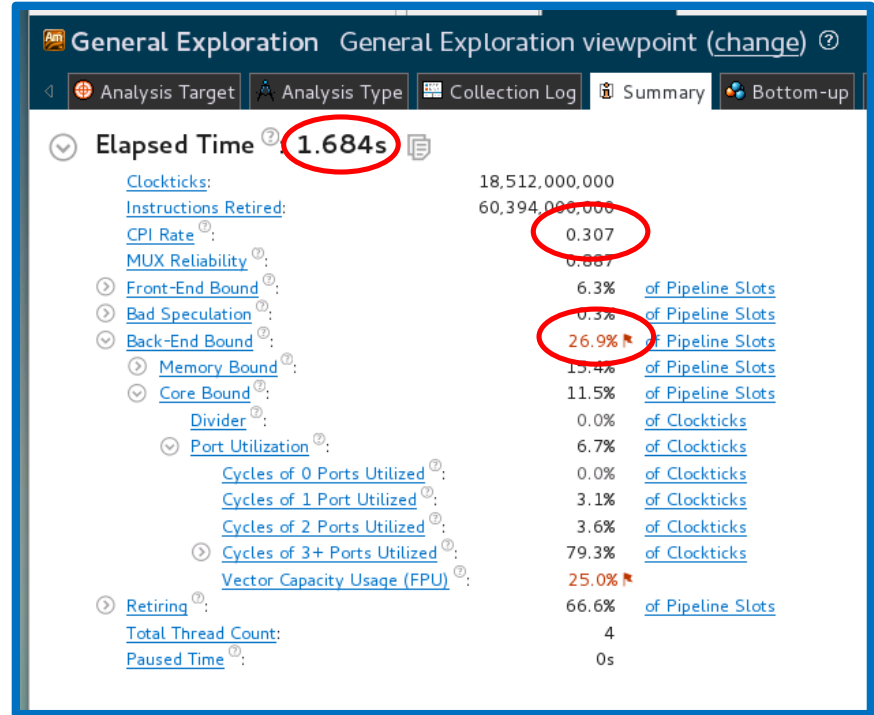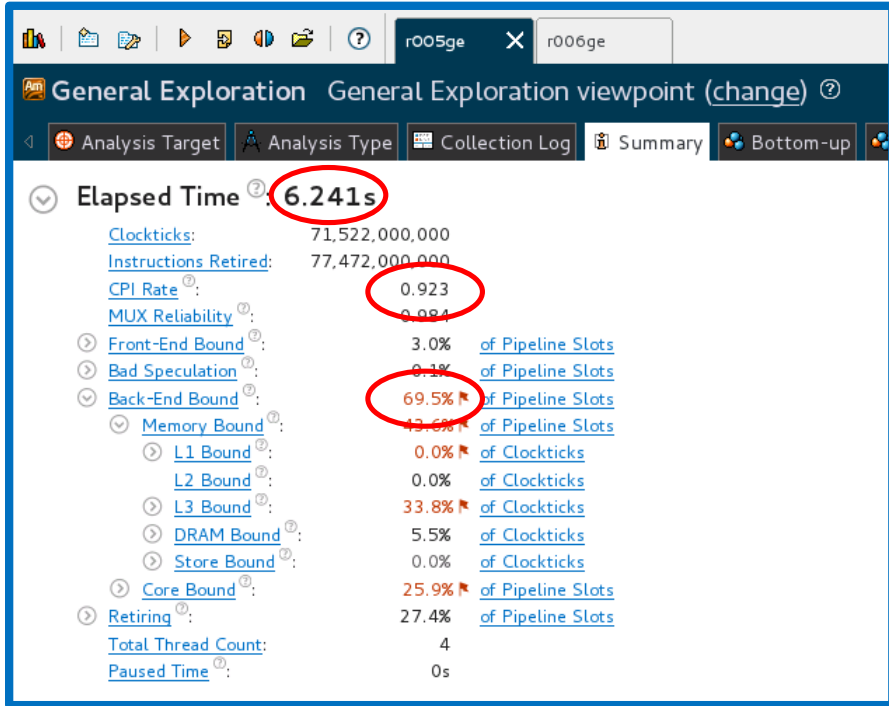
Elapsed Time ?: 1.684s

