**ATPESC 21**

# Intel® VTune Profiler and Intel® Advisor Overview

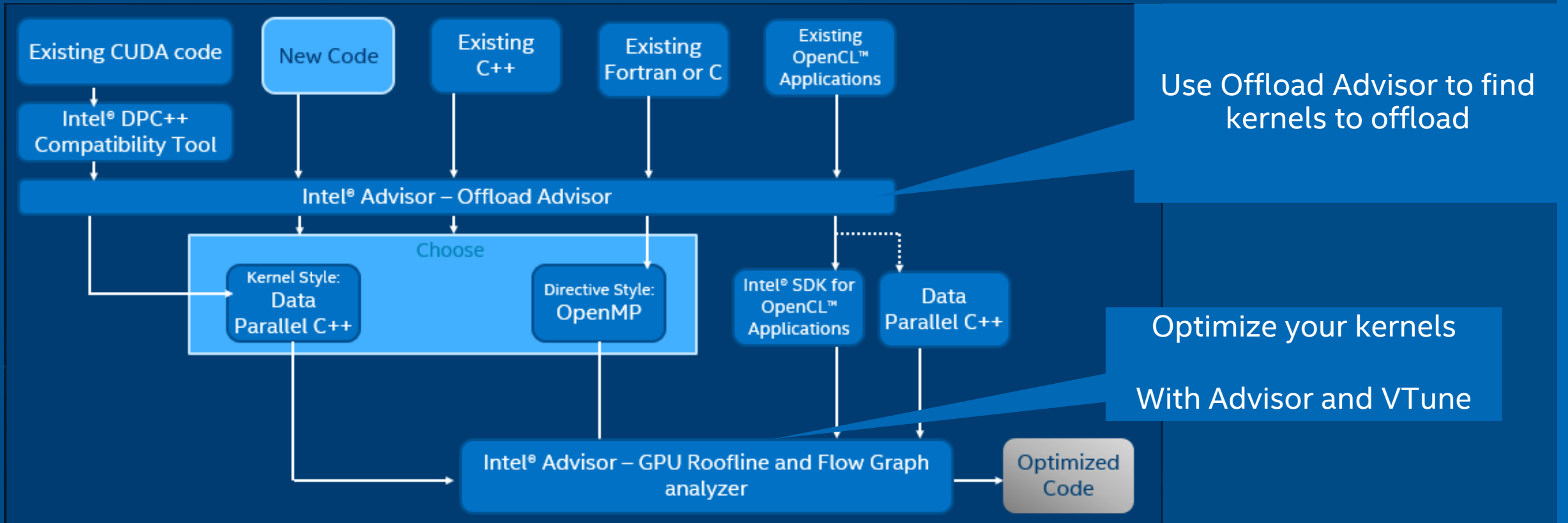Kevin O'Leary

intel.

# Agenda

| | |
|---|---|
| **1** | Intel® VTune Profiler overview |
| **2** | Intel® VTune Profiler for GPU |
| **3** | Offload Modeling with Intel® Advisor |
| **4** | GPU Roofline with Intel® Advisor |

# Using Intel® Analyzers to increase performance

# Intel Analysis Tools for GPU Compute

## Intel® Advisor

- **Offload Advisor**
  - Identify high-impact opportunities to offload
  - Detect bottlenecks and key bounding factors
  - Get your code ready even before you have the hardware by modeling performance, headroom and bottlenecks

- **Roofline Analysis**
  - See performance headroom against hardware limitations
  - Determine performance optimization strategy by identifying bottlenecks and which optimizations will payoff the most
  - Visualize optimization progress

- **Flow Graph Analyzer**
  - Visualize your CPU/GPU code and get recommendations for the CPU device

## Intel® VTune™ Profiler

- **Offload Performance Tuning**
  - Explore code execution on various CPU and GPU cores on your platform
  - Correlate CPU and GPU activity
  - Identify whether your application is GPU or CPU bound

- **HPC Performance Characterization**
  - Identify whether the OpenMP application offloads work to GPU effectively

