



Intel[®] oneAPI

IXPUG WEBINAR 2/20/20
ANALYZER TRAINING

AGENDA

- Introduction to Intel® oneAPI
- Intel® Advisor
 - Intel® Advisor – Offload Advisor
 - Intel® Advisor – GPU Roofline Analysis
- Intel® VTune™ Profiler
 - GPU Analysis

PROGRAMMING CHALLENGES FOR MULTIPLE ARCHITECTURES

Growth in specialized workloads

No common programming language or APIs

Inconsistent tool support across platforms

Each platform requires unique software investment

Diverse set of data-centric hardware required

Application Workloads Need Diverse Hardware



SCALAR



VECTOR



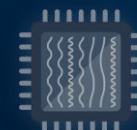
MATRIX



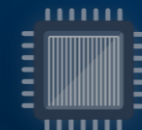
SPATIAL

Middleware / Frameworks

Language & Libraries

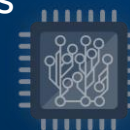


CPU

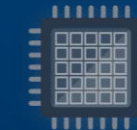


GPU

XPU^s

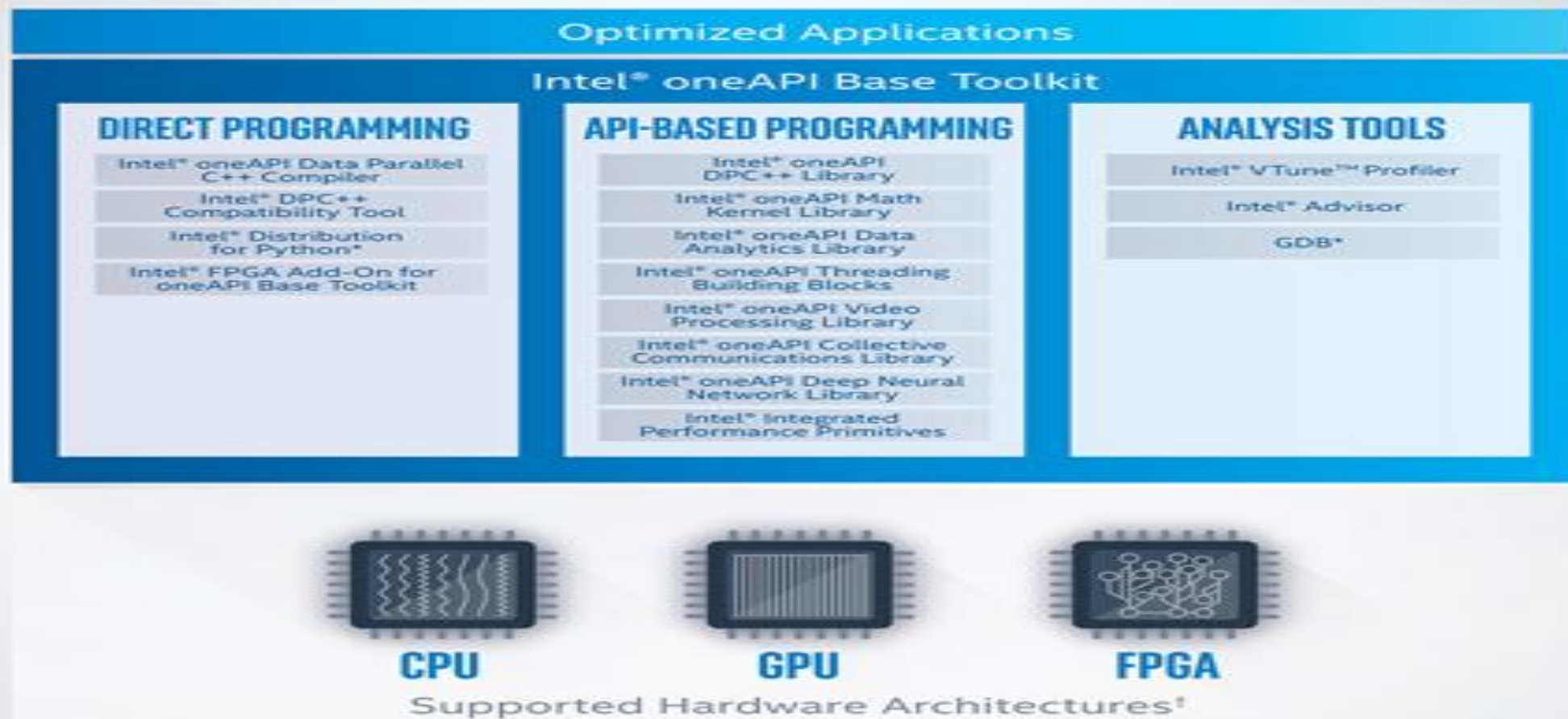


FPGA



OTHER ACCEL.

INTEL® ONEAPI



†Hardware support varies by individual oneAPI tool. Additional architecture support will be expanded over time.
*Other names and brands may be claimed as the property of others.

Intel VTune™ Profiler vs Intel® Advisor

GPU ANALYSIS

Intel® Advisor

- Offload Advisor in design phase
 - Design code for efficient offloading to accelerators—even before you have hardware. Estimate code performance and compare it with data transfer costs. No recompilation required.
- GPU Roofline
 - Optimize your CPU or GPU code for compute and memory. Locate bottlenecks and identify performance headroom for each loop or kernel to prioritize which optimizations will deliver the highest performance payoff.

• Intel® VTune™ Profiler

- GPU Offload to determine bottleneck
 - Explore code execution on various CPU and GPU cores on your platform, correlate CPU and GPU activity, and identify whether your application is GPU or CPU bound.
- 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

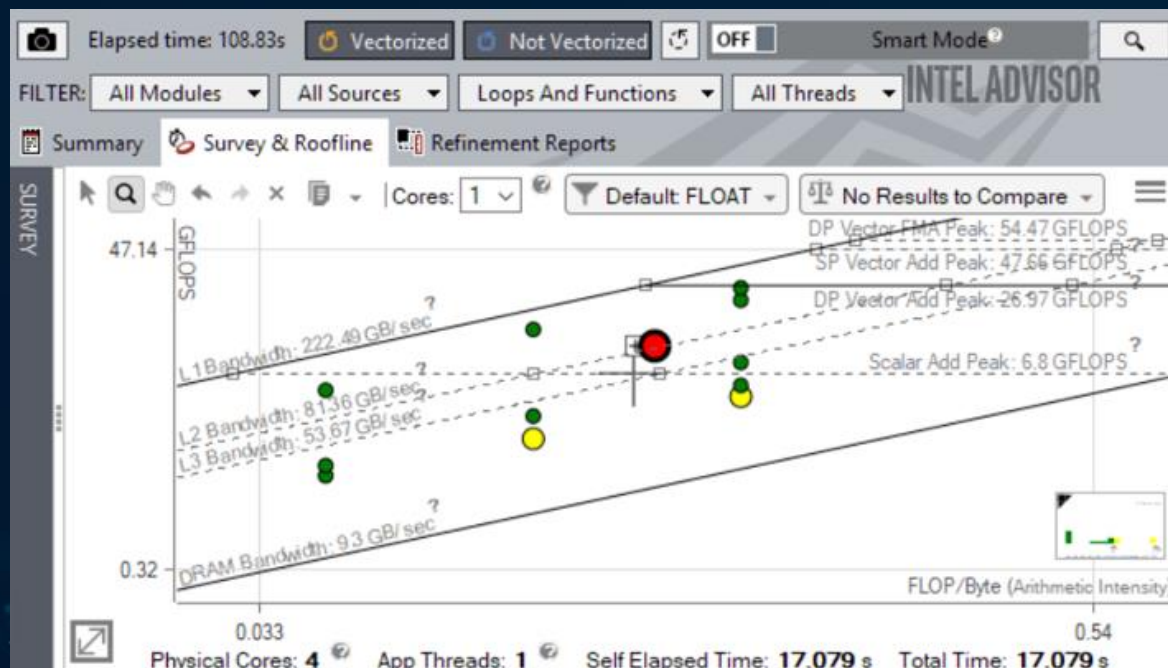


INTEL® ADVISOR

Vectorization Optimization, Thread Prototyping & Flow Graph Analysis

Design Modern Software with Intel® Advisor

Optimize Vectorization, Prototype Threading, Create & Analyze Flow Graphs



Capabilities

- Adds & optimizes vectorization
- Analyzes memory patterns
- Quickly prototypes threading
- Shows performance headroom for each loop
- Creates and analyzes flow graphs

New for 2020 Release (partial list)

- AVX-512 VNNI support
- Roofline Guidance
 - Loop specific actionable optimization recommendations
- Roofline Compare
 - Visualize multiple Roofline results in one chart to track optimization progress
- Integrated Roofline (Technical Preview)
 - Pinpoint the cache level causing bottlenecks

Learn More: software.intel.com/advisor

INTEL® ADVISOR

OFFLOAD ADVISOR

[Optimization Notice](#)

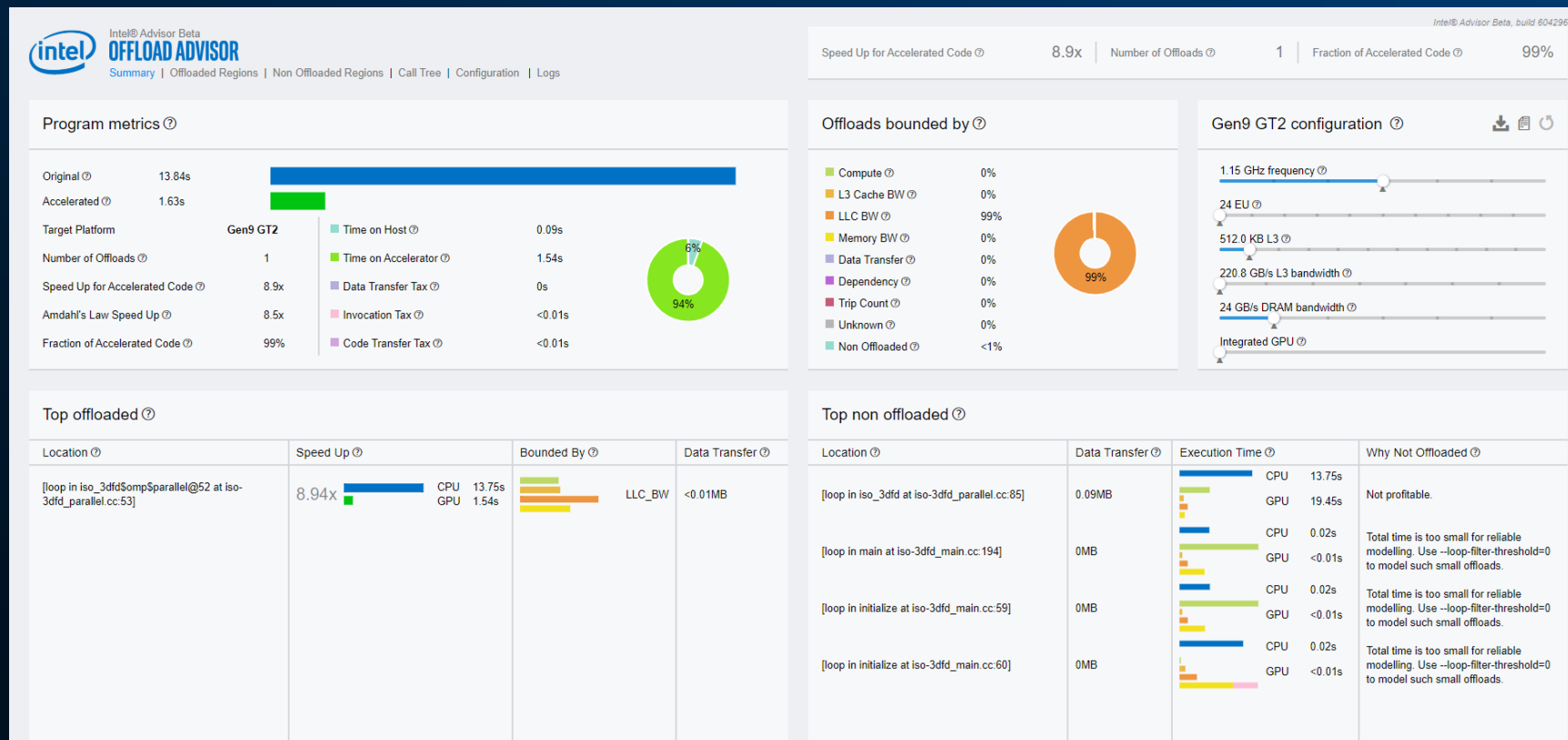
Copyright © 2019, Intel Corporation. All rights reserved.