| | | |
|---|---|---|
| Clockticks: | 18,512,000,000 | |
| Instructions Retired: | 60,394,000,000 | |
| CPI Rate ?: | 0.307 | |
| MUX Reliability ?: | 0.887 | |
| Front-End Bound ?: | 6.3% | of Pipeline Slots |
| Bad Speculation ?: | 0.3% | of Pipeline Slots |
| Back-End Bound ?: | 26.9% ⚑ | of Pipeline Slots |
| Memory Bound ?: | 15.4% | of Pipeline Slots |
| Core Bound ?: | 11.5% | of Pipeline Slots |
| Divider ?: | 0.0% | of Clockticks |
| Port Utilization ?: | 6.7% | of Clockticks |
| Cycles of 0 Ports Utilized ?: | 0.0% | of Clockticks |
| Cycles of 1 Port Utilized ?: | 3.1% | of Clockticks |
| Cycles of 2 Ports Utilized ?: | 3.6% | of Clockticks |
| Cycles of 3+ Ports Utilized ?: | 79.3% | of Clockticks |
| Vector Capacity Usage (FPU) ?: | 25.0% ⚑ | |
| Retiring ?: | 66.6% | of Pipeline Slots |
| Total Thread Count: | 4 | |
| Paused Time ?: | 0s | |

(intel)

# Example 1 – Matrix Multiply

# Example 2 – Calculating Prime Numbers
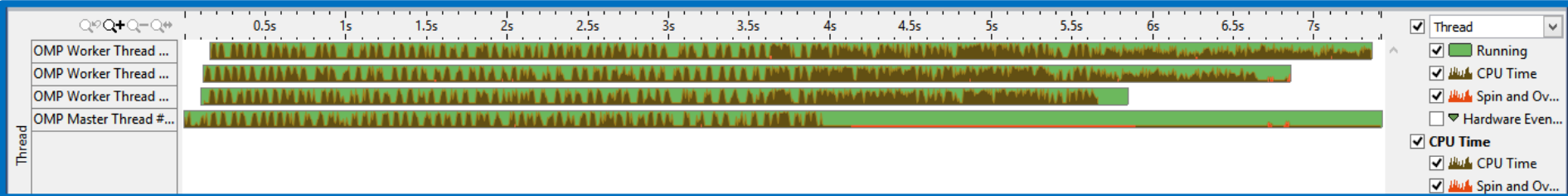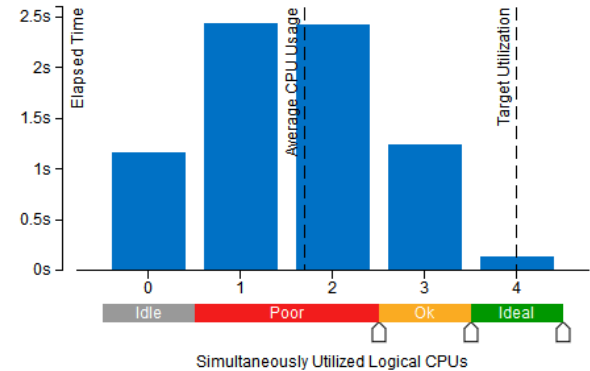
```
41  int _tmain(int argc, _TCHAR* argv[])
42  {
43      DWORD msBegin = timeGetTime();
44
45      #pragma omp parallel for
46      for(int p = 3; p <= limit; p += 2) {
47          if (IsPrime(p)) Tick();
48      }
49      DWORD msDuration = timeGetTime() - msBegin;
50
51      printf("MS: %d\n", msDuration);
52      printf("primes = %d\n", primes);
53      return primes != correctCount;
54  }
55
```

