

COMPREHENSIVE GENERAL LUMI COURSE

26/4/2024

INTRODUCTION TO OMNIPERF

GINA SITARAMAN, SUYASH TANDON, GEORGE MARKOMANOLIS,
JONATHAN MADSEN, AUSTIN ELLIS, BOB ROBEY, XIAOMIN LU,
NOAH WOLFE, SAMUEL ANTAO

JAKUB KURZAK - PRESENTER

ADVANCED MICRO DEVICES, INC.

AMD 
together we advance_

slides on LUMI in /project/project_465001098/Slides/AMD/

hands-on exercises: https://hackmd.io/@gmarkoma/lumi_finland

hands-on source code: /project/project_465001098/Exercises/AMD/HPCTrainingExamples/

Background – AMD Profilers

ROC-profiler (rocprof)

Hardware Counters

- Raw collection of GPU counters and traces
- Counter collection with user input files
- Counter results printed to a CSV

Traces and timelines

- Trace collection support for CPU copy, HIP API, HSA API, GPU Kernels

Visualisation

- Traces visualized with Perfetto

	A	B	C	D	E
1	Name	Calls	TotalDura	AverageN	Percentage
2	hipMemcpyAsync	99	3.22E+10	3.25E+08	44.14872
3	hipEventSynchronize	330	2.42E+10	73394557	33.225
4	hipMemsetAsync	87	7.76E+09	89232696	10.64953
5	hipHostMalloc	9	5.41E+09	6.01E+08	7.415198
6	hipDeviceSynchronize	28	1.32E+09	47006288	1.805515
7	hipHostFree	17	1.05E+09	61534688	1.435014
8	hipMemcpy	41	8.11E+08	19791876	1.113161
9	hipLaunchKernel	1856	58082083	31294	0.079676
10	hipStreamCreate	2	46380834	23190417	0.063625
11	hipMemset	2	18847246	9423623	0.025854
12	hipStreamDestroy	2	15183338	7591669	0.020828
13	hipFree	38	8269713	217624	0.011344
14	hipEventRecord	330	2520035	7636	0.003457
15	hipMalloc	30	1484804	49493	0.002037
16	__hipPopCallConfigura	1856	229159	123	0.000314
17	__hipPushCallConfigur	1856	224177	120	0.000308
18	hipGetLastError	1494	100458	67	0.000138
19	hipEventCreate	330	76675	232	0.000105
20	hipEventDestroy	330	64671	195	8.87E-05
21	hipGetDevicePropertie	47	51808	1102	7.11E-05
22	hipGetDevice	64	11611	181	1.59E-05
23	hipSetDevice	1	401	401	5.50E-07
24	hipGetDeviceCount	1	220	220	3.02E-07

Omnitrace

Trace collection

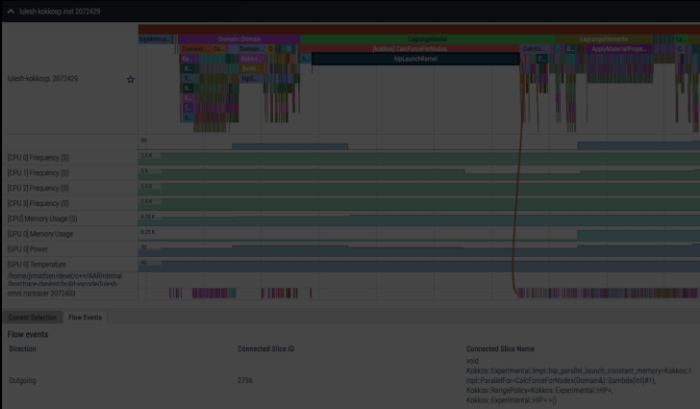
- Comprehensive trace collection
- CPU, GPU

Supports

- CPU copy, HIP API, HSA API, GPU Kernels
- OpenMP®, MPI, Kokkos, p-threads, multi-GPU

