![Intel Software logo](intel Software)

# OPENMP* ANALYSIS IN INTEL® VTUNE™ AMPLIFIER XE -
# TALKING TO A USER ABOUT OPENMP* PERFORMANCE  IN THE LANGUAGE THE
# PROGRAM WAS WRITTEN IN,
# WITH A LITTLE WANDER INTO VECTORIZATION TOO

Material from  Dmitry Prohorov, VTune HPC lead
Zakhar Matveev (Intel® Advisor Architect) and Alex Shinsel
Presented by Jim Cownie

# Agenda

VTune Amplifier XE OpenMP* Analysis: answering customers' questions about performance in the same language their program was written in

- Concepts, metrics and technology inside

- VTune Amplifier XE OpenMP Analysis Workflow

A short introduction to Intel® Advisor's Roofline Analysis

Summary

intel

# Typical customer questions on parallelization efficiency of OpenMP* applications

"I put pragmas but why is my speed up so poor?"

- Parallelization inefficiency

"I ran my app on a system with more cores but it doesn't run as efficiently as on a smaller one"

- Scalability issues

Decomposing the questions →

Is the serial time of my application significant in preventing scaling?
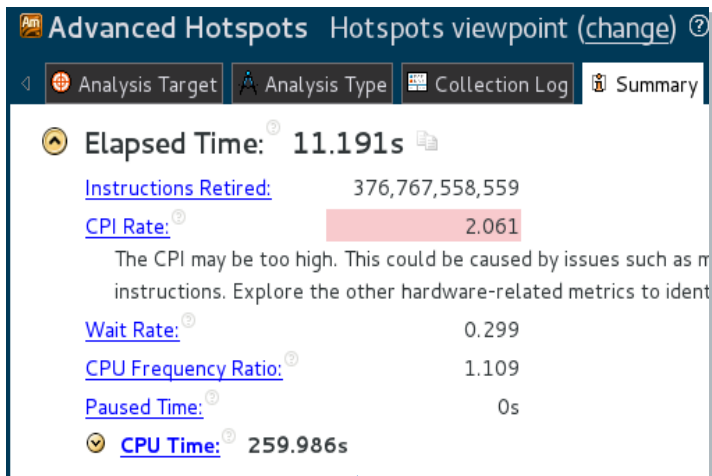
How efficient is my OpenMP parallelization?

- If inefficient, how much gain can be achieved if I invest in fighting the inefficiencies?

Which OpenMP regions/loops/barriers are worth tuning?

- What are their particular problems?

# If Performance Information is OpenMP "Unaware"..

The questions are tied to OpenMP program structure – #pragmas
Answers should be given the same way to be understandable and actionable



OpenMP "unaware" views of VTune Amplifier XE
Difficult to detect problems, customers blame the OpenMP
runtime seeing CPU time consumption there and not
understanding that this is a result of parallelization inefficiency

# Overview of summary pane

# Key to OpenMP awareness in VTune – Region based views and metrics

## Definition of Region **Potential Gain** (elapsed time metric)



Actual **Parallel Region Elapsed Time**

Fork

Join

Effective time (sampling)

Lock spinning (sampling)

Imbalance (tracing)

Scheduling (sampling)

Work forking (sampling)

Atomics (sampling)

**Estimated Ideal Time** =
Effective time / Number of Threads

**Potential Gain** as a sum of inefficiencies normalized by num of threads

(intel)

# Technology used by VTune Amplifier XE OpenMP Analysis

**Tracing** of OpenMP constructs to provide region/work sharing context and precise imbalance at barriers

- Provided to VTune by LLVM/Intel OpenMP Runtime under profiling

  - Fork-Join points of parallel regions with number of working threads (Intel Compiler 14 and later)

  - OpenMP construct barrier points with imbalance info and OpenMP loop metadata

  - `-parallel-source-info=2` Intel compiler option to embed source file name in region name

Looking at transition to OMPT, working with John M-C on interface enrichments for low overhead analysis

**Sampling** to define and classify CPU time – user's code and OpenMP RTL work

- Classification: Locking, Scheduling, Work Forking

(intel)

# VTune Amplifier XE OpenMP Analysis Workflow

Start with HPC Performance Characterization analysis

Explore CPU Utilization metrics related to OpenMP in summary, grid, and source views