*Other names and brands may be claimed as the property of others.



INTEL® ADVISOR - OFFLOAD ADVISOR (BETA)

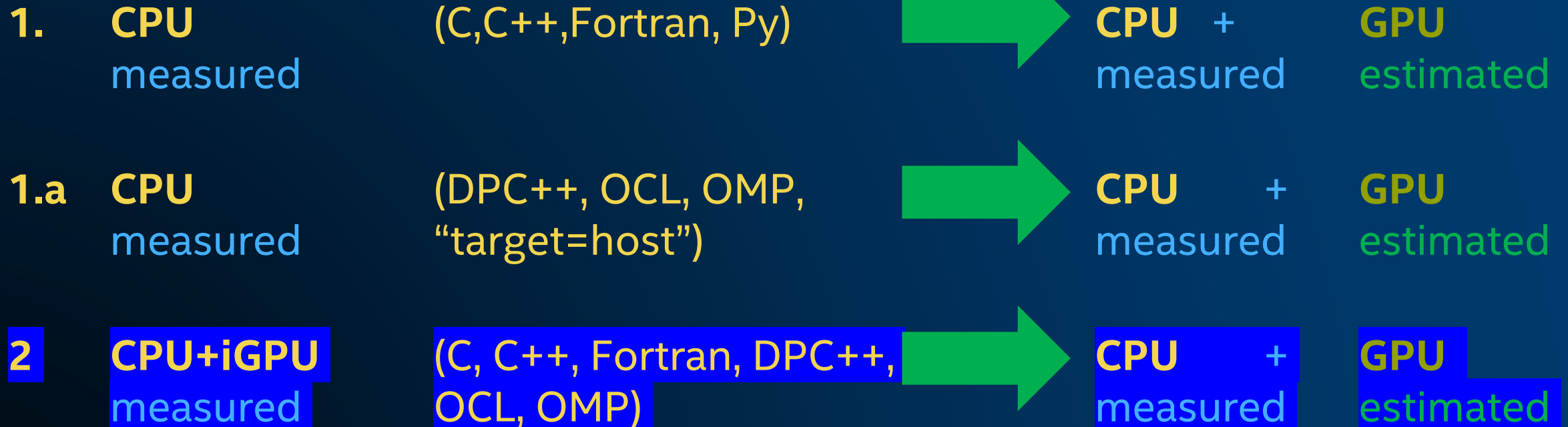
- Starting from a optimized binary (running on CPU):
 - Helps define which sections of the code should run on a given accelerator
 - Provides performance projection on accelerators (currently gen9 and gen11)



MODELING FLOWS SUPPORTED

Baseline HW (Programming model)

Target HW



Coming soon!

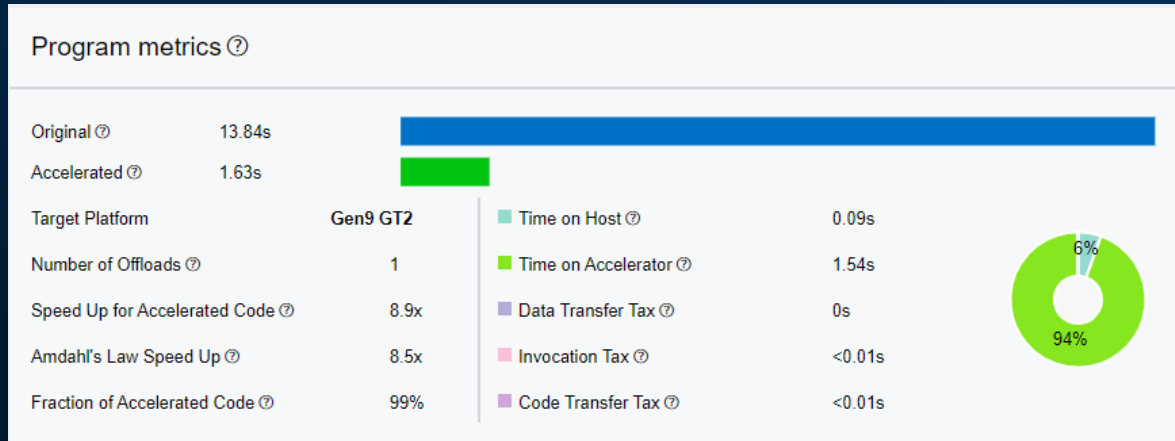
[Optimization Notice](#)

Copyright © 2019, Intel Corporation. All rights reserved.

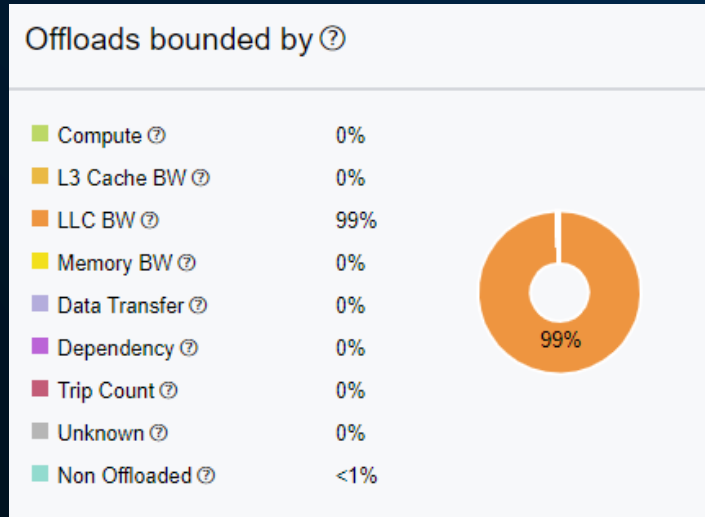
*Other names and brands may be claimed as the property of others.

WILL OFFLOAD INCREASE PERFORMANCE?

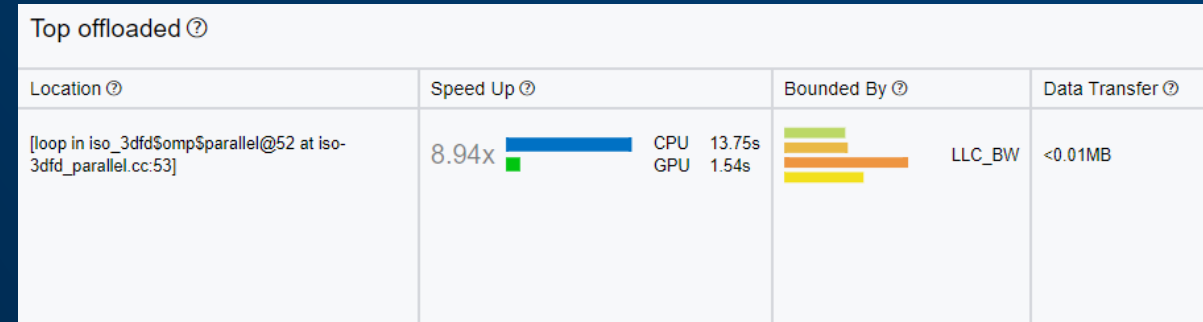
- How your code might perform on an accelerator ?



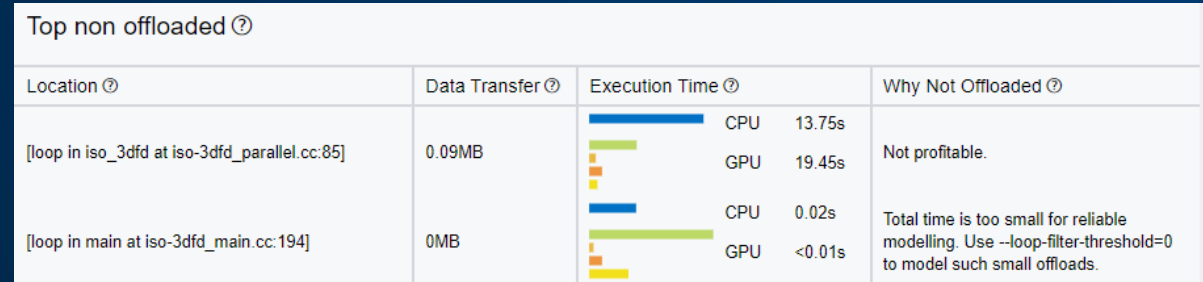
- What might be limiting your performance on the accelerator ?



- What should you offload ?



- What are the bad candidates for offload and Why ?



IN DEPTH ANALYSIS OF TOP OFFLOAD REGIONS

- Provides a detailed description of each loop interesting for offload
 - Timings (total time, time on the accelerator, speedup)
 - Offload metrics (offload tax, data transfers)
 - Memory traffic (DRAM, L3, L2, L1), trip count
 - Highlight which part of the code should run on the accelerator

This is where you will use DPCPP or OMP target for offload


The screenshot displays the Intel Advisor Beta Offload Advisor interface. The top navigation bar includes links for Summary, Offloaded Regions, Non Offloaded Regions, Call Tree, Configuration, and Logs. The main content area features a table with columns for Hierarchy, Trip Counts, L3 Cache, LLC, Memory, Instruction & Traffic Counts, and Diagnostics. The table lists three offload regions, with the first two highlighted in blue. To the right of the table, a code snippet is shown, highlighting a loop that is the subject of the offload analysis. A red arrow points from the text box above to this code snippet.