Visualisation

- Traces visualized with Perfetto



Omniperf

Performance Analysis

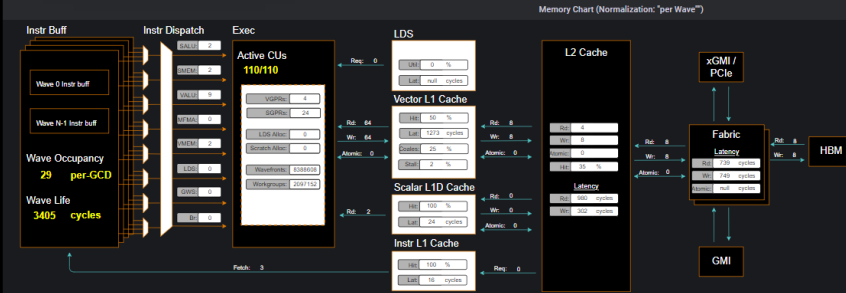
- Automated collection of hardware counters
- Analysis, Visualisation

Supports

- Speed of Light, Memory chart, Rooflines, Kernel comparison

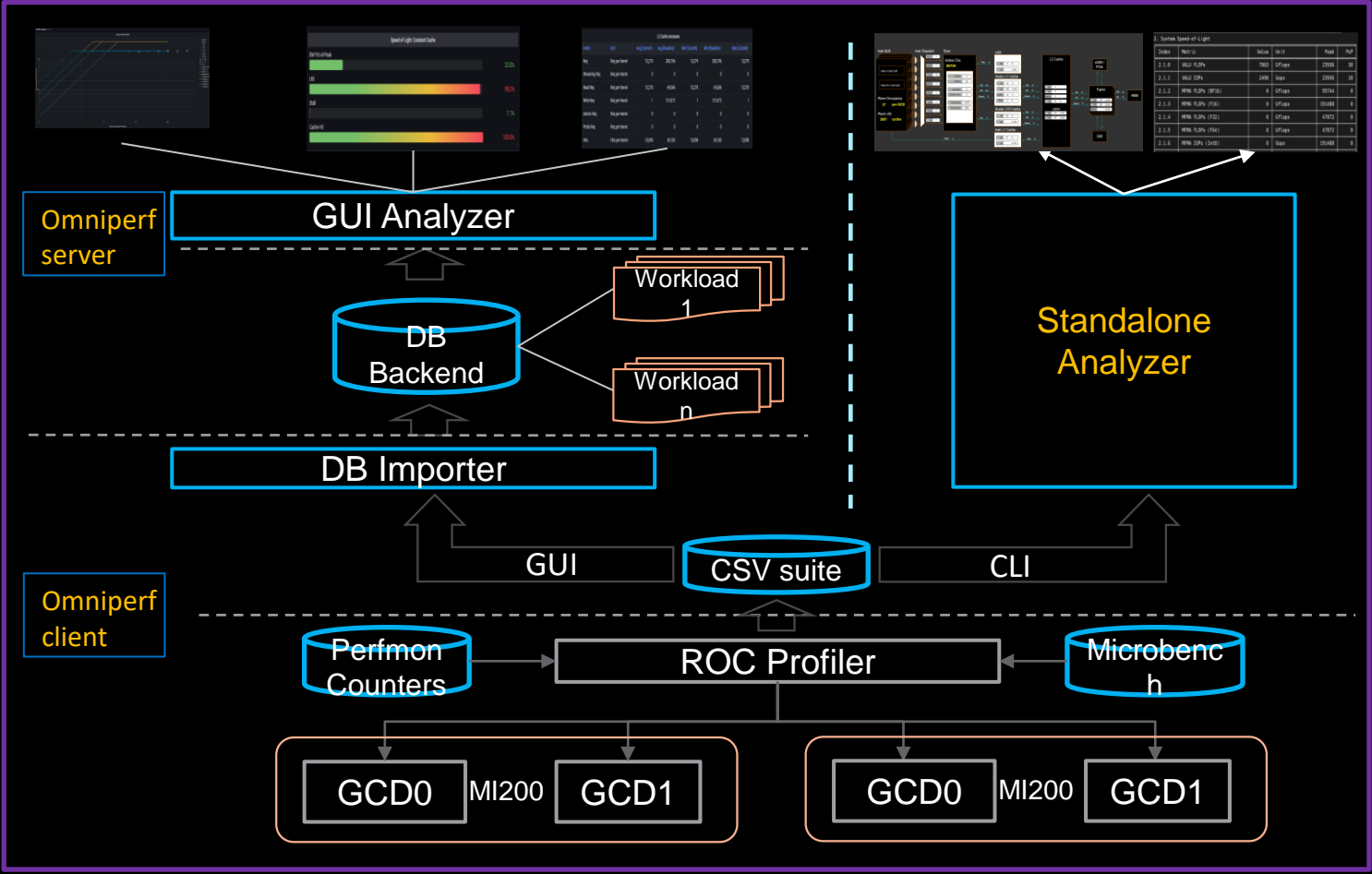
Visualisation

- With Grafana or standalone GUI



Omniperf: Automated Collection of Hardware Counters and Analysis

AMD Research Tool	Repository: https://github.com/AMDResearch/omniperf			
	X Not part of ROCm stack	Built on top of ROC-profiler		
Integrated Performance Analyzer for AMD GPUs	Speed-of-Light	Roofline	Memory chart	Baseline comparison
	Sub-system performance analysis			
	LDS	vL1D	L2 Cache	HBM
	Shader Compute	Wavefront	Instruction mix	Latencies
INSTINCT™ Support	MI200	MI100		
User Interfaces	Grafana™ GUI	Standalone GUI	Command Line (CLI)	



Refer to [current documentation](#) for recent updates

Omniperf features

Omniperf Features	
MI200 support	Roofline Analysis Panel (<i>Supported on MI200 only, SLES 15 SP3 or RHEL8</i>)
MI100 support	Command Processor (CP) Panel
Standalone GUI Analyzer	Shader Processing Input (SPI) Panel
Grafana/MongoDB GUI Analyzer	Wavefront Launch Panel
Dispatch Filtering	Compute Unit - Instruction Mix Panel
Kernel Filtering	Compute Unit - Pipeline Panel
GPU ID Filtering	Local Data Share (LDS) Panel
Baseline Comparison	Instruction Cache Panel
Multi-Normalizations	Scalar L1D Cache Panel
System Info Panel	Texture Addresser and Data Panel
System Speed-of-Light Panel	Vector L1D Cache Panel
Kernel Statistic Panel	L2 Cache Panel
Memory Chart Analysis Panel	L2 Cache (per-Channel) Panel

Omniperf

- Omniperf is an integrated performance analyzer for AMD GPUs built on ROCprofiler
- Omniperf executes the code many times to collect various hardware counters (over 100 counters default behavior)
- Using specific filtering options (kernel, dispatch ID, metric group), the overhead of profiling can be reduced