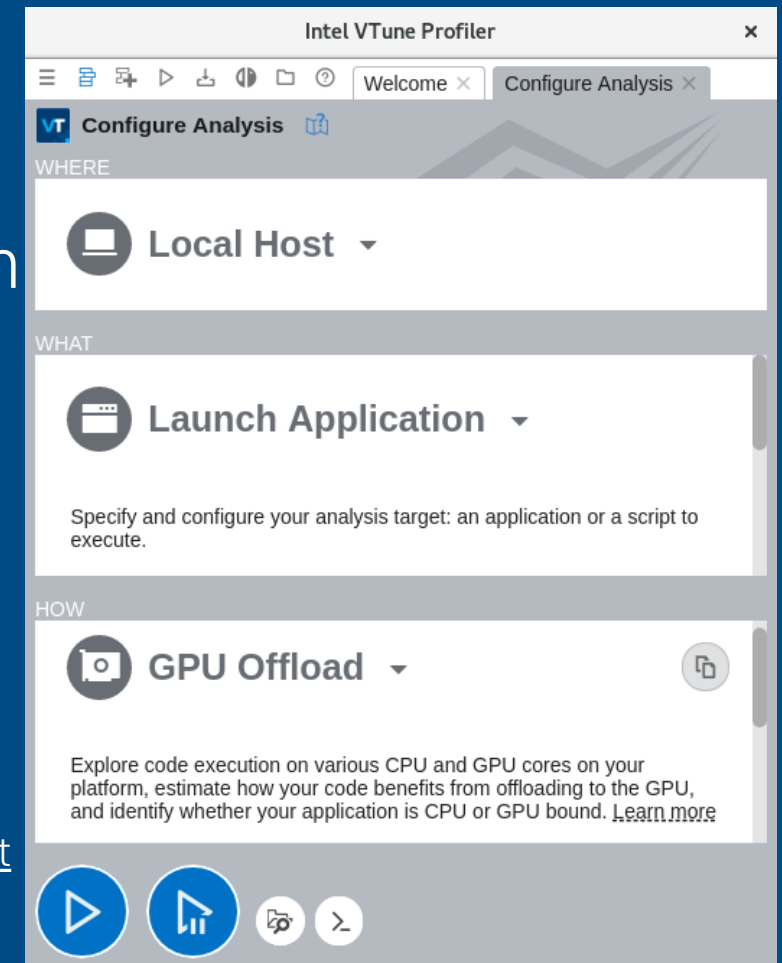- **GPU Compute/Media Hotspots**
  - Analyze the most time-consuming GPU kernels, characterize GPU usage based on GPU hardware metrics
  - GPU code performance at the source-line level and kernel assembly level

# Intel® VTune™ Profiler

# Intel® VTune™ Profiler GUI: quick overview

- GUI provides 3 panes to configure the analysis:

  - **WHERE** is used to specify an analysis system

  - **WHAT** is used to specify an analysis target

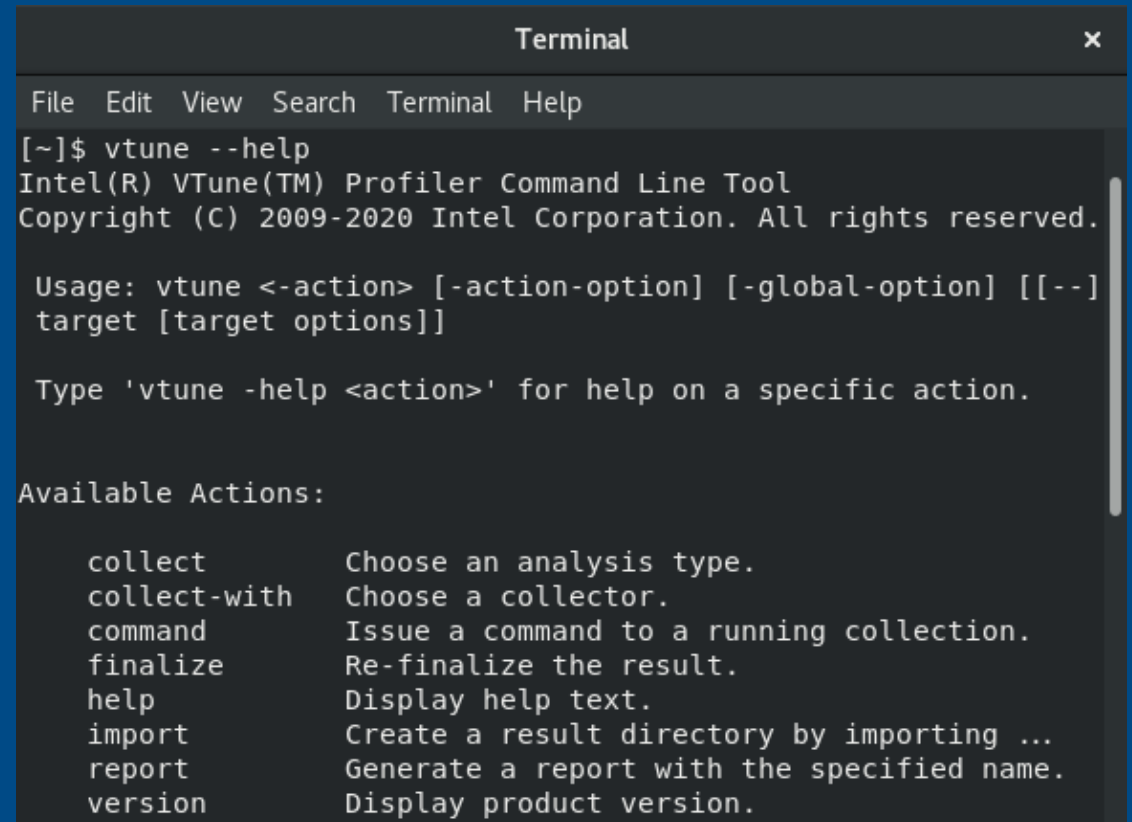  - **HOW** is used to select an analysis type

VTune Profiler documentation: <u>WHERE: Analysis system</u>, <u>WHAT: Analysis Target</u> and <u>HOW: Analysis Types</u>

intel.

# VTune CLI: quick overview

- CLI has its own help with several levels:
  - vtune –help
  - vtune –help collect
  - vtune –help collect gpu-offload
- Run collection:
  - vtune –collect *<analysis_type>* *<target>*
- Generate a report:
  - vtune –report *<report_name>* –r *<result_dir>*

  VTune Profiler documentation: Command Line Interface

```
Terminal                                                    ✕

File  Edit  View  Search  Terminal  Help

[~]$ vtune --help
Intel(R) VTune(TM) Profiler Command Line Tool
Copyright (C) 2009-2020 Intel Corporation. All rights reserved.

Usage: vtune <-action> [-action-option] [-global-option] [[--]
target [target options]]

Type 'vtune -help <action>' for help on a specific action.


Available Actions:

    collect        Choose an analysis type.
    collect-with   Choose a collector.
    command        Issue a command to a running collection.
    finalize       Re-finalize the result.
    help           Display help text.
    import         Create a result directory by importing ...
    report         Generate a report with the specified name.
    version        Display product version.
```

# GPU offload

- Helps to identify whether the application offloads work to GPU effectively.

- Can be used to profile OpenCL, Level0 and Intel Media SDK based applications or DPC++ and OpenMP applications that offload work on Intel GPU.

VTune Profiler documentation: GPU Offload Analysis

# Optimize your GPU usage using Intel® VTune Profiler
## GPU offload



This analysis enables you to:

Identify how effectively your application uses DPC++ or OpenCL kernels.

Explore GPU usage and analyze a software queue for GPU engines at each moment of time

# Optimize your GPU usage using Intel® VTune™ Profiler
## GPU offload



Use the GPU offload features Intel® VTune™ Profiler to see how effectively we are using our GPU.

VTune Profiler shows a synchronized time line between the CPU and GPU. GPU offload does indicate that our GPU execution units are stalling as indicated by the dark read bar in our timeline.

# HPC Performance Characterization

- Helps to identify whether the OpenMP application offloads work to GPU effectively.