Hierarchy	Total Data Transferred from GPU to CPU (MB)	Average Trip Count	Call Count	Total L3 Traffic (GB)	Total LLC Access (GB)	Total Memory Traffic (GB)	FPU Util (GFLOP/s)	FLOP per Cycle	Diagnostics
[loop in iso_3dfd\$omp\$parallel@52 at i	<0.01	57600	102	174.250	113.259	23.637	7.896	7.896	In whole loop
[loop in iso_3dfd\$omp\$parallel@52 at	0	30	5875200	173.894	113.257	23.637	7.947	7.947	
[loop in iso_3dfd\$omp\$parallel@52	<1	<1	<1	0	0	0	0	0	Aggregated e

```
51 #pragma omp parallel for OMP_SCHEDULE num_threads(1) c
52 for(int iz=HALF_LENGTH; iz<n3-HALF_LENGTH; iz++) {
53     for(int iy=HALF_LENGTH; iy<n2-HALF_LENGTH; iy++) {
54         #pragma omp simd
55             for(int ix=HALF_LENGTH; ix<n1-HALF_LENGTH; ix+
56                 int offset = iz*dimn1n2 + iy*n1 + ix;
57                 float value = 0.0;
58                 value += ptr_prev[offset]*coeff[0];
59                 for(int ir=1; ir<=HALF_LENGTH; ir++) {
60                     value += coeff[ir] * (ptr_prev[offset
61                     value += coeff[ir] * (ptr_prev[offset
62                     value += coeff[ir] * (ptr_prev[offset
```

WHAT SHOULD NOT BE OFFLOADED

- Explains why Intel® Advisor doesn't recommend a given loop for offload
 - Dependency issues
 - Not profitable
 - Total time is too small



Intel® Advisor Beta

OFFLOAD ADVISOR

[Summary](#) |
 [Offloaded Regions](#) |
 [Non Offloaded Regions](#) |
 [Call Tree](#) |
 [Configuration](#) |
 [Logs](#)

Speed Up for Accelerated Code ⓘ

8.9x

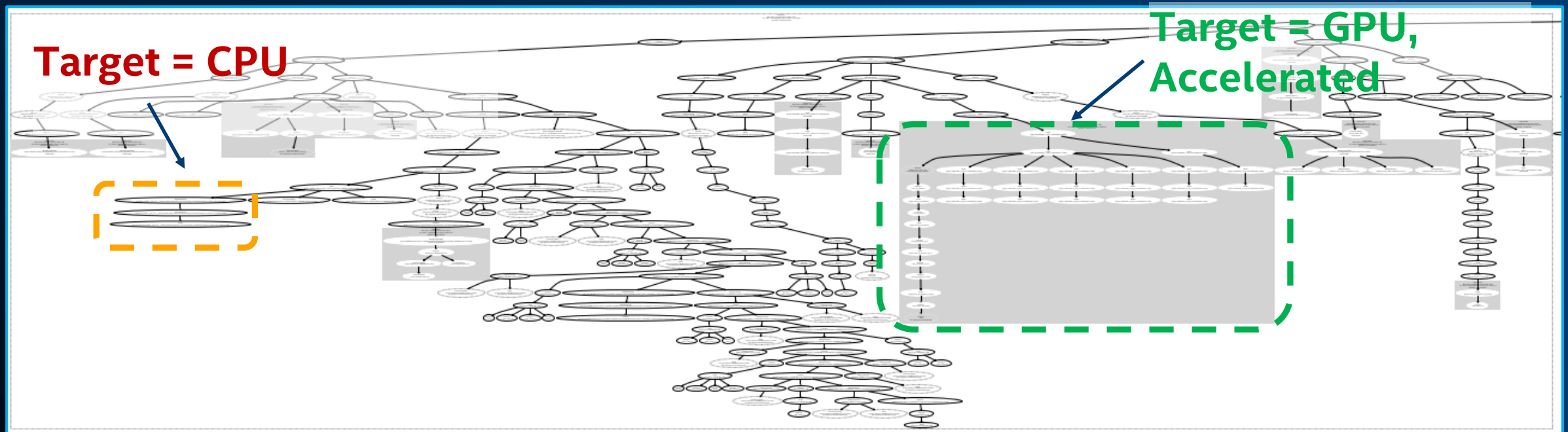
Number of O

Hierarchy	Information >				Potential Offload <		Column configurator
	Loop ion ed?	Estimated Execution Time on Accelerator (+Host) (s)	Bounded By	Fraction Offloaded (%)	Why Not Offloaded	Potential Speed Up for Whole Region	
[loop in iso_3dfd at iso-3dfd_parallel.cc]		1.539		100.00	Not profitable.	0.7068x	
[loop in main at iso-3dfd_main.cc:194]		0.020		0	Total time is too small for reliable modelling. Use --loop-filter-threshold=0 to model such small offloads.	41338.1565x	
> [loop in initialize at iso-3dfd_main.cc:59]					Total time is too small for reliable modelling. Use --loop-filter-threshold=0 to model such small offloads.	640.4016x	Custom file

ANOTHER VIEW OF YOUR CODE

PROGRAM TREE

The program tree offers another view of the proportion of code that can be offloaded to the accelerator.



BEFORE YOU START

INTEL® ADVISOR OFFLOAD ADVISOR

- The only strict requirement for compilation and linking is full debug information:
 - g:** Requests full debug information (compiler and linker)
- Offload Advisor supports any optimization level, but the following settings are considered the optimal requirements:
 - O2:** Requests moderate optimization
 - no-ipo:** Disables inter-procedural optimizations that may inhibit Offload Advisor to collect performance data (Intel® C++ & Fortran Compiler specific)

SETTING UP YOUR ENVIRONMENT

- To set up the Intel® Advisor Beta environment, run one of the shell script:

```
source <ONEAPI_INSTALL_DIR>/setvars.sh
```

or

```
source <ADV_INSTALL_DIR>/env/vars.sh
```

- This script sets all required Intel Advisor environment variables, including APM, which points to **<ADV_INSTALL_DIR>/perfmodes**
- This is the location of the Offload Advisor scripts in the Intel® Advisor Beta installation directory



The performance modeling functionality
is available on Linux* OS only

HOW TO RUN

- Easy to collect data and generate output with batch mode:

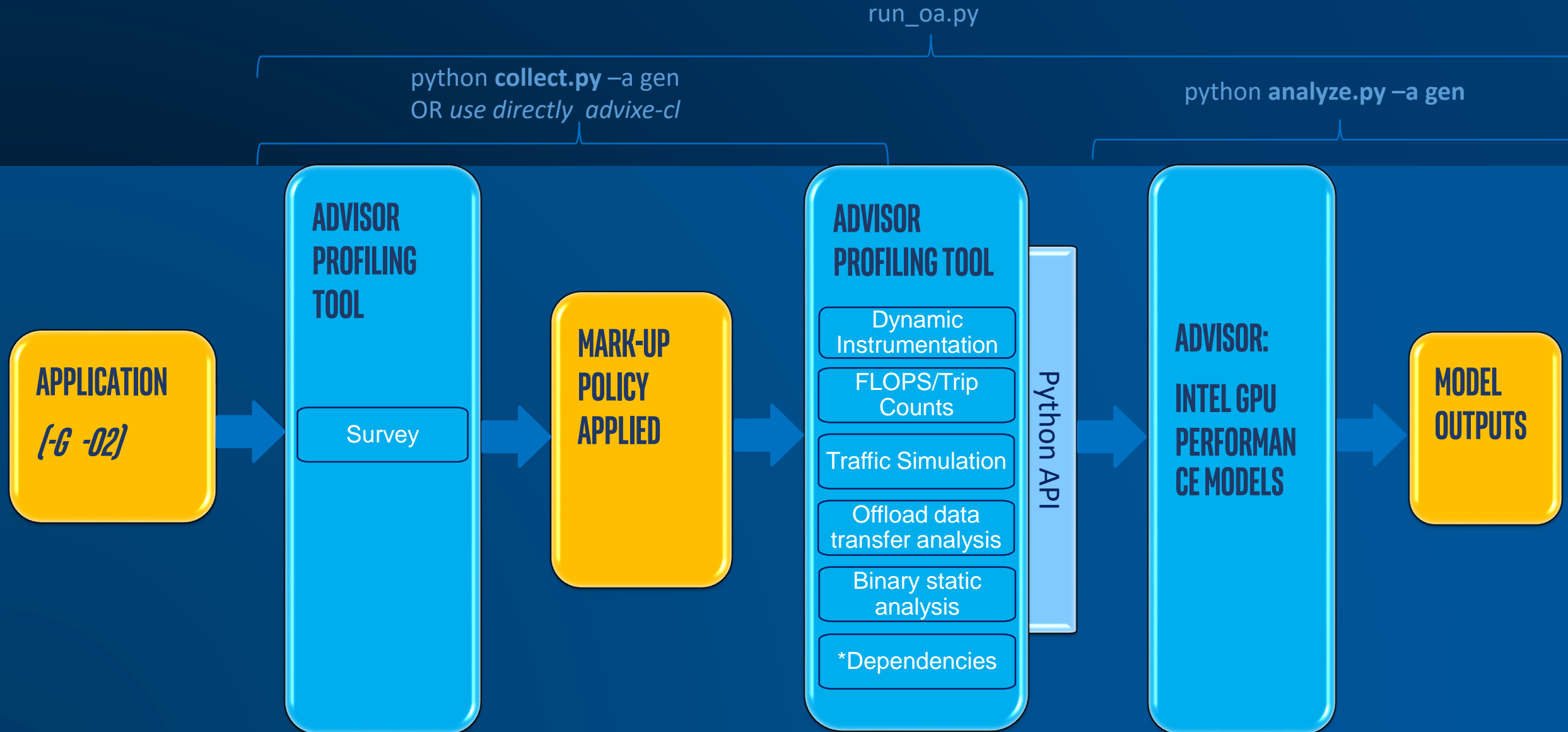
```
advixe-python <ADV_INSTALL_DIR>/perfmodels/run_oa.py <path_to_result_dir>  
-config gen9 --out-dir <path_to_result_dir> [--options] -- <app>
```

- By default, **run_oa.py** marks up all regions and only selects the most profitable ones for analysis
- To generate the report.html, uses the following command:

```
advixe-python $APM/analyse.py <project_dir> --config gen9 [--options] --  
<app_binary> [app_options]
```

```
u31313@s001-n004:/opt/intel/inteloneapi/advisor/latest/perfmodels$ ls  
accelerators    analyze.py      collect.py      debug.so        environ.py     oa_wrapper.so  shared.so      toml  
analyze_impl.so collect_impl.so compute_stats.py dot_graph.so    helpers        run_oa.py      template       tree.so
```

HOW TO RUN



[Optimization Notice](#)

Copyright © 2019, Intel Corporation. All rights reserved.

*Other names and brands may be claimed as the property of others.



INTEL ADVISOR® - OFFLOAD ADVISOR OUTPUT

- **report.html**: Main report in HTML format
- **report.csv** and **whole_app_metric.csv**: Comma-separated CSV files
- **program_tree.dot**: A graphical representation of the call tree showing the offloadable and accelerated regions

- **program_tree.pdf**: A graphical representation of the call tree

Generated if the DOT(GraphViz*) utility is installed

1:1 conversion from the **program_tree.dot** file

- **JSON** and **LOG** files that contain data used to generate the HTML report and logs, primarily used for debugging and reporting bugs and issues

HOW TO SPEED UP COLLECTION

AVOID DEPENDENCY CHECKING

- Dependency adds a lot of time to the collection and you might want to remove it.
- Add the option `-c basic` for the collection:

```
advixe-python <ADV_INSTALL_DIR>/perfmodels/run_oa.py <path_to_result_dir>  
-config gen9 -c basic --out-dir <path_to_result_dir> [--options] -- <app>
```

- Add the option `--assume-parallel` for the analysis:

```
advixe-python $APM/analyse.py <project_dir> --assume-parallel --config  
gen9 [--options] -- <app_binary> [app_options]
```

GET BETTER CONTROL OF YOUR RUN

INDIVIDUAL LOW LEVEL COMMANDS

- You might want to run the command lines independently to tweak the parameters
- A good start is to use run_oa.py script with --dry-run to see the list of command lines and retrieve the cache configuration of the target accelerator.
- The next command will output the different command lines for doing separate analyses without running advisor collection.
- **advixe-python <ADV_INSTALL_DIR>/perfmodels/run_oa.py
<path_to_result_dir> --dry-run -config gen9 -c basic --out-dir
<path_to_result_dir> [--options] -- <app>**

GET BETTER CONTROL OF YOUR RUN

INDIVIDUAL LOW LEVEL COMMANDS

- We start with the survey
- **`advixe-cl --collect=survey --auto-finalize --stackwalk-mode=online --static-instruction-mix --project-dir=./oa_report - my_app`**
- The survey times your application and run some static analysis on the binary without impact on the application's performance.
 - Sampling
 - Binary static analysis
 - Static code analysis (compiler and debug infos)