- Roofline analysis is supported on MI200 GPUs
- Omniperf shows many panels of metrics based on hardware counters, we will show a few here
- Typical Omniperf workflows:
 - Profile + Analyze with CLI or visualize with standalone GUI
 - Profile + Import to database and visualize with Grafana
- Omniperf targets MI100 and MI200 and future generation AMD GPUs
- Omniperf requires to use just 1 MPI process
- For problems, create an issue here: <https://github.com/AMDResearch/omniperf/issues>

Client-side installation (if required)

 Download the latest version from here: <https://github.com/AMDRResearch/omniperf/releases>

 Full documentation: <https://amdresearch.github.io/omniperf/>

```
wget https://github.com/AMDRResearch/omniperf/releases/download/v1.0.8-PR2/omniperf-v1.0.8-PR2.tar.gz

tar zxvf omniperf-v1.0.8-PR2.tar.gz

cd omniperf-v1.0.8-PR2/
python3 -m pip install -t ${INSTALL_DIR}/python-libs -r requirements.txt
mkdir build
cd build
export PYTHONPATH=${INSTALL_DIR}/python-libs:${PYTHONPATH}
cmake -DCMAKE_INSTALL_PREFIX=${INSTALL_DIR}/1.0.8 \
      -DPYTHON_DEPS=${INSTALL_DIR}/python-libs \
      -DMOD_INSTALL_PATH=${INSTALL_DIR}/modulefiles ..
make install
export PATH=${INSTALL_DIR}/1.0.8/bin:$PATH
```

Omniperf modes

Profile	Target application is launched using AMD ROC-profiler		
	Kernels	Dispatches	IP Blocks
Analyze	Profiled data is loaded to omniperf CLI		
	Immediate access to metrics	Lightweight standalone GUI	
Database	Profiled data is imported to Grafana™ database		
	Grafana™ GUI is based on MongoDB	Interact with saved workload database	

Basic command-line syntax:

Profile:

```
$ omniperf profile -n workload_name [profile options]
                    [roofline options] -- <CMD> <ARGS>
```

Analyze:

```
$ omniperf analyze -p
<path/to/workloads/workload_name/mi200/>
```

To use a lightweight standalone GUI with CLI analyzer:

```
$ omniperf analyze -p
<path/to/workloads/workload_name/mi200/> --gui
```