CL: >amplxe-cl –collect hpc-performance <my_app>

# Per-region Details in grid view: inefficiencies in elapsed time are classified and highlighted

# Details in Grid View: Serial Time Hotspots



Serial hotspots under Master Thread

Time Filter to exclude initialization phase

# Details on Scalable Timeline

## Super tiny timeline display mode – a bird's eye view showing all data without scrolling

**Region frames on the ruler**



**More green => more efficient multithreaded execution**

Intel® Xeon Phi™ profiling result with 288 threads

# Details for a Region at source file level

# Intel® Vtune Amplifier Summary

VTune Amplifier XE OpenMP analysis answers customers' questions about performance in the language of OpenMP constructs

The analysis scales well for many-core systems with good balance of tracing and sampling collection technologies

The full feature set is available in VTune Amplifier XE with Intel OpenMP and Intel MPI runtimes as a part of Intel® Parallel Studio XE

(intel)

# ADD OPENMP SIMD WITH INTEL® (VECTOR) ADVISOR

Slides by Zakhar Matveev, Intel® Advisor Architect

# Intel Advisor: 5 tools for Efficient Vectorization and Memory utilization

**1. Compiler diagnostics + Performance Data + SIMD efficiency information**

**2. Guidance: detect problem and recommend how to fix it**

**3. "Precise" Trip Counts & FLOPs. Roofline analysis. Characterize your application.**

**4. Loop-Carried Dependency Analysis**

**5. Memory Access Pattern Analysis**

# Advisor Survey: vectorize and improve SIMD code performance!



- **Efficiency** – my performance thermometer

- **Recommendations** – get tips on how to improve performance, in particular using OpenMP 4.* and later!

# Advisor Dependencies: The Answer to Tough SIMD/Threading Question #1!

## Is it safe to force the compiler to vectorize?

```
DO I = 1, N
    A(I) = A(I-1) * B(I)
ENDDO
```

```
void scale(int *a, int *b)
{
    for (int i = 0; i < 1000; i++)
        b[i] = z * a[i];
}
```

**Issue: Assumed dependency present**

The compiler assumed there is an anti-dependency (Write after read – WAR) or true dependency (Read after write – RAW) in the loop. Improve performance by investigating the assumption and handling accordingly.

○ **Enable vectorization**

Potential performance gain: Information not available until Beta Update release
Confidence this recommendation applies to your code: Information not available until Beta Update release
The Correctness analysis shows there is no real dependency in the loop for the given workload. Tell the compiler it is safe to vectorize using the `restrict` keyword or a directive.

| ICL/ICC/ICPC Directive | IFORT Directive | Outcome |
|---|---|---|
| #pragma simd or #pragma omp simd | !DIR$ SIMD or !$OMP SIMD | Ignores all dependencies in the loop |
| #pragma ivdep | !DIR$ IVDEP | Ignores only vector dependencies (which is safest) |

**Read More:**

- User and Reference Guide for the Intel C++ Compiler 15.0 > Compiler Reference > Pragmas > Intel-specific Pragma Reference >
  ○ ivdep
  ○ omp simd

| 📋 Summary | 🕙 Survey & Roofline | 📊 Refinement Reports |
|---|---|---|

| Site Location | Loop-Carried Dependencies ▲ |
|---|---|
| ⊞ 📄 [loop in find_bg at lm.c:962] | 🟢 No Dependencies Found |
| ⊞ 📄 [loop in live_get_partialhyp at live.c:192] | 🟢 No Dependencies Found |
| ⊞ 📄 [loop in live_get_partialhyp at live.c:177] | 🟢 No Dependencies Found |
| ⊞ 📄 [loop in feat_s2mfc2feat_block at feat.c:1045] | 🟢 No Dependencies Found |
| ⊞ 📄 [loop in parse_tmat_senmap at mdef.c:389] | ⚠ Potential Reduction:1 |
| ⊞ 📄 [loop in subvq_init at subvq.c:271] | ⚠ Potential Reduction:1 |
| ⊞ 📄 [loop in utt_decode_block at utt.c:1125] | ⚠ Potential Reduction:1 |
| ⊞ 📄 [loop in utt_decode_block at utt.c:1024] | ⚠ Potential Reduction:2 |
| ⊞ 📄 [loop in vector_gautbl_eval_logs3 at vector.c:524] | ⚠ Potential Reduction:2 |
| ⊞ 📄 [loop in vithist_backtrace at vithist.c:775] | ❌ RAW:1 |
| 📄 [loop in glist_free at glist.c:230] | ❌ RAW:1 |
| ⊞ 📄 [loop in feat_s2mfc2feat_block at feat.c:1041] | ❌ RAW:1 |
| ⊞ 📄 [loop in vithist_lmstate_reset at vithist.c:462] | ❌ RAW:1 ⚠ Potential Reduction:1 |
| ⊞ 📄 [loop in utt_end at utt.c:219] | ❌ RAW:1 ⚠ Potential Reduction:1 |
| ⊞ 📄 [loop in utt_decode_block at utt.c:944] | ❌ RAW:1 ⚠ WAR:1 ⚠ Potential Reduction:1 |