OpenMP uses Static Scheduling

Load Imbalance



**CPU Usage Histogram**

This histogram displays a percentage of the wall time the specific number of CPUs

Simultaneously Utilized Logical CPUs

# Example 2 – Calculating Prime Numbers
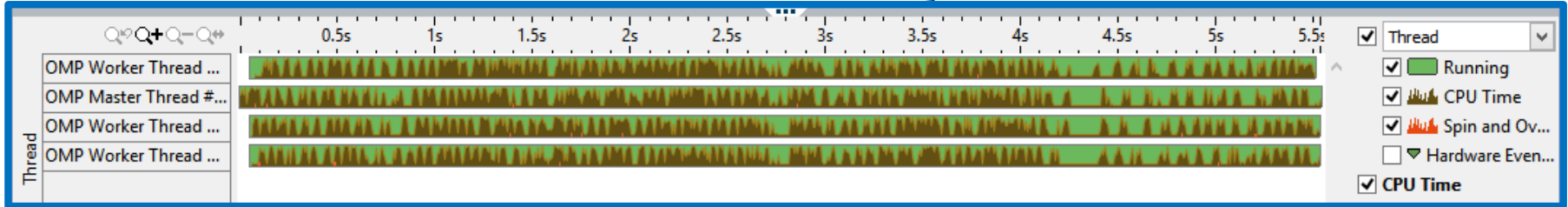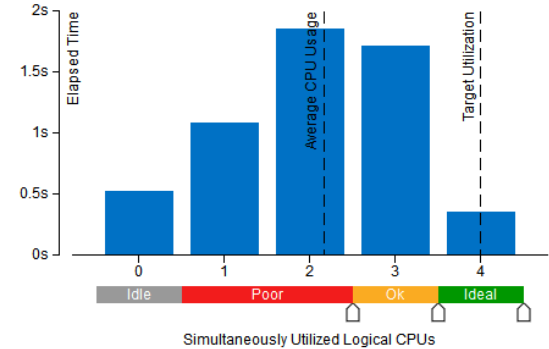
```
41  int _tmain(int argc, _TCHAR* argv[])
42  {
43      DWORD msBegin = timeGetTime();
44
45  #pragma omp parallel for schedule (dynamic, 1000)
46      for(int p = 3; p <= limit; p += 2) {
47          if (IsPrime(p)) Tick();
48      }
49      DWORD msDuration = timeGetTime() - msBegin;
50
51      printf("MS: %d\n", msDuration);
52      printf("primes = %d\n", primes);
53      return primes != correctCount;
54  }
55
56
```

**Switch to Dynamic Scheduling**

**More Balanced Threads**

**CPU Usage Histogram**

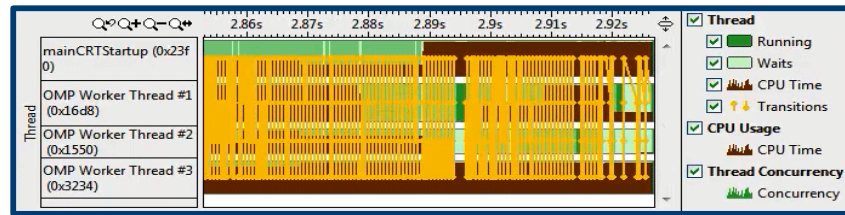This histogram displays a percentage of the wall time the specific number of CPUs



Simultaneously Utilized Logical CPUs

# Summary: Top Down Tuning Method

- **Make system-level optimizations**
- **Make algorithmic optimizations**
  - Use Threading Advisor to add threading
  - Use *Concurrency Analysis* and *Locks & Waits Analysis* to tune threading
- **Make microarchitectural optimizations**
  - **Find your hotspots**
    - Use *Hotspots Analysis* or *Advanced Hotspots Analysis*
  - **For each hotspot, determine efficiency.**
    - Use *General Exploration Analysis* to identify inefficient hotspots.
    - If inefficient: Categorize the bottleneck, identify the cause, and optimize it!
      - Hierarchical metrics in *General Exploration Analysis* focus your attention where it's needed most and allow you to easily identify the issue.
      - *Memory Access Analysis* can help with Back End Bound code.
      - Vectorization Advisor can help improve the efficiency of Retiring code.