Database:

```
$ omniperf database <interaction type> [connection options]
```

For more information or help use -h/--help/? flags:

```
$ omniperf profile --help
```

For problems, create an issue here: <https://github.com/AMDResearch/omniperf/issues>

Documentation: <https://amdresearch.github.io/omniperf>

Omniperf profiling

We use the example sample/vcopy.cpp from the Omniperf installation folder:

```
$ wget https://github.com/AMDRResearch/omniperf/raw/main/sample/vcopy.cpp
```

Compile with hipcc:

```
$ hipcc -o vcopy vcopy.cpp
```

Profile with Omniperf:

```
$ omniperf profile -n vcopy_all -- ./vcopy -n 1048576 -b 256
```

```
...
```

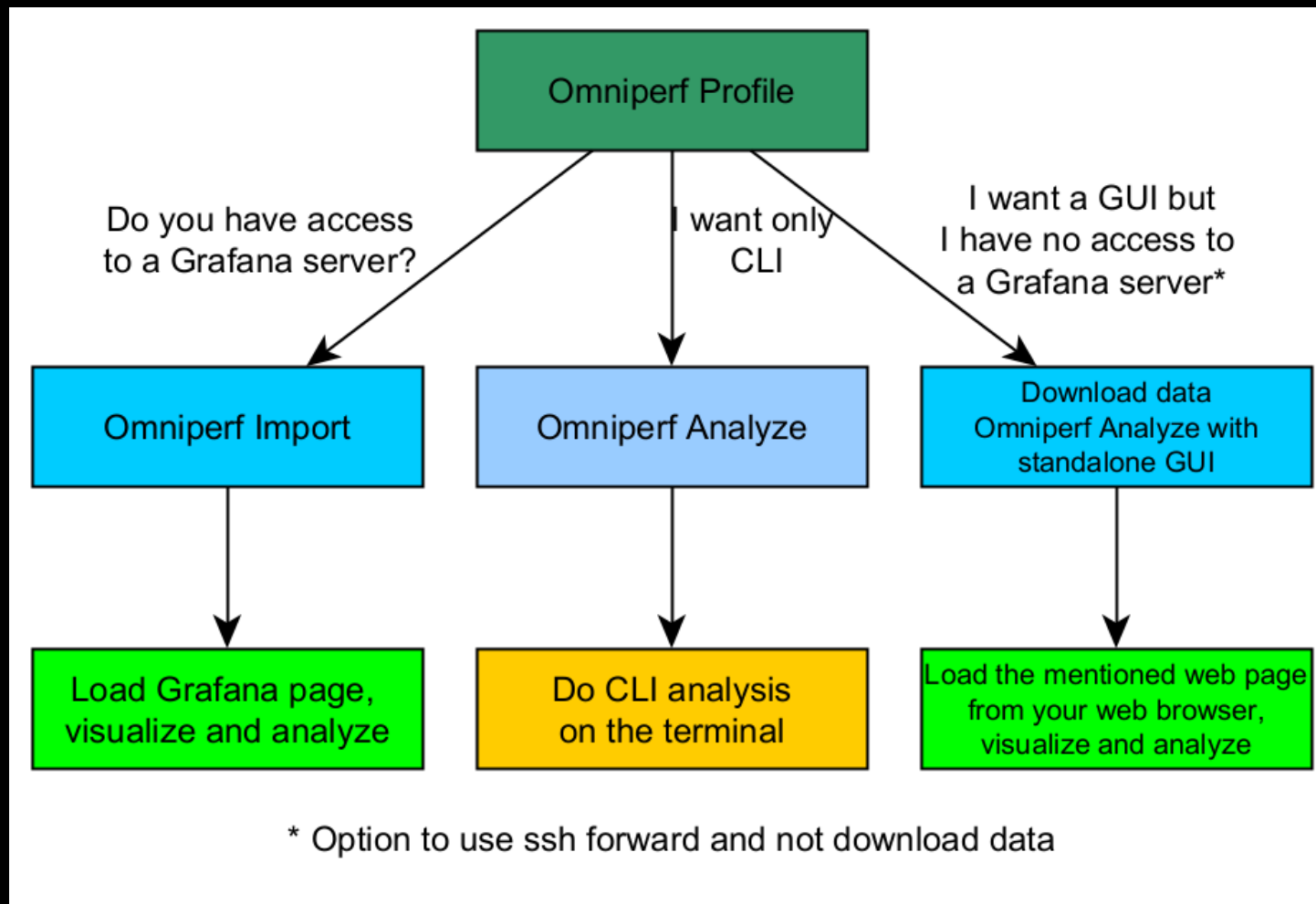
```
-----  
Profile only  
-----
```

```
omniperf ver: 1.0.4  
Path: /pfs/lustrep4/scratch/project_462000075/markoman/omniperf-  
1.0.4/build/workloads  
Target: mi200  
Command: ./vcopy 1048576 256  
Kernel Selection: None  
Dispatch Selection: None  
IP Blocks: All
```