> Safe to vectorize (at least for given workload), use **OMP SIMD**!

> Use **OMP reduction** (also for threading)!

> True dependence proven, not way to parallelize without extra work

(intel)

# Intel Advisor gives recommendations to guide you where and how to add OpenMP SIMD!

**What issues must to be fixed?**

**Which loops are vectorized?**



**Read more about the issue and what to modify in your code to fix it (to enable vector parallelism!)**

**In more detail: which function is causing the problem?**

**Add OpenMP SIMD: reference example**

# More examples: also threading aware..



**Issue: OpenMP function call(s) present** *Confidence level: low*

OpenMP* function call(s) in the loop body are preventing the compiler from effectively vectorizing the loop.

**Recommendation: Move OpenMP call(s) outside the loop body** *Confidence level: low*

OpenMP calls prevent automatic vectorization when the compiler cannot move the calls outside the loop body, such as when OpenMP calls are not invariant. To fix:

1. Split the OpenMP parallel loop
2. Move the OpenMP calls outside the loop when possible.

Example:

**Original code:**

```
#pragma omp parallel for private(tid, nthreads)
for (int k = 0; k < N; k++)
{
    tid = omp_get_thread_num(); // this call inside loop
    nthreads = omp_get_num_threads(); // this call inside
    ...
}
```

**Revised code:**

```
#pragma omp parallel private(tid, nthreads)
{
    // Move OpenMP calls here
    tid = omp_get_thread_num();
    nthreads = omp_get_num_threads();

    #pragma omp for nowait
    for (int k = 0; k < N; k++)
    {
        ...
    }
}
```

**Read More:**

- omp for, omp parallel recommendations
- Getting Started with Intel Compiler Pragmas and Dir

**Recommendation: Remove OpenMP lock functions** *Confidence level: low*

**Recommendation: Change the floating point model** *Confidence level: low*

Your application calls serialized versions of math functions when you use the `strict` floating point model. To fix: Do one of the following:

- Use the `fast` floating point model to enable more aggressive optimizations or the `precise` floating point model to disable optimizations that are not value-safe on fast transcendental functions.

| Windows* OS | Linux* OS |
|---|---|
| /fp:fast | -fp-model fast |
| /fp:precise /Qfast-transcendentals | -fp-model precise -fast-transcendentals |

CAUTION: This may reduce floating point accuracy.

- Use the `precise` floating point model and enforce vectorization of the source loop using a directive: `#pragma simd` or `#pragma omp simd`

Example: ⊙

```
gcc program.c -O2 -fopenmp -fp-model precise -fast-transcendentals

#pragma omp simd collapse(2)
for (i=0; i<N; i++)
{
    a[i] = b[i] * c[i];
    for (i=0; i<N; i++)
    {
        d[i] = e[i] * f[i];
    }
}
```

(intel)

# 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

# Roofline Metrics

Roofline is based on FLOPS and Arithmetic Intensity (AI).

- **FLOPS**: **Fl**oating-Point **Op**erations / **S**econd

- **Arithmetic Intensity**: FLOP / Byte Accessed

SpMV    FFTs    N-body

Low AI    High AI

Collecting this information in Intel® Advisor requires two analyses.

**Run Roofline**
▶ Collect

**1. Survey Target**
⏱ Collect

**1.1 Find Trip Counts and FLOP**
↻ Collect
☐ Trip Counts
☑ FLOP

Shortcut to run Survey followed by Trip Counts + FLOPs

Runs system benchmarks and collects timing data.

Collects memory traffic and FLOP data.
Must be run separately due to higher overhead that would interfere with timing measurements.

(intel)

# Classic vs. Cache-Aware Roofline

Intel® Advisor uses the Cache-Aware Roofline model, which has a different definition of Arithmetic Intensity than the original ("Classic") model.