GET BETTER CONTROL OF YOUR RUN

INDIVIDUAL LOW LEVEL COMMANDS

- We continue with the trip count and cache simulation
- **`advixe-cl --collect=tripcounts -return-app-exitcode -flop -stacks -auto-finalize -ignore-checksums -enable-data-transfer-analysis -track-heap-objects -profile-jit -cache-sources -enable-cache-simulation -cache-config=1:8w:32k/1:64w:512k/1:16w:8m --project-dir=./oa_report - my_app`**
- The tripcounts with `-flop` and `-cache-simulation` counts:
 - The number of iterations in your loops
 - The number of operations
 - Evaluate the data transfers between memory subsystems configured with `-cache-config`
- This analysis has usually $\approx 10\times$ slowdown

GET BETTER CONTROL OF YOUR RUN

INDIVIDUAL LOW LEVEL COMMANDS

- Optional step: Dependency analysis
- **`advixe-cl --collect=dependencies --loops="total-time>5" --filter-reductions --loop-call-count-limit=16 --project-dir=./oa_report - my_app`**
- Detects data dependencies in your loop by checking your memory accesses
- This analysis has an important impact on the performance
- It is up to the user to define how loops will be selected for this analysis, here we use `--loops="total-time>5"` which select all loops impacting more than 5% of the overall time

GET BETTER CONTROL OF YOUR RUN

INDIVIDUAL LOW LEVEL COMMANDS

- Last step: Generating the report
- 2 Cases:
 - You ran the dependency analysis:

```
advixe-python $APM/analyse.py ./oa_report --config gen9 --out-dir  
oa_report - my_app
```

- You didn't run the dependency analysis

```
advixe-python $APM/analyse.py ./oa_report --config gen9 --assume-  
parallel --out-dir oa_report - my_app
```

MODELING PERFORMANCE

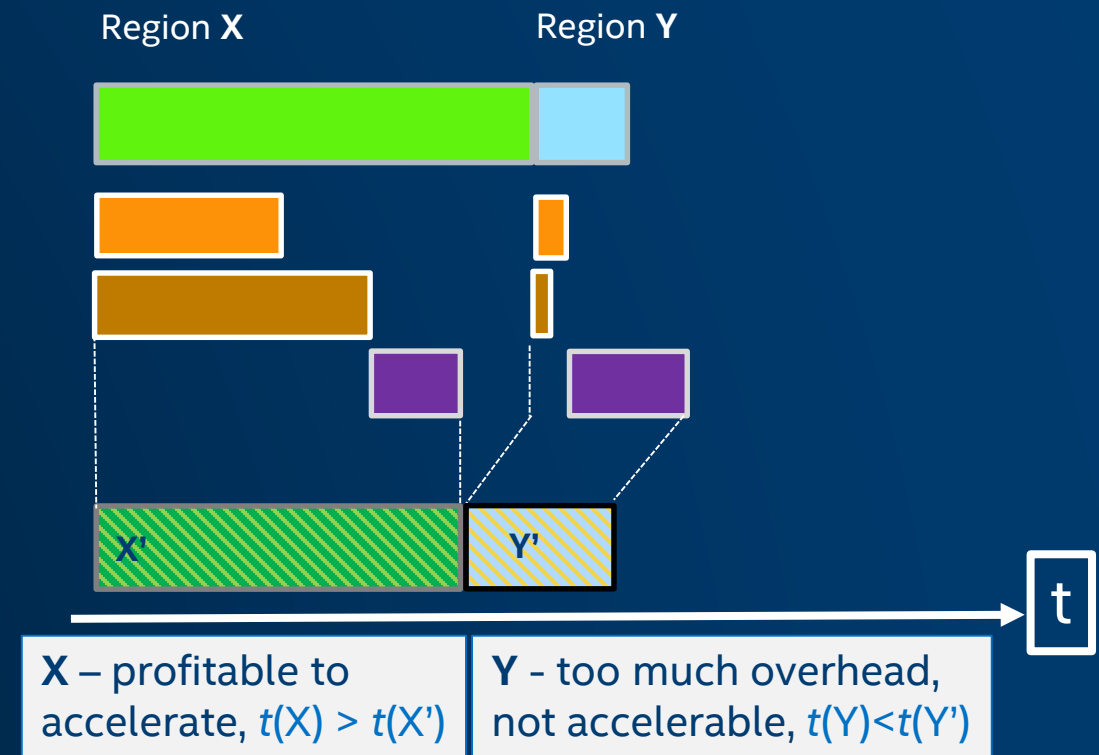
First order analytical modeling pillars:

- Compute throughput model
- Memory sub-system model
- Offload data transfer modeling

Execution time on baseline platform (CPU)

- Execution time on accelerator. Estimate assuming bound exclusively by Compute
- Execution time on accelerator. Estimate assuming bound exclusively by caches/memory
- Offload Tax estimate (data transfer + invoke)

Final estimated time on target platform (eg GPU)



$$t_{\text{region}} = \max(t_{\text{compute}}, t_{\text{memory subsystem}}) + t_{\text{data transfer tax}} + t_{\text{kernel launch}}$$

MODELING PERFORMANCE

We minimize the total time spent in this loop hierarchy by varying offload strategies U (offload/non-offload, #threads for each component $loop_i$ of loopnest)

$$\textbf{Objective function} : T_{all} = \min_{U=\{uf_1, uf_2, \dots\}} (\sum_i T_i + t_{data\ transfer} + t_{invoke} + T_{cpu})$$

Reject loopnests for which
 $T(x86) / T_{all}(x86+\"X\") < 1.0$

$$T_i = \max \begin{cases} T_i^{Comp_only} () \\ T_i^{M_k_only} (M_i^k) = \frac{M_i^k}{BW_k} \end{cases}$$

This is effectively “balance”
(throughput) model

Under algorithmic constraints (Dependencies and TripCount/Granularity)

INTEL[®] ADVISOR

GPU ROOFLINE ANALYSIS

[Optimization Notice](#)

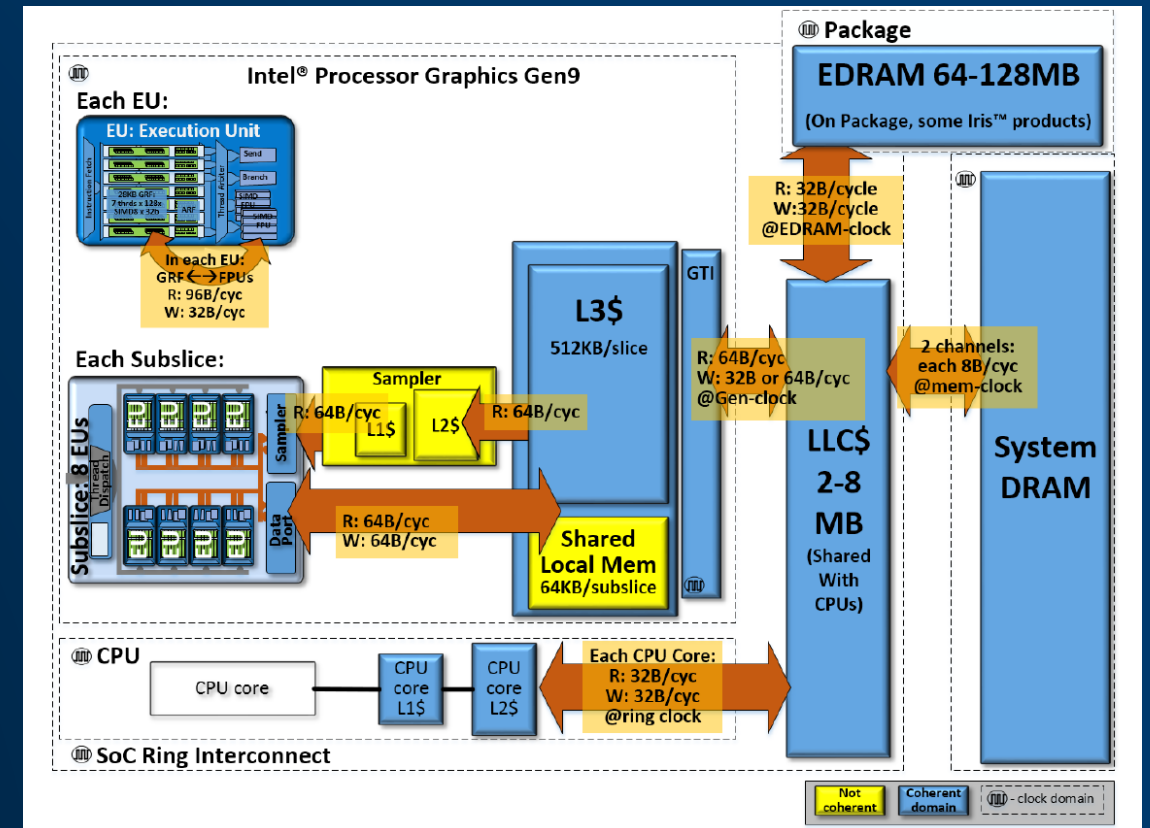
Copyright © 2019, Intel Corporation. All rights reserved.

*Other names and brands may be claimed as the property of others.



INTEL® GEN9 MEMORY HIERARCHY

- Intel® Graphics Compute Architecture uses the same DRAM with the CPU
- Level-3 (L3) data cache: slice-shared asset
- Shared Local Memory (SLM): a dedicated structure within the L3 that supports the work-group local memory address space
- Graphics Technology Interface (GTI): a dedicated interface unit connects the entire architecture interfaces to the rest of the SoC components
- The rest of SoC memory hierarchy includes the large Last-Level Cache (LLC, which is shared between CPU and GPU), possibly embedded DRAM and finally the system DRAM



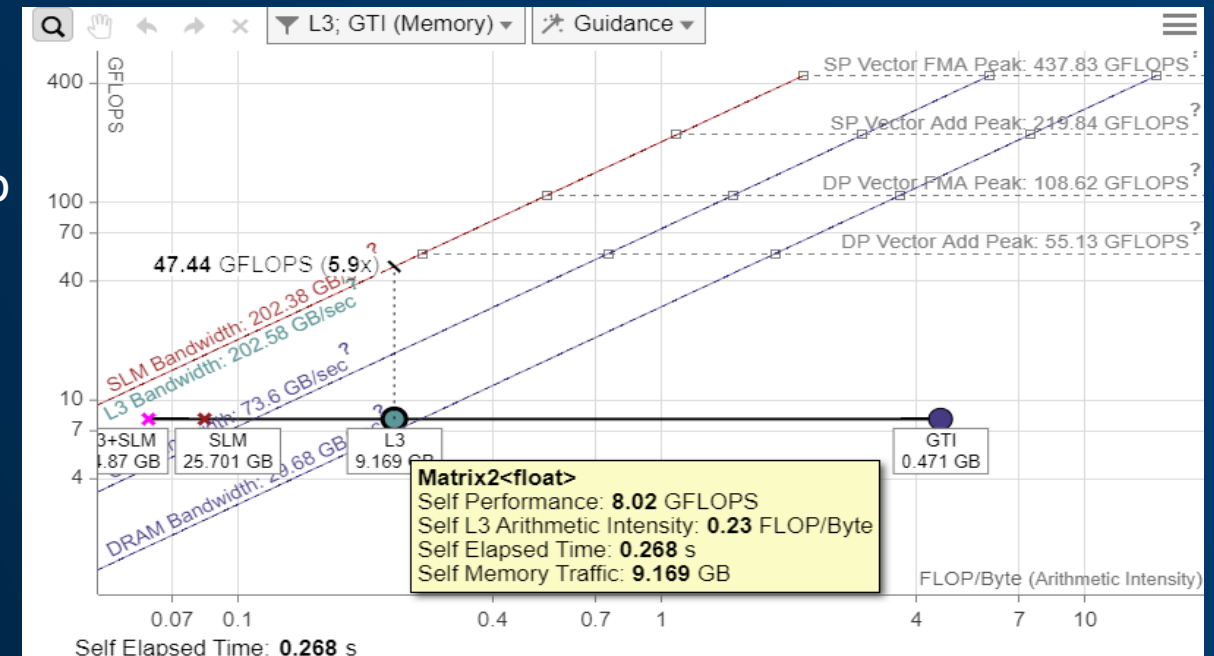
A view of the SoC chip level memory hierarchy and its theoretical peak bandwidths for the compute architecture of Intel processor graphics gen9