A new directory will be created called workloads/vcopy_all

Note: Omniperf executes the code as many times as required to collect all HW metrics. Use kernel/dispatch filters especially when trying to collect roofline analysis.

Omniperf workflows



Omniperf analyze

We use the example sample/vcopy.cpp from the Omniperf installation folder:

```
$ wget https://github.com/AMDRResearch/omniperf/raw/main/sample/vcopy.cpp
```

Compile with hipcc:

```
$ hipcc --offload-arch=gfx90a -o vcopy vcopy.cpp
```

Profile with Omniperf:

```
$ omniperf profile -n vcopy_all -- ./vcopy -n 1048576 -b 256
```

A new directory will be created called workloads/vcopy_all

Analyze the profiled workload:

```
$ omniperf analyze -p workloads/vcopy_all/mi200/ &> vcopy_analyze.txt
```

0. Top Stat

	KernelName	Count	Sum(ns)	Mean(ns)	Median(ns)	Pc
0	vecCopy(double*, double*, double*, int, int) [clone .kd]	1	341123.00	341123.00	341123.00	100.00

2. System Speed-of-Light

Index	Metric	Value	Unit	Peak	PoP
2.1.0	VALU FLOPs	0.00	Gflop	23936.0	0.0
2.1.1	VALU IOPs	89.14	Giop	23936.0	0.37242200388114116
2.1.2	MFMA FLOPs (BF16)	0.00	Gflop	95744.0	0.0
2.1.3	MFMA FLOPs (F16)	0.00	Gflop	191488.0	0.0
2.1.4	MFMA FLOPs (F32)	0.00	Gflop	47872.0	0.0
2.1.5	MFMA FLOPs (F64)	0.00	Gflop	47872.0	0.0
2.1.6	MFMA IOPs (Int8)	0.00	Giop	191488.0	0.0
2.1.7	Active CUs	58.00	Cus	110	52.72727272727273
2.1.8	SALU Util	3.69	Pct	100	3.6862586934167525
2.1.9	VALU Util	5.90	Pct	100	5.895531580380328
2.1.10	MFMA Util	0.00	Pct	100	0.0
2.1.11	VALU Active Threads/Wave	32.71	Threads	64	51.10526315789473
2.1.12	IPC = Iops	0.08	Instr/cycle	5	10.576640821020212

7.1 Wavefront Launch Stats

Index	Metric	Avg	Min	Max	Unit
7.1.0	Grid Size	1048576.00	1048576.00	1048576.00	Work items
7.1.1	Workgroup Size	256.00	256.00	256.00	Work items
7.1.2	Total Wavefronts	16384.00	16384.00	16384.00	Wavefronts
7.1.3	Saved Wavefronts	0.00	0.00	0.00	Wavefronts
7.1.4	Restored Wavefronts	0.00	0.00	0.00	Wavefronts
7.1.5	VGPRs	44.00	44.00	44.00	Registers
7.1.6	SGPRs	48.00	48.00	48.00	Registers
7.1.7	LDS Allocation	0.00	0.00	0.00	Bytes
7.1.8	Scratch Allocation	16496.00	16496.00	16496.00	Bytes

Omniperf Analyze

- Execute omniperf analyze -h to see various options
- Use specific IP block (-b)

Top kernels:

```
$ srun -n 1 --gpus 1 omniperf analyze -p workloads/vcopy_all/mi200/ -b 0
```

IP Block of wavefronts

```
$ srun -n 1 --gpus 1 omniperf analyze -p workloads/vcopy_all/mi200/ -b 7.1.2
```