## Classical Roofline

- Traffic measured from one level of memory (usually DRAM)
- AI may change with data set size
- AI changes as a result of memory optimizations

## Cache-Aware Roofline

- Traffic measured from all levels of memory
- AI is tied to the algorithm and will not change with data set size
- Optimization does not change AI*, only the performance

*Compiler optimizations may modify the algorithm, which may change the AI.*

(intel)

# Plotting a Roofline Chart

The maximum FLOPS as a product of ops/byte (AI) and maximum bytes supplied per second is a diagonal line.

**FLOPS**

The CPU's maximum FLOPS can be plotted as a horizontal line.

A loop or function can be plotted as a point on the graph.

**Arithmetic Intensity**
FLOP/Byte

A Roofline Chart uses AI as its X axis and FLOPS as its Y axis.

(intel)

# Ultimate Performance Limits



Performance cannot exceed the machine's capabilities, so each loop is ultimately limited by either compute or memory capacity.

FLOPS

Ultimately Memory-Bound

Ultimately Compute-Bound

Arithmetic Intensity
FLOP/Byte

# Sub-Roofs and Current Limits



FLOPS

L1 Cache

L2 Cache

L3 Cache

Vector with FMAs

Vector

Scalar

Additional roofs can be plotted for specific computation types or cache levels.

These sub-roofs can be used to help diagnose bottlenecks.

Arithmetic Intensity
FLOP/Byte

(intel)

# The Intel® Advisor Roofline Interface

Roofs are based on benchmarks run before the application.

- Roofs can be hidden, highlighted, or adjusted.

Intel® Advisor has size- and colour-coding for dots.

- Colour code by duration or vectorization status

- Categories, cutoffs, and visual style can be modified.

# Identifying Good Optimization Candidates

Focus optimization effort where it makes the most difference.

- Large, red loops have the most impact.

- Loops far from the upper roofs have more room to improve.

# Identifying Potential Bottlenecks

Final roofs *do* apply; sub-roofs *may* apply.

- Roofs above indicate *potential* bottlenecks

- Closer roofs are the most likely suspects

- Roofs below may contribute but are generally not primary bottlenecks

# Feature Synergy
## Overcoming the Scalar Add Peak

Survey and Code Analytics tabs indicate vectorization status with colored icons.

↻ = Scalar    ↻ = Vectorized

"Why No Vectorization" tab and column in Survey explain what prevented vectorization.

Recommendations tab may help you vectorize the loop.

Dependencies determines if it's safe to force vectorization.

# Feature Synergy
## Overcoming the Vector Add Peak

Survey and Code Analytics display the vector efficiency and presence of FMAs.

- Recommendations may help improve efficiency or induce FMA usage.



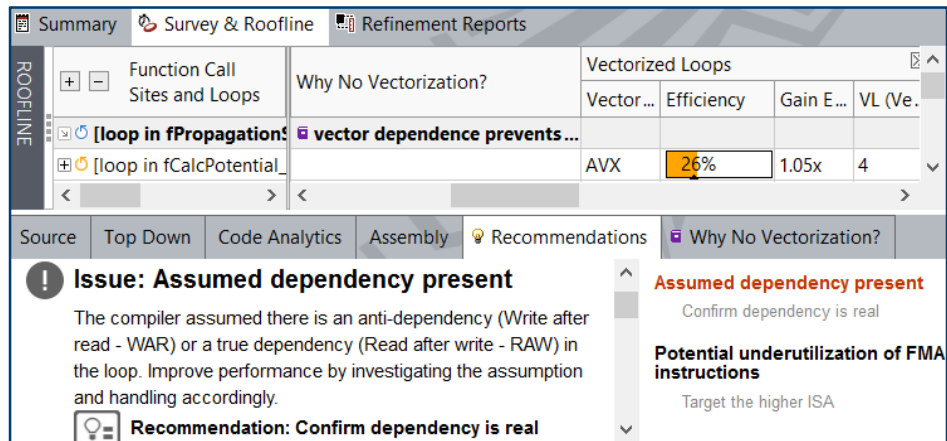| Address | Line | Assembly |
|---|---|---|
| 0x140001550 | | Block 1: 1660000000 |
| 0x140001550 | 262 | vmovupd ymm3, ymmword ptr [rsi+rcx*8+0x26400] |
| 0x140001559 | 262 | vmovdqa ymm1, ymm0 |
| 0x14000155d | 262 | vfmadd132pd ymm1, ymm3, ymmword ptr [rsi+rcx*8+0x23a80] |
| 0x140001567 | 262 | vaddpd ymm2, ymm1, ymm3 |
| 0x14000156b | 262 | vmovupd ymm1, ymmword ptr [rsi+rcx*8+0x26420] |
| 0x140001574 | 262 | vaddpd ymm4, ymm2, ymm3 |
| 0x140001578 | 262 | vmovupd ymm5, ymm0 |
| 0x14000157c | 262 | vfmadd132pd ymm5, ymm1, ymmword ptr [rsi+rcx*8+0x23aa0] |
| 0x140001586 | 262 | vmovupd ymmword ptr [rsi+rcx*8+0x21100], ymm4 |
| 0x14000158f | 262 | vaddpd ymm5, ymm5, ymm1 |
| 0x140001593 | 262 | vaddpd ymm2, ymm5, ymm1 |
| 0x140001597 | 260 | add rcx, 0x8 |
| 0x14000159b | 262 | vmovupd ymmword ptr [rsi+rdx*8+0x21100], ymm2 |
| 0x1400015a4 | 260 | add rdx, 0x8 |
| 0x1400015a8 | 260 | cmp rcx, 0x530 |
| 0x1400015af | 260 | jb 0x140001550 <Block 1> |