FIND EFFECTIVE OPTIMIZATION STRATEGIES

INTEL® ADVISOR - GPU ROOFLINE

GPU Roofline Performance Insights

- Highlights poor performing loops
- Shows performance 'headroom' for each loop
 - Which can be improved
 - Which are worth improving
- Shows likely causes of bottlenecks
 - Memory bound vs. compute bound
- Suggests next optimization steps



RUNNING INTEL® ADVISOR GPU ROOFLINE ANALYSIS

The Roofline model on GPU is a technical preview feature and is not available by default. To enable it:

```
export ADVIXE_EXPERIMENTAL=gpu-profiling
```

To run the GPU Roofline analysis in the Intel® Advisor CLI:

Run the Survey analysis with the **--enable-gpu-profiling** option:

```
advixe-cl -collect=survey --enable-gpu-profiling --project-dir=<my_project_directory> --search-dir  
src:r=<my_source_directory> -- ./myapp [app_parameters]
```

Run the Trip Counts and FLOP analysis with **--enable-gpu-profiling** option:

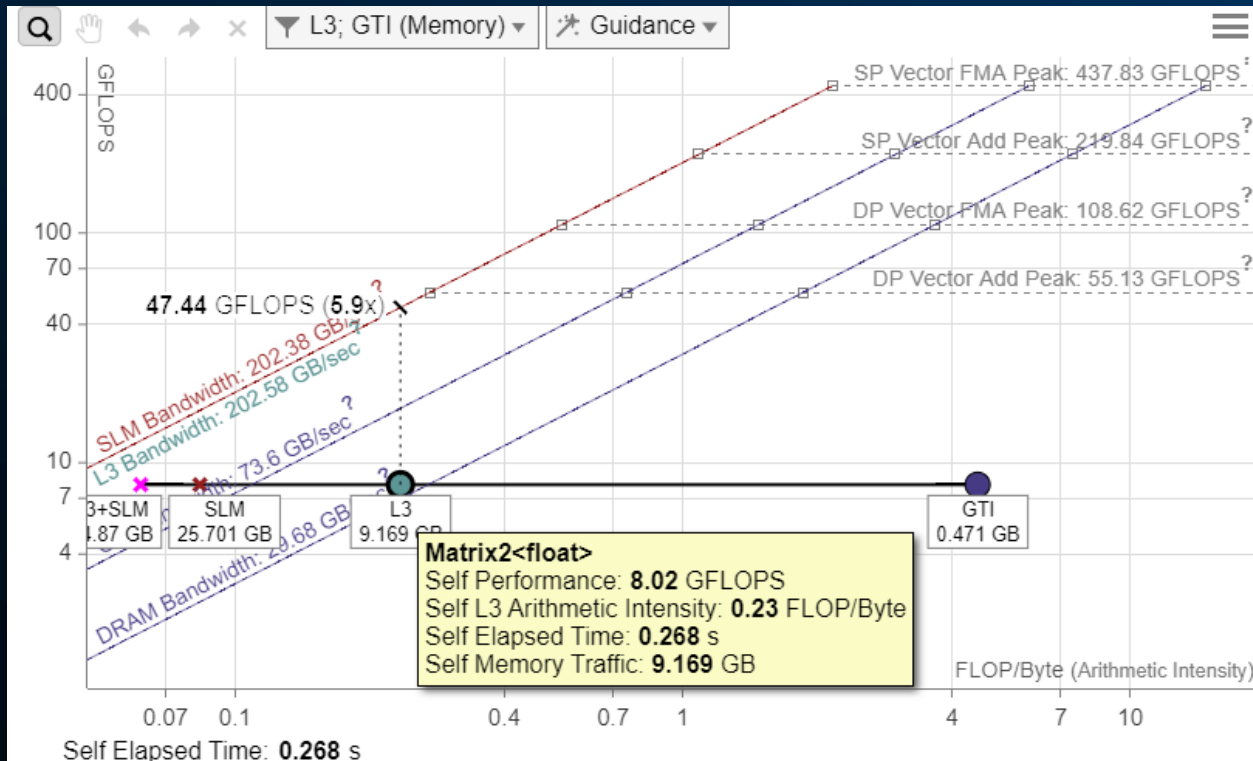
```
advixe-cl -collect=tripcounts --stacks --flop --enable-gpu-profiling --project-  
dir=<my_project_directory> --search-dir src:r=<my_source_directory> -- ./myapp [app_parameters]
```

Generate a GPU Roofline report:

```
advixe-cl --report=roofline --gpu --project-dir=<my_project_directory> --report-output=roofline.html
```

Open the generated roofline.html in a web browser to visualize GPU performance.

ROOFLINE ANALYSIS ON INTEL® GPU



As an example you can see from the roofline chart, our L3 dot is very close to the L3 maximum bandwidth, to get more FLOPS we need to optimize our caches further. A cache blocking optimization strategy can make better use of memory and should increase our performance. The GTI (traffic between our GPU, GPU uncore (LLC) and main memory) is far from the GTI roofline, transfer costs between our CPU to GPU does not seem to be an issue.

INTEL® VTUNE® PROFILER

GPU ANALYSIS

[Optimization Notice](#)

Copyright © 2019, Intel Corporation. All rights reserved.

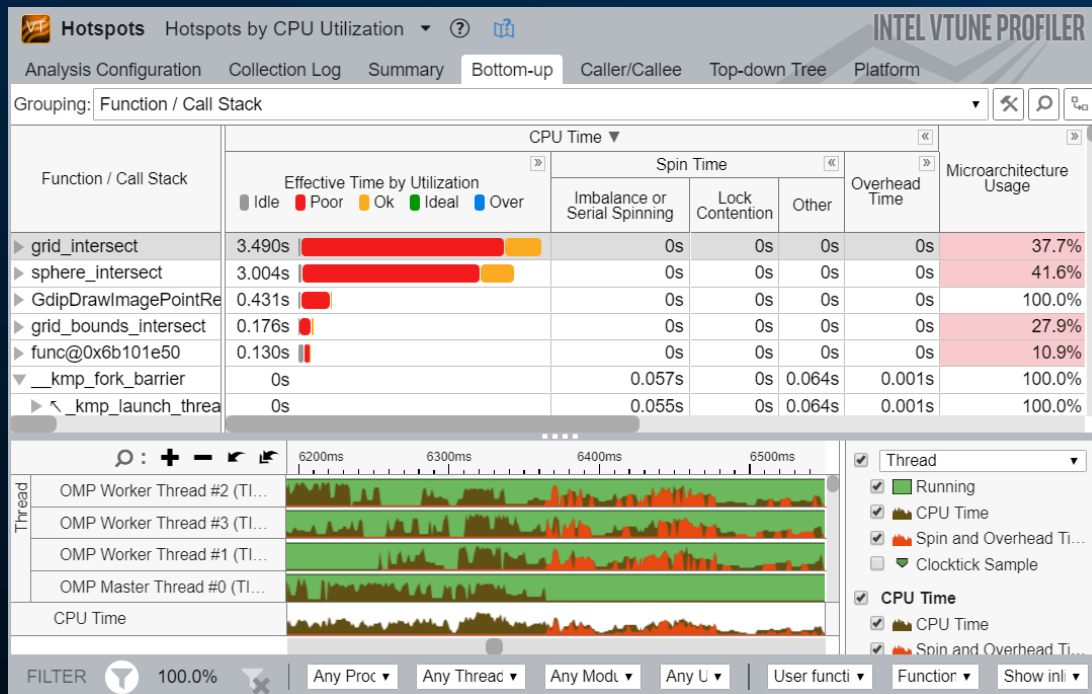
*Other names and brands may be claimed as the property of others.



Analyze & Tune Application Performance

Intel® VTune™ Profiler — Performance Profiler

New more descriptive name!



Save Time Optimizing Code

- Accurately profile C, C++, Fortran*, Python*, Go*, Java*, or any mix
- Optimize CPU, threading, memory, cache, storage & more
- Save time: rich analysis leads to insight
- Take advantage of [Priority Support](#)
 - Connects customers to Intel engineers for confidential inquiries (paid versions)


What's New in 2020 Release (partial list)

- Production release of Platform Profiler - Longer Data Collection
- Design & optimize for Intel® Optane™ DC Persistent Memory
- Application Performance Snapshot – adds communication pattern diagnosis, profiles more ranks
- Linux* – extensive additional use of Perf enables analysis without adding drivers

Learn More: software.intel.com/vtune

OPTIMIZE YOUR GPU USAGE USING INTEL® VTUNE PROFILER

GPU OFFLOAD

 GPU Offload (Preview) GPU Offload (Preview) ? ⓘ

Analysis Configuration Collection Log Summary Graphics Platform

⌵

Elapsed Time ⓘ: 2.017s

⌵

GPU Usage ⓘ: 47.8% ⓘ

Use this section to understand whether the GPU was utilized properly and which of the engines were utilized. Identify the amount of gaps in the GPU utilization that potentially could be loaded with some work. This metric is calculated for the engines that had at least one piece of work scheduled to them.

⌵

GPU Usage

GPU Usage breakdown by GPU engines and work types.

GPU Engine / Packet Type	GPU Time	GPU Utilization ⓘ
Render and GPGPU	0.964s	47.8% ⓘ
Unknown	0.964s	47.8%

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

⌵

Hottest GPU Computing Tasks

This section lists the most active computing tasks running on the GPU, sorted by the Total Time. Focus on the computing tasks flagged as performance-critical.