0. Top Stat

	KernelName	Count	Sum(ns)	Mean(ns)	Median(ns)	Pct
0	vecCopy(double*, double*, double*, int, int) [clone .kd]	1	20960.00	20960.00	20960.00	100.00

7. Wavefront

7.1 Wavefront Launch Stats

Index	Metric	Avg	Min	Max	Unit
7.1.2	Total Wavefronts	16384.00	16384.00	16384.00	Wavefronts

Omniperf analyze

To see available options and usage instructions:

```
$ omniperf analyze -h
...
Help:
  -h, --help            show this help message and exit

General Options:
  -v, --version          show program's version number and exit
  -V, --verbose          Increase output verbosity

Analyze Options:
  -p [ ...], --path [ ...]  Specify the raw data root dirs or desired results directory.
  -o, --output            Specify the output file.
  --list-kernels          List kernels. Top 10 kernels sorted by duration (descending order).
  --list-metrics          List metrics can be customized to analyze on specific arch:
                          gfx906
                          gfx908
                          gfx90a
  -b [ ...], --metric [ ...] Specify IP block/metric id(s) from --list-metrics for filtering.
  -k [ ...], --kernel [ ...] Specify kernel id(s) from --list-kernels for filtering.
  --dispatch [ ...]       Specify dispatch id(s) for filtering.
  --gpu-id [ ...]         Specify GPU id(s) for filtering.
  -n, --normal-unit       Specify the normalization unit: (DEFAULT: per_wave)
                          per_wave
                          per_cycle
                          per_second
                          per_kernel
  --config-dir            Specify the directory of customized configs.
  -t, --time-unit         Specify display time unit in kernel top stats: (DEFAULT: ns)
                          s
                          ms
                          us
                          ns
  --decimal               Specify the decimal to display. (DEFAULT: 2)
  --cols [ ...]           Specify column indices to display.
  -g                       Debug single metric.
  --dependency            List the installation dependency.
  --gui [GUI]             Activate a GUI to interate with Omniperf metrics.
                          Optionally, specify port to launch application (DEFAULT: 8050)
```

Easy things you can check

- Are all the CUs being used?
 - If not, more parallelism is required (for most of the cases)
- Are all the VGPRs being spilled?
 - Try smaller workgroup sizes
- Is the code Integer limited?
 - Try reducing the integer ops, usually in the index calculation

Omniperf analyze with standalone GUI

We use the example sample/vcopy.cpp from the Omniperf installation folder:

```
$ wget https://github.com/AMDRResearch/omniperf/raw/main/sample/vcopy.cpp
```

Compile with hipcc:

```
$ hipcc --offload-arch=gfx90a -o vcopy vcopy.cpp
```

Profile with Omniperf:

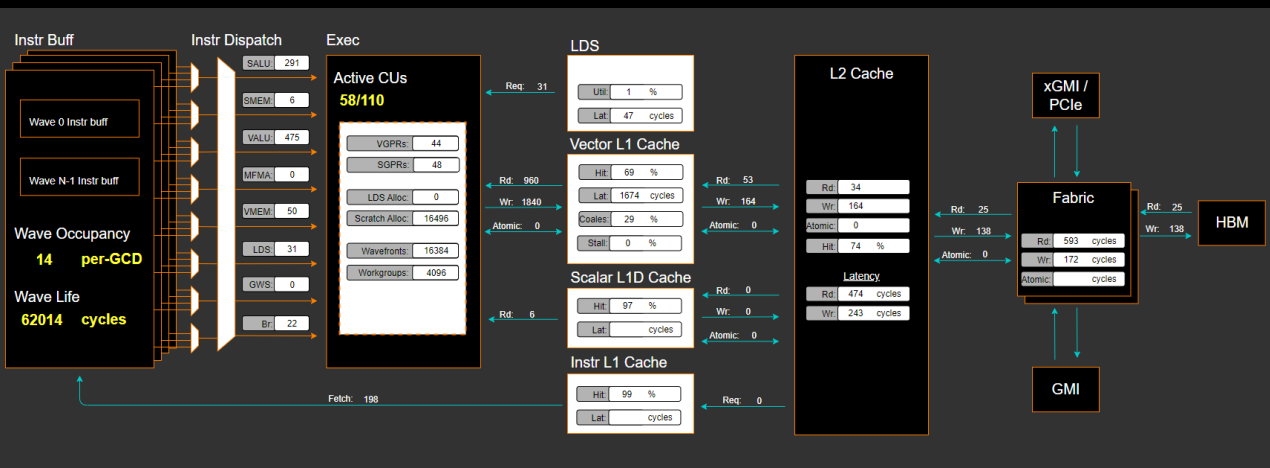
```
$ omniperf profile -n vcopy_all -- ./vcopy 1048576 256
```