- The main difference with GPU Offload analysis is that the data is collected through OMPT interface.

VTune Profiler documentation: HPC Performance Characterization Analysis

# GPU Compute/Media Hotspots

- Allows to analyze the most time-consuming GPU kernels, characterize GPU usage based on GPU hardware metrics, identify performance issues caused by memory latency or inefficient kernel algorithms, and analyze GPU instruction frequency per certain instruction types.

VTune Profiler documentation: GPU Compute/Media Hotspots Analysis



intel

# Optimize your GPU usage using Intel® VTune Profiler
## GPU hotspots

Active vs. Idle EU activity

Bandwidth at multiple levels

Run VTune Profiler GPU Hotspots to try to identify the source of our low GPU utilization and stalls. Click on the graphics tab in GPU Hotspots and you can see a high-level diagram of your architecture.

# Intel® VTune™ Profiler for Intel GPUs – Timelines for Correlation



Identify too much or too little kernel activity

Correlate GPU activity with kernels and threads

# Intel® VTune™ Profiler for DPC++ Code



Performance metrics at the DPC++ statement level

GPU assembly available for compute kernels

# Intel® Advisor

# Rich Set of Capabilities for High Performance Code Design
Intel® Advisor

**Offload Advisor**

Design offload strategy and model performance on GPU.

**Roofline Analysis**

Optimize your application for memory and compute.

**Vectorization Optimization**

Enable more vector parallelism and improve its efficiency.

**Thread Prototyping**

Model, tune, and test multiple threading designs.

**Build Heterogeneous Algorithms**

Create and analyze data flow and dependency computation graphs.

# Offload Modeling
# With Intel® Advisor

# Intel® Advisor – Offload Modeling

- Run on CPU or GPU – Predict for GPU

- Helps to define which sections of the code should run on given accelerator

- Provides performance projection on accelerators

# Intel® Advisor – Offload Modeling
## Find code that can be profitably offloaded



**Top Metrics**

| 3.3x | 3.3x | 100% | 2 |
|------|------|------|---|
| Speed Up for Accelerated Code | Amdahl's Law Speed Up | Fraction of Accelerated Code | Number of Offloads |

Loop takes 100% of the whole app execution time

Loop on GPU is 3.3x faster than on CPU

The whole app is 3.3x faster

**Program Metrics**

| Original | 23.30s |
| Accelerated | 7.00s |

| Program Time on Host... | 0s | Target Platform | Gen9 GT2 |
| Non Accelerated Time | 0s | Number of Offloads | 2 |
| Time in MPI calls | 0s | Speed Up for Accelerated Code | 3.3x |
| Time on Target | 7.00s | Amdahl's Law Speed Up | 3.3x |
| | | Fraction of Accelerated Code | 100% |

# In-Depth Analysis of Top Offload Regions

Provides a detailed description of modeling for each loop

- Timings(total time, time on the accelerator, speedup)

- Offload metrics (offload tax data transfers)

Loop at multiply.c:53 is recommended for offloading
- LLC BW bound
- Estimated to run on GPU in 6.997s
- Transfers 101MB of data

## Offload Modeling

Summary > Accelerated Regions > Logs

| | 3.3x | % | 2 |
|---|---|---|---|
| | Speed Up for Accelerated | Fraction of Accelerated Code | Number of Offloads |

- Memory traffic (DRAM, L3, L2, L1), trip count

CPU+GPU Highlight which part of the code should be run on the accelerator

| Loop/Function | Measured » | Basic Estimated Metrics » | | | Estimated Bounded By » | | | Estimated Data Transfer With Reuse » |
|---|---|---|---|---|---|---|---|---|
| | Time | Speed-Up | Time | Offload Summary | Throughput | Taxes With Reuse | Latencies | |
| ▼ [loop in multiply1 at multiply.c:53] | 23.27s | 3.326x | 6.997s | 🔴 Offloaded | LLC BW 6.997s<br>L3 BW 3.215s | Launch Tax < 0.1ms<br>All Taxes < 0.1ms | Load < 0.1ms | Read 101MB<br>Write 0B |
| ▼ [loop in multiply1 at multiply.c:54] | 23.27s | | | | | | | |
| [loop in multiply1 at multiply.c:55] | 23.26s | | | | | | | |
| ▼ [loop in multiply1 at multiply.c:54] | 28.5ms | 19.002x | 1.5ms | 🔴 Offloaded | LLC BW 1.5ms<br>L3 BW 0.7ms | Launch Tax < 0.1ms<br>All Taxes < 0.1ms | Load < 0.1ms | Read 67.9kB<br>Write 0B |
| [loop in multiply1 at multiply.c:55] | 19.7ms | | | | | | | |