Computing Task	Total Time	Total Compute Time	Total Transfer Time (f)
Matrix<float> ⓘ	3.980s	0.961s	3.019s
clEnqueueReadBufferRect ⓘ	0.000s	0.000s	0.000s

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

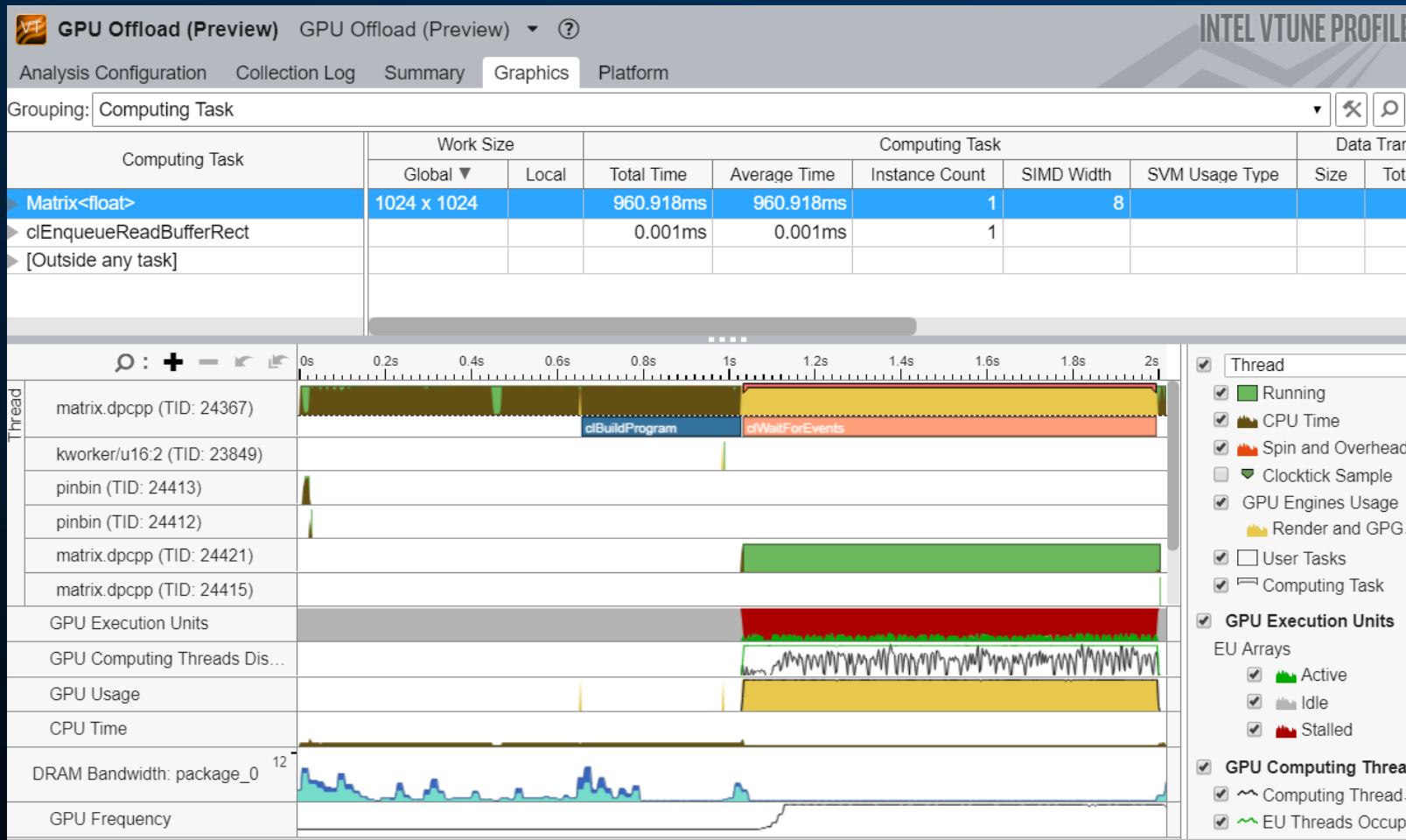
This analysis enables you to:

- Identify how effectively your application uses DPC++ or OpenCL kernels.
- Analyze execution of Intel Media SDK tasks over time (for Linux targets only)
- 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 of Intel® VTune™ Profiler to see how effectively we are using our GPU.

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

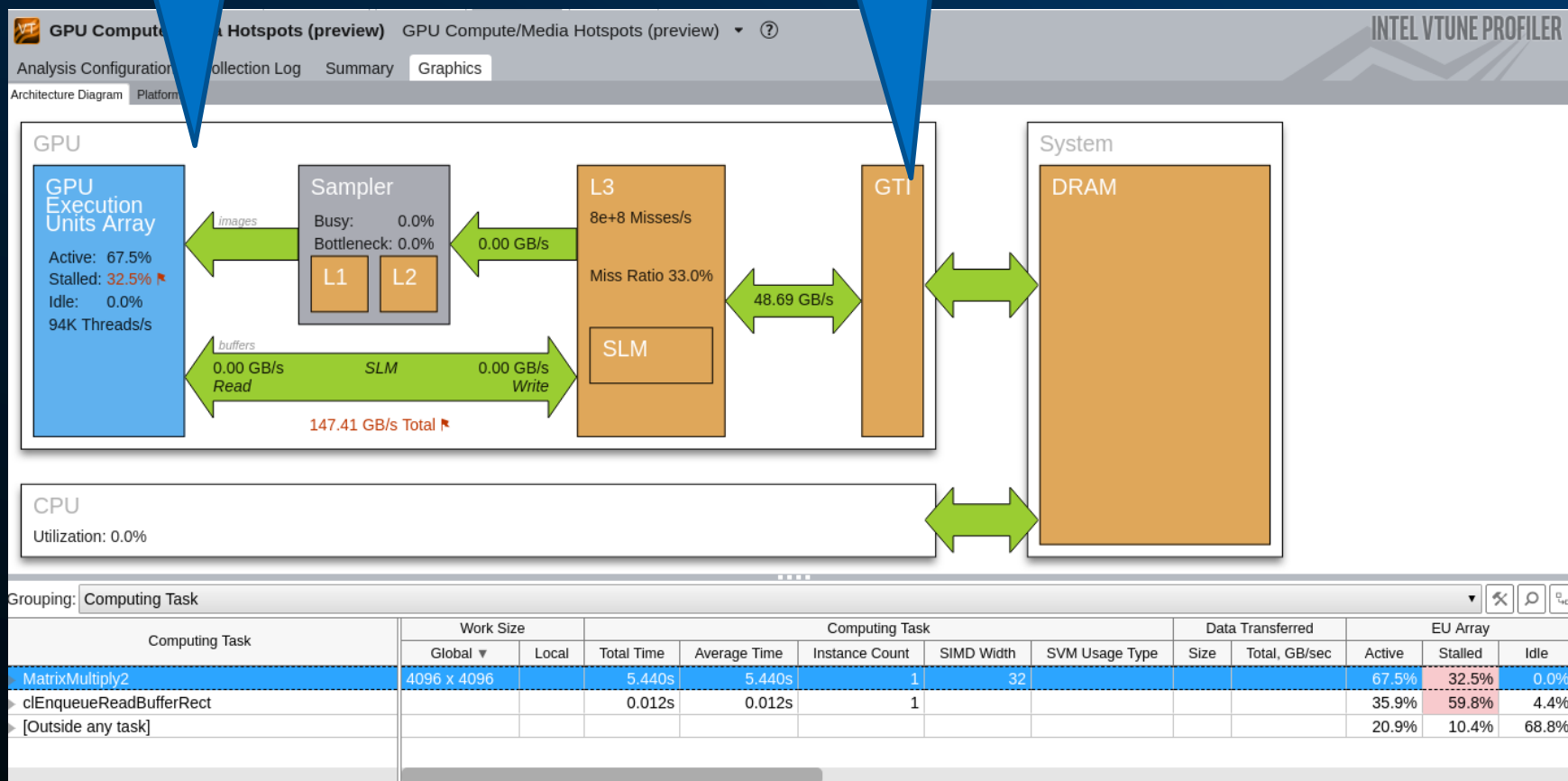
OPTIMIZE YOUR GPU USAGE USING INTEL® VTUNE PROFILER

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.



[Optimization Notice](#)

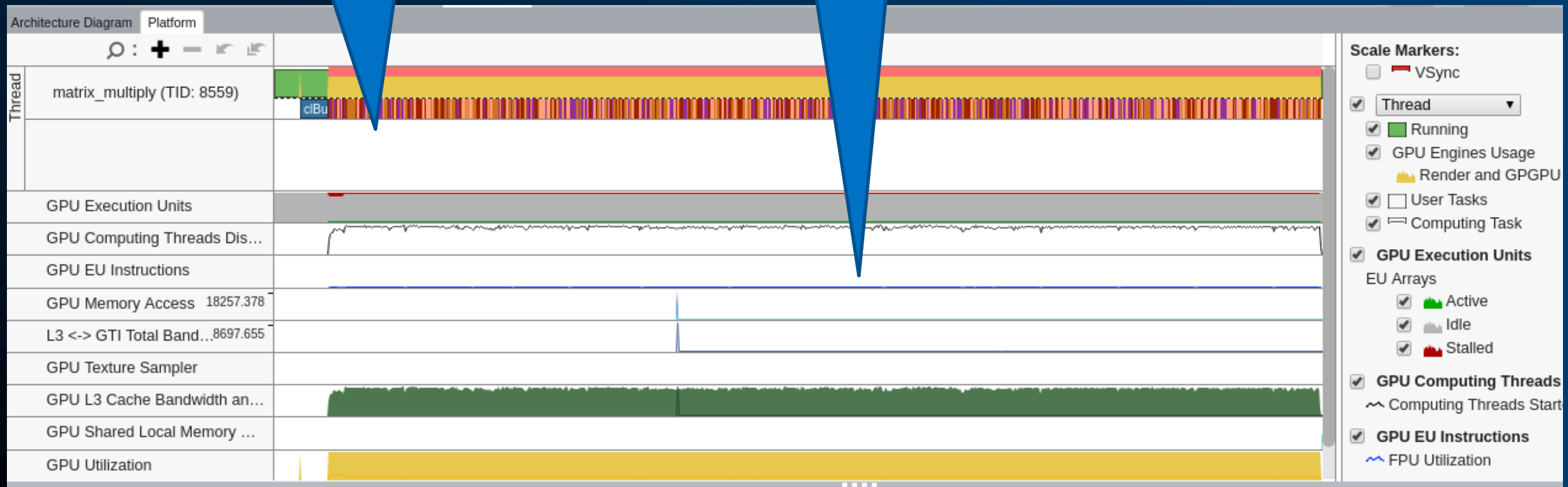
Copyright © 2019, Intel Corporation. All rights reserved.

*Other names and brands may be claimed as the property of others.

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

Source		🔥 Computing Task			Data Transferred		EU Array		
		Total Time	Average Time	Instance Count	Size	Total, GB/sec	Active	Stalled	Idle
183	hdlr.parallel_for<class MatrixMultiply3>(matrixRange, [=](dpcpp::id<2> id)	5.484s	5.484s	1			53.9%	46.1%	0.0%
184	{								
185	size_t i = id[0], j = id[1];								
186									
187	T c = {};								
188	for (size_t k = 0; k < w; k++)								
189	{								
190	c += ra[i][k] * rb[k][j];								
191	}								

Performance metrics
at the DPC++
statement level

Source		Assembly grouping: Address	
		Address	Source
179	auto ra = bA.template get_access<dpcpp::access::mode::read>(hd	0x10	183 (W) or (1 M0) cr0.0<1>:ud cr0.0<0>:1,0>:ud 0x4c0:uw {Switch}
180	auto rb = bB.template get_access<dpcpp::access::mode::read>(hd	0x20	183 (W) mul (1 M0) r5.0<1>:d r14.1<0>:1,0>:d r97.6<0>:1,0>:d
181	auto rc = bC.template get_access<dpcpp::access::mode::discard_w	0x30	183 (W) mul (1 M0) r23.0<1>:d r14.0<0>:1,0>:d r97.1<0>:1,0>:d {Compacted}
182		0x38	mov (16 M0) r98.0<1>:f 0x0:f
183	hdlr.parallel_for<class MatrixMultiply3>(matrixRange, [=](dpcpp	0x48	mov (16 M16) r100.0<1>:f 0x0:f
184	{	0x58	188 add (8 M0) r15.0<1>:q r3.0<8>:8,1>:uw r5.0<0>:1,0>:ud
185	size_t i = id[0], j = id[1];	0x68	188 add (8 M8) r17.0<1>:q r3.8<8>:8,1>:uw r5.0<0>:1,0>:ud
186		0x78	188 add (8 M16) r19.0<1>:q r4.0<8>:8,1>:uw r5.0<0>:1,0>:ud
187	T c = {};	0x88	188 add (8 M24) r21.0<1>:q r4.8<8>:8,1>:uw r5.0<0>:1,0>:ud
188	for (size_t k = 0; k < w; k++)	0x98	188 add (8 M0) r24.0<1>:q r1.0<8>:8,1>:uw r23.0<0>:1,0>:ud
189	{	0xa8	188 add (8 M8) r26.0<1>:q r1.8<8>:8,1>:uw r23.0<0>:1,0>:ud
190	c += ra[i][k] * rb[k][j];	0xb8	188 add (8 M16) r28.0<1>:q r2.0<8>:8,1>:uw r23.0<0>:1,0>:ud
191	}	0xc8	188 add (8 M0) r103.0<1>:q r15.0<4>:4,1>:q r7.1<0>:1,0>:ud
192	rc[i][j] = c;	0xd8	188 add (8 M8) r105.0<1>:q r17.0<4>:4,1>:q r7.1<0>:1,0>:ud
193	});	0xe8	188 add (8 M16) r107.0<1>:q r19.0<4>:4,1>:q r7.1<0>:1,0>:ud
194	});	0xf8	188 add (8 M24) r109.0<1>:q r21.0<4>:4,1>:q r7.1<0>:1,0>:ud
195	}	0x108	188 add (8 M24) r30.0<1>:q r2.8<8>:8,1>:uw r23.0<0>:1,0>:ud
		0x118	(W) mov (1 M0) r102.0<1>:q 0:w