# Intel® VTune™ Amplifier

Faster, Scalable Code Faster

## Get the Data You Need

- Hotspot (Statistical call tree), Call counts (Statistical)
- Thread Profiling – Concurrency and Lock & Waits Analysis
- Cache miss, Bandwidth analysis…[1]
- GPU Offload and OpenCL™ Kernel Tracing

## Find Answers Fast

- View Results on the Source / Assembly
- OpenMP Scalability Analysis, Graphical Frame Analysis
- Filter Out Extraneous Data – Organize Data with Viewpoints
- Visualize Thread & Task Activity on the Timeline

## Easy to Use

- No Special Compiles – C, C++, C#, Fortran, Java, ASM
- Visual Studio* Integration or Stand Alone
- Local & Remote Data Collection, Command Line
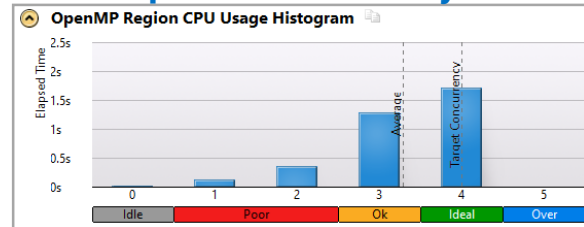- Analyze Windows* & Linux* data on OS X*[2]

### Quickly Find Tuning Opportunities

| Function / Call Stack | CPU Time | | | |
|---|---|---|---|---|
| | Effective Time by Utilization<br>Idle Poor Ok Ideal Over | Spin Time | Overhead Time | |
| ⊞ FireObject::checkCollision | 4.507s | 0s | 0s | |
| ⊞ FireObject::ProcessFireCollisionsRange | 3.444s | 0s | 0s | |
| ⊞ NtWaitForSingleObject | 0s | 3.406s | 0s | |
| ⊟ std::basic_ifstream<char,struct std::char_traits | 3.359s | 0s | 0s | |
| ⊞ ↰ Ogre::FileSystemArchive::open | 3.359s | 0s | 0s | |
| ⊞ CBaseDevice::Present | 2.335s | 0.671s | 0s | |
| Selected 1 row(s): | 1.151s | 0.728s | 0s | |

### See Results On The Source Code

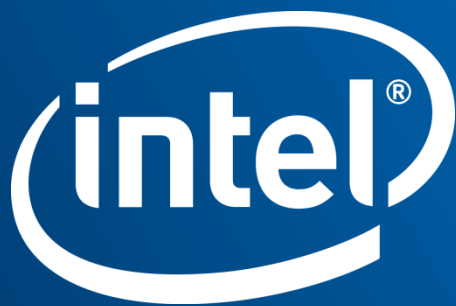| Source Line | Source | CPU Time: Total by Utilization<br>Idle Poor Ok Ideal Over |
|---|---|---|
| 81 | for (int i = 0; i < mem_array_i_max; i++) | 0.300s |
| 82 | { | |
| 83 | for (int j = 0; j < mem_array_j_max; j++) | 4.936s |
| 84 | { | |
| 85 | mem_array [j*mem_array_j_max+i] = *fill_val | 7.207s |

### Tune OpenMP Scalability

OpenMP Region CPU Usage Histogram

### Visualize & Filter Data

[1] Events vary by processor.   [2] No data collection on OS X*

# Legal Disclaimer & Optimization Notice

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products.  For more complete information visit www.intel.com/benchmarks.

INFORMATION IN THIS DOCUMENT IS PROVIDED "AS IS". NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO THIS INFORMATION INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.

**Optimization Notice**