# In-Depth Analysis of Top Offload Regions

Loop metrics are matched with Source and Call Tree



| Loop/Function | Measured » Time | Basic Estimated Metrics » | | | Estimated Bounded By » | | | Estimated Data Transfer With Reuse » |
|---|---|---|---|---|---|---|---|---|
| | | Speed-Up | Time | Offload Summary | Throughput | Taxes With Reuse | Latencies | |
| ▼ Total | 23.28s | | | | | | | |
| ▼ func@0x4b2e8759 | 23.27s | | | | | | | |
| ▼ func@0x4b2e8775 | 23.27s | | | | | | | |
| ▼ BaseThreadInitThunk | 23.27s | | | | | | | |
| ▼ ThreadFunction | 23.27s | | | | | | | |
| ▼ multiply1 | 23.27s | | | | | | | |
| ▶ [loop in multiply1 at multiply.c:53] | 23.27s | 3.326x | 6.... | 🔲 Offloaded | LLC... 6.... L3 ... 3.2... | Launch Tax < 0.1ms All Taxes < 0.1ms | L... < 0.... | Read 101MB Write 0B |
| ▶ _scrt_common_main_seh | 98.5ms | | | | | | | |

# GPU Roofline
# With Intel® Advisor

- 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?

# What is a Roofline Chart?

## Compare application performance against hardware limitations



Performance (GFLOPS) — Use Single-Threaded Roofs

42.16 — DP Vector FMA Peak (single-threaded): 42.16 GFLOPS
DP Vector Add Peak (single-threaded): 22.89 GFLOPS

L1 Bandwidth (single-threaded): 207.75 GB/sec
L2 Bandwidth (single-threaded): 75.98 GB/sec
L3 Bandwidth (single-threaded): 49.91 GB/sec

Scalar Add Peak (single-threaded): 5.37 GFLOPS

DRAM Bandwidth (single-threaded): 9.64 GB/sec

0.8

0.04                    0.72
Arithmetic Intensity (FLOP/Byte)

# 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

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

# Find Effective Optimization Strategies

Intel® Advisor– GPU Roofline

- GPU Roofline Performance

- Highlights poor performing loops

- Shows likely causes of bottlenecks

- Suggests next optimization steps

- Shows performance 'headroom' for each loop

- Which can be improved

- Which are worth improving

# Find Effective Optimization Strategies

Intel® Advisor– GPU Roofline



Configure levels to display

Shows performance headroom for each loop

Likely bottlenecks

Suggests optimization next steps

# Customize to Display Only Desired Roofs



Click on the top-right corner and remove unused roofs

GPU Roofline Insights

**Summary metrics**

**Create a snapshot**

53.0% FPU Utilization
99.4% EU Threading Occupancy
1.639 EU IPC Rate

Summary > GPU Roofline Regions · Logs ▾

**Switch between report tabs**

**Customizable GPU Roofline chart**

**GPU performance of compute tasks**

GPU Roofline

INT; CARM; GTI (Memory) ▾ | Compare ▾ | Guidance ▾

Int16 Vector Add Peak: 436.78 GINTOPS

L3 Bandwidth: 202.55 GB/sec
SLM Bandwidth: 202.45 GB/sec
GTI Bandwidth: 76.8 GB/sec
DRAM Bandwidth: 35.05 GB/sec

INTOP/Byte (Arithmetic Intensity)