GPU assembly
available for
compute kernels

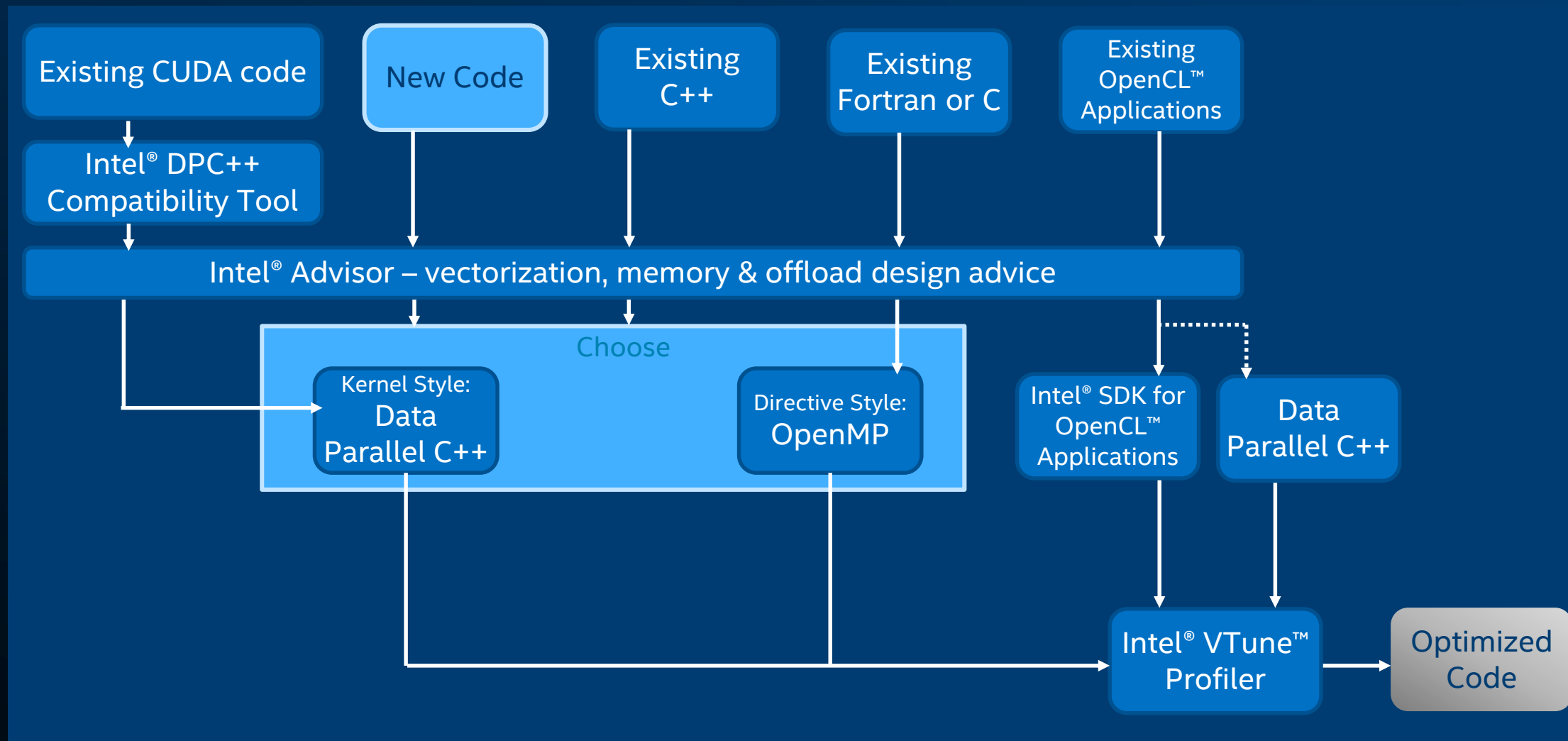
CASE STUDY/DEMO

PUTTING IT ALL TOGETHER

GPU ANALYSIS

HPC ONEAPI SINGLE NODE WORKFLOW

I WANT TO ACCELERATE USING DIRECT PROGRAMMING ON A GPU...



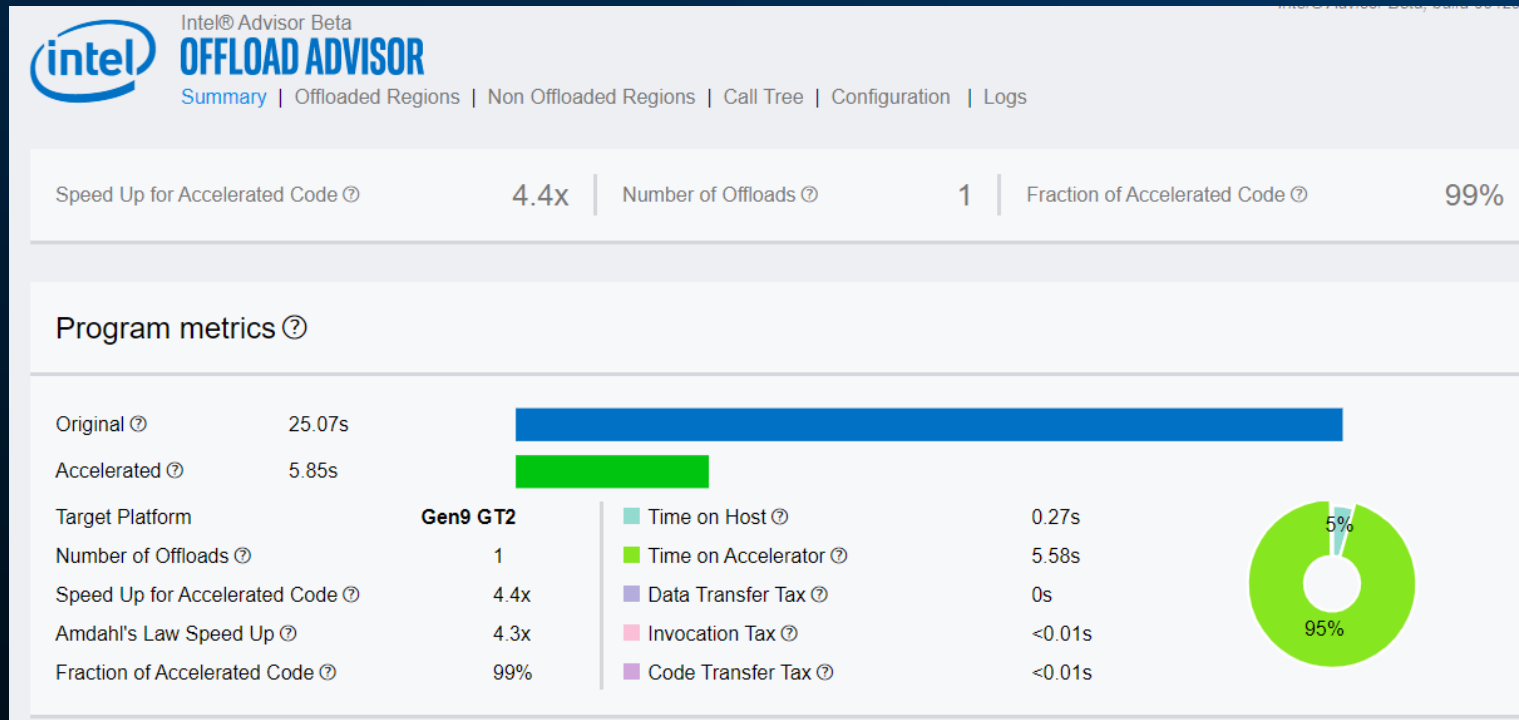
Purpose: Demo performance analysis tools for Intel® oneAPI

Steps:

1. Run Intel® Advisor - offload advisor on your CPU. Generate a report that indicates what should be offloaded to your accelerator.
2. After porting code to DPC++. Run Intel® VTune™ Profiler offload analysis and GPU hotspots. Show VTune Profiler results indicating places for improvement.
3. Re-Run a second VTune Profiler gpu-hotspot analysis on our optimized gpu kernel.
4. Run Intel Advisor GPU Roofline on our GPU workload to see where we can improve based on system maximums(roofs).

STEP 1 - RUN INTEL® ADVISOR - OFFLOAD ADVISOR

TO FIND CODE THAT CAN BE PROFITABLY OFFLOADED



Some key observations

The workload was accelerated 4.4x
You can see in program metrics that the original workload ran in 25.07s and the accelerated workload ran in 5.85s

STEP 1 - RUN INTEL® ADVISOR - OFFLOAD ADVISOR

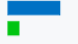

TO FIND CODE THAT CAN BE PROFITABLY OFFLOADED TO YOUR ACCELERATOR


Click on the offloaded regions link.
Additional details are displayed including:

Specific line of code to be offloaded with the source.