The Assembly tab* is useful for determining how well you are making use of FMAs.
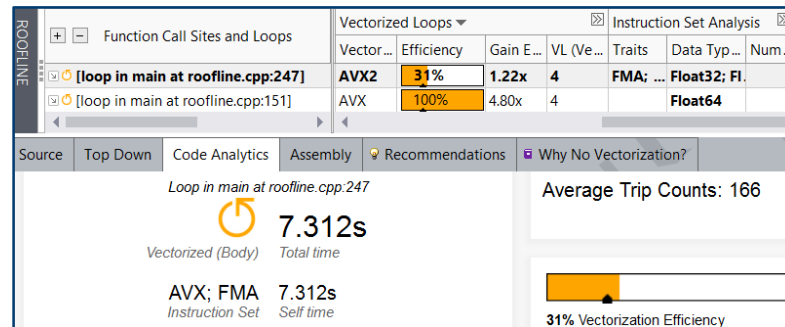
*Color coding added for clarity.*

| Address | Line | Assembly |
|---|---|---|
| 0x1400015f0 | | Block 1: 1660000000 |
| 0x1400015f0 | 275 | vmovupd ymm1, ymmword ptr [rsi+rcx*8+0x26400] |
| 0x1400015f9 | 275 | vmovupd ymm2, ymmword ptr [rsi+rcx*8+0x26420] |
| 0x140001602 | 275 | vfmadd231pd ymm1, ymm1, ymm0 |
| 0x140001607 | 275 | vfmadd231pd ymm2, ymm2, ymm0 |
| 0x14000160c | 275 | vfmadd231pd ymm1, ymm0, ymmword ptr [rsi+rcx*8+0x23a80] |
| 0x140001616 | 275 | vfmadd231pd ymm2, ymm0, ymmword ptr [rsi+rcx*8+0x23aa0] |
| 0x140001620 | 275 | vmovupd ymmword ptr [rsi+rcx*8+0x21100], ymm1 |
| 0x140001629 | 275 | vmovupd ymmword ptr [rsi+rdx*8+0x21100], ymm2 |
| 0x140001632 | 273 | add rcx, 0x8 |
| 0x140001636 | 273 | add rdx, 0x8 |
| 0x14000163a | 273 | cmp rcx, 0x530 |
| 0x140001641 | 273 | jb 0x1400015f0 <Block 1> |

# Feature Synergy
## Overcoming the Memory Bandwidth Roofs
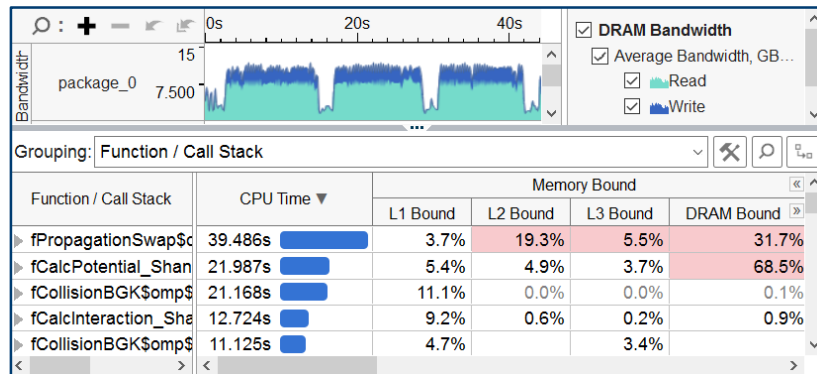


Memory Access Patterns (MAP) identifies inefficient access patterns.