Self Elapsed Time: 0.150 s

**Point Info**
Matrix1<float>
Self Performance: 45.095 GINTOPS
Self GTI (Memory) Arithmetic Intensity
Bounded by: L3 Bandwidth
Self Elapsed Time: 0.150 s
Self Memory Traffic: 0.835 GB

Copy To Clipboard

**Memory Metrics**
Impacts
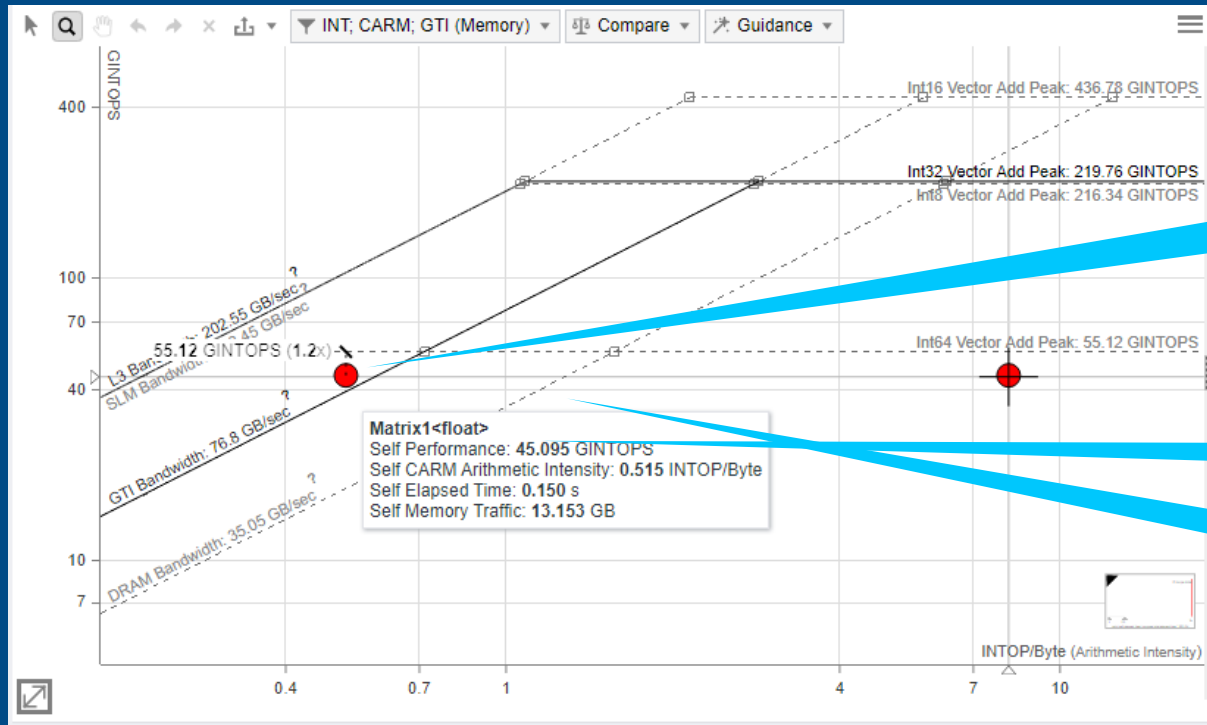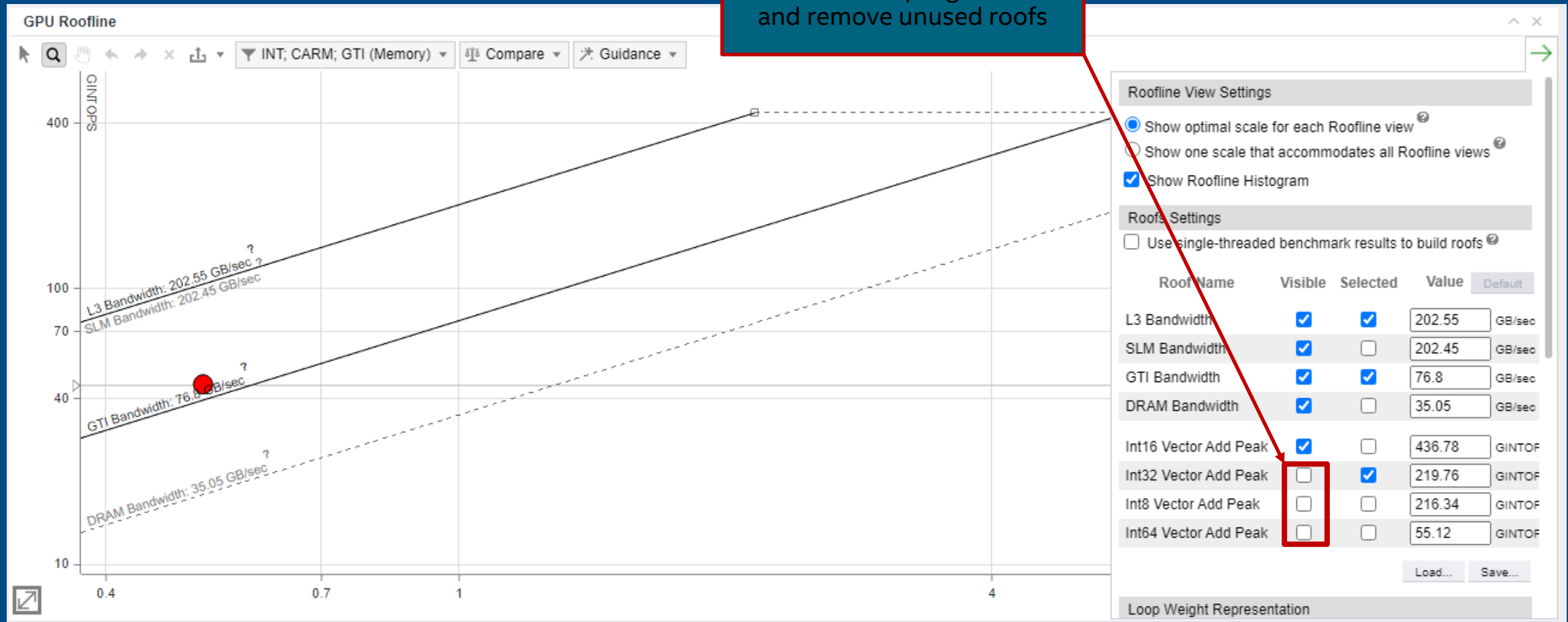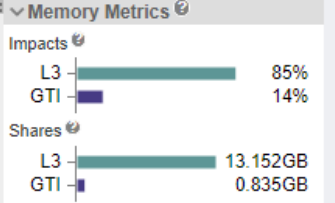L3 — 85%
GTI — 14%
Shares
L3 — 13.152GB
GTI — 0.835GB