Offload information

Data transfer information including the amount of data transferred from CPU->GPU. We can also see we are bounded by LLC_BW

Top offloaded ?			
Location ?	Speed Up ?	Bounded By ?	Data Transfer ?
[loop in multiply1\$omp\$parallel@201 at multiply.c:202]	4.44x  CPU 24.80s GPU 5.58s	 LLC_BW	100.68MB

 Intel® Advisor Beta
OFFLOAD ADVISOR
Summary | [Offloaded Regions](#) | Non Offloaded Regions | Call Tree | Configuration | Logs

Speed Up for Accelerated Code ? 4.4x

Number of Offloads ? 1

Fraction of Accelerated Code ? 99%

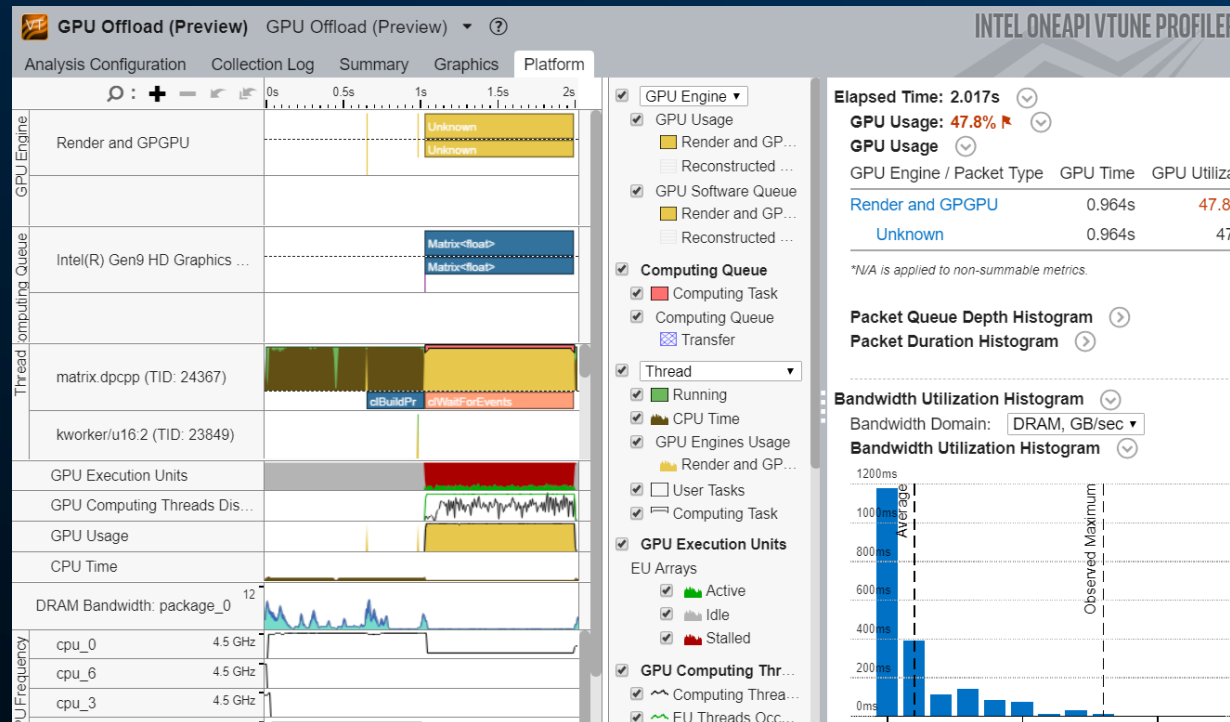
Hierarchy	Loop/Function >		Offload Information >					Overhead
	Elapsed Time (s)	Total Time ↓	Dependency Type	Estimated Speed Up	Estimated Time on Accelerator	Parallel Threads	Bounded By	
> [loop in multiply1\$omp\$parallel@201 a	24.80s	196.325s (99.21%)	Parallel: Explicit	4.44x	5.584s (95.39%)	48	LLC_BW	<0.1

Source Name: [loop in multiply1\$omp\$parallel@201 at multiply.c:202]

```
200 #pragma omp parallel for
201   for(i=0; i<msize; i++) {
202     for(j=0; j<msize; j++) {
203       for(k=0; k<msize; k++) {
204         c[i][j] = c[i][j] + a[i][k]
205       }
```

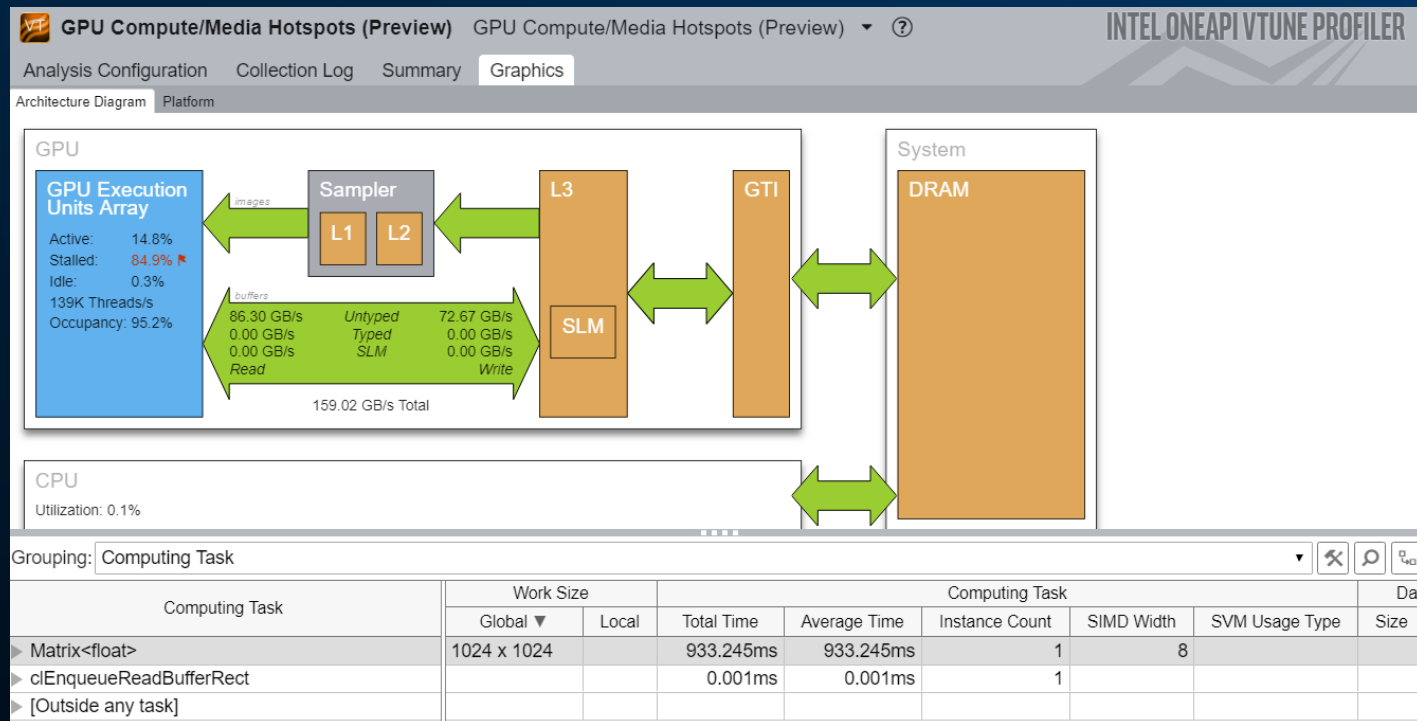
STEP 2: RUN INTEL® VTUNE™ PROFILER

TO OPTIMIZE YOUR USE OF GPU (OFFLOAD)




STEP 2: RUN INTEL® VTUNE™ PROFILER

TO OPTIMIZE YOUR USE OF GPU (GPU HOTSPOTS)



STEP 3: RERUN INTEL® VTUNE™ PROFILER

AFTER MAKING AN OPTIMIZATION, RE-RUN GPU-HOTSPOTS TO SEE WHAT HAS IMPROVED

 GPU Offload (Preview) GPU Offload (Preview) ?

INTEL ONEAPI VTUNE PROFILER

Analysis Configuration Collection Log Summary Graphics Platform

⌵

Elapsed Time ⓘ: 1.182s

⌵

GPU Usage ⓘ: 15.7% 🚩

Use this section to understand whether the GPU was utilized properly and which of the engines were utilized. Identify the amount of gaps in the GPU utilization that potentially could be loaded with some work. This metric is calculated for the engines that had at least one piece of work scheduled to them.

⌵

GPU Usage

GPU Usage breakdown by GPU engines and work types.

GPU Engine / Packet Type	GPU Time	GPU Utilization ⓘ
Render and GPGPU	0.185s	15.7% 🚩
Unknown	0.185s	15.7%

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

⌵

Packet Queue Depth Histogram

⌵

Packet Duration Histogram

⌵

Hottest GPU Computing Tasks

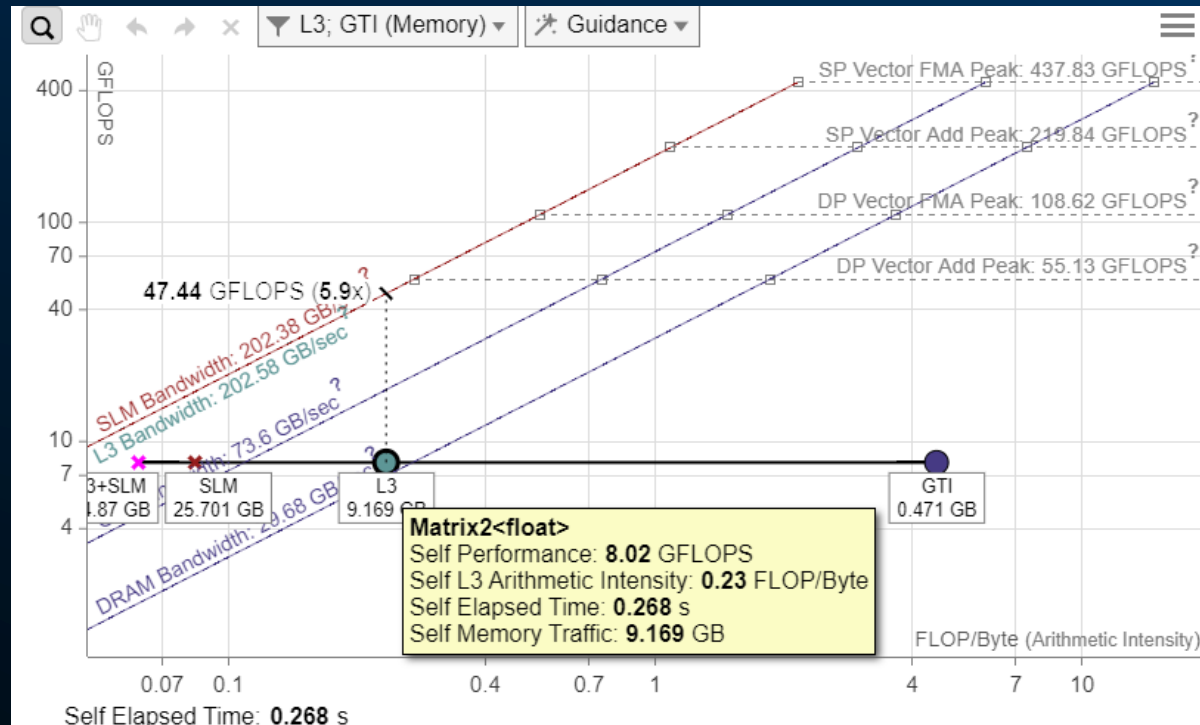
This section lists the most active computing tasks running on the GPU, sorted by the Total Time. Focus on the computing tasks flagged as performance-critical.

Computing Task	Total Time	Total Compute Time	Total Transfer Time (f)
Matrix2<float> 🚩	0.759s	0.183s	0.575s
clEnqueueReadBufferRect 🚩	0.000s	0.000s	0.000s

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

STEP 4: RUN INTEL® ADVISOR GPU ROOFLINE

TO SEE HOW CLOSE YOU ARE TO THE SYSTEM MAXIMUMS (ROOFLINES)



SUMMARY

Intel® oneAPI includes some powerful analysis tools!

- Intel® Advisor and Intel® VTune™ Profiler now have features for analyzing your GPU.
- Intel® Advisor and Intel® VTune™ Profiler provide an effective workflow for optimizing your workloads for running on Intel® Accelerators.
- The workflow presented here shows how the tools complement each other.

NOTICES AND DISCLAIMERS

This document contains information on products, services and/or processes in development. All information provided here is subject to change without notice. Contact your Intel representative to obtain the latest forecast, schedule, specifications and roadmaps.

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. No product or component can be absolutely secure. Check with your system manufacturer or retailer or learn more at [intel.com].

Intel, the Intel logo is trademark of Intel Corporation or its subsidiaries in the U.S. and/or other countries.

Other names and brands may be claimed as the property of others

© Intel Corporation.

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