Intel® SIMD Data Layout Templates (Intel® SDLT) allows code written as AOS to be stored as efficient SOA.

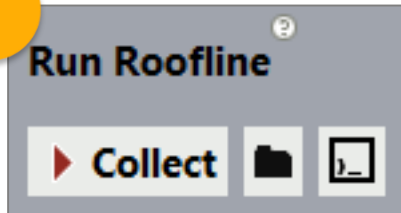Intel® VTune™ Amplifier can be used to further optimize cache usage.

If cache usage cannot be improved, try re-working the algorithm to increase the AI (and slide up the roof)



| Site Location | Strides Distribution | Access Pattern | Max. Site Footprint | Recommendations |
|---|---|---|---|---|
| [loop in main at roofline.cpp:1 … | 0% / 100% / 0% | All const strides | 38KB | 1 Inefficient me… |
| [loop in main at roofline.cpp:1 … | 50% / 50% / 0% | Mixed strides | 10KB | 1 Inefficient me… |
| [loop in main at roofline.cpp:1 … | 100% / 0% / 0% | All unit strides | 9KB | |

Memory Access Patterns Report | Dependencies Report | Recommendations

| ID | | Stride | Type | Source | Variable references | Max. Site Footprint | Access Type |
|---|---|---|---|---|---|---|---|
| P1 | | 8 | Constant stride | roofline.cpp:127 | AoS1_Y | 38KB | Read |
| P2 | | 2 | Constant stride | roofline.cpp:127 | AoS1_X | 10KB | Write |

DRAM Bandwidth
☑ Average Bandwidth, GB…
☑ Read
☑ Write

Grouping: Function / Call Stack

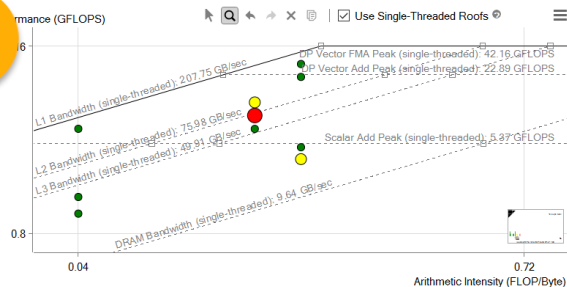| Function / Call Stack | CPU Time ▼ | Memory Bound | | | |
|---|---|---|---|---|---|
| | | L1 Bound | L2 Bound | L3 Bound | DRAM Bound |
| fPropagationSwap$c | 39.486s | 3.7% | 19.3% | 5.5% | 31.7% |
| fCalcPotential_Shan | 21.987s | 5.4% | 4.9% | 3.7% | 68.5% |
| fCollisionBGK$omp$ | 21.168s | 11.1% | 0.0% | 0.0% | 0.1% |
| fCalcInteraction_Sha | 12.724s | 9.2% | 0.6% | 0.2% | 0.9% |
| fCollisionBGK$omp$ | 11.125s | 4.7% | | 3.4% | |

(intel)

# Intel® Advisor Roofline Summary



Intel® Advisor's Roofline Chart is highly customizable and easy to generate.

Lets you identify the best optimization candidates by focusing on low, large loops.

Use the chart to identify the most likely bottlenecks.

Intel® Advisor's many other features allow deep analysis of suspected problems and provide advice on how to overcome them.

# Overall Conclusions:

Tuning requires analysis of what's really happening in the code

- Our guesses are frequently wrong

We also need to understand the potential for improvement

Intel tools (Vtune Amplifier, Advisor and Inspector) can help

Download Intel Parallel Studio XE (includes all the tools and other goodies) and try it.

- Free for students, teachers, Open Source contributors.

- Free evaluation licenses for everyone else.

- Or you can even pay money!

intel

# Legal Disclaimer & Optimization Notice