Details × | GPU Source × | G ∧ ×

**View Details, GPU Source, and GPU Assembly info**

Matrix1<float>

SUMMARY
Elapsed Time 0.15s | GINTOPS
| GFLOPS
Work Size 1024 x 1024 | Local 256 x 1

ROOFLINE
Bounded by Int32 Vector Add Peak

Int32 Vector Add Peak

104.31

(↑2.3x)

45.1

L3 13.152 GB | GTI (Memory) 0.835 GB

GPU

| Compute Task | Elapsed Time | GPU Compute Performance | | | | | | Work Size | | Compute Task Purpose | Compute Task |
| | | GFLOPS | GINTOPS | FP AI | INT AI | GFLOP | GINTOP | Global | Local | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| [Outside any task] | 3.213s | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | | | [Unknown] | 0s |
| zeCommandListAppendMemoryCopyRegion | 0.002s | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | | | Transfer Out | 0.002s |
| zeCommandListAppendBarrier | 0.000s | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | | | Synchroniz... | 0.000s |
| Matrix1<float> | 0.150s | 14.298 | 45.095 | 2.573 | 8.114 | 2.147 | 6.773 | 1024 x 1024 | 256 x 1 | Compute | 0.150s |

# Summary

- You can use the Advisor and VTune GUI & CLI to run the collection and to generate the reports.

- Advisor and VTune both provide several analysis types to profile GPU workload.

# 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**

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.

Notice revision #20110804

intel