A new directory will be created called workloads/vcopy_all

Analyze the profiled workload:

```
$ omniperf analyze -p workloads/vcopy_all/mi200/ --gui
```

Open web page <http://IP:8050/>

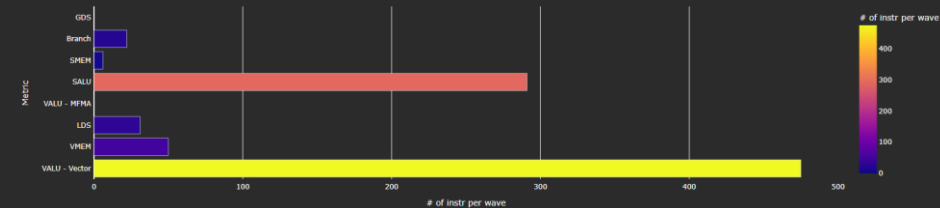


2. System Speed-of-Light

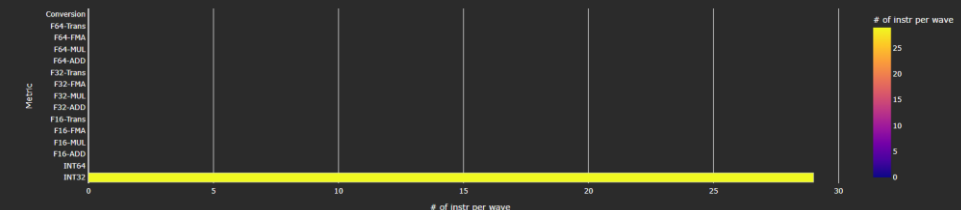
Metric	Value	Unit	Peak	Pop
VALU FLOPs	0.00	Gflop	23936.00	0.00
VALU TOPs	89.14	Gflop	23936.00	0.37
MFMA FLOPs (BF16)	0.00	Gflop	95744.00	0.00
MFMA FLOPs (F16)	0.00	Gflop	191488.00	0.00
MFMA FLOPs (F32)	0.00	Gflop	47872.00	0.00
MFMA FLOPs (F64)	0.00	Gflop	47872.00	0.00
MFMA TOPs (Int8)	0.00	Gflop	191488.00	0.00
Active CUs	58.00	Cus	110.00	52.73

10. Compute Units - Instruction Mix

10.1 Instruction Mix



10.2 VALU Arithmetic Instr Mix



Omniperf analyze with Grafana™ GUI

We use the example sample/vcopy.cpp from the Omniperf installation folder:

```
$ wget https://github.com/AMDRResearch/omniperf/raw/main/sample/vcopy.cpp
```

Compile with hipcc:

```
$ hipcc --offload-arch=gfx90a -o vcopy vcopy.cpp
```

Profile with Omniperf:

```
$ omniperf profile -n vcopy_all -- ./vcopy -n 1048576 -b 256
```

A new directory will be created called workloads/vcopy_all

Import the database to analyze in Grafana™ GUI:

```
$ omniperf database --import [connection options] -w workloads/vcopy_demo/mi200/  
ROC Profiler: /usr/bin/rocprof
```

```
-----  
Import Profiling Results  
-----
```

```
Pulling data from /root/test/workloads/vcopy_demo/mi200
```

```
The directory exists
```

```
Found sysinfo file
```

```
KernelName shortening enabled
```

```
Kernel name verbose level: 2
```

```
Password:
```

```
Password recieved
```

```
-- Conversion & Upload in Progress --
```

```
... ..
```

```
9 collections added.
```

```
Workload name uploaded
```

```
-- Complete! --
```

General / MIPerf_v1.0

Normalization: per Wave | Workload: miperf_asw_vcopy_demo_mi200 | Dispatch Filter: Enter variable value | GCD: 0 | Kernels: All

Baseline Workload: miperf_asw_vcopy_mi200 | Baseline Dispatch Filter: Enter variable value | Baseline GCD: 0 | Baseline Kernels: All

TopN: 5

> System Info (1 panel)

System Speed-of-Light

Metric	Speed of Light		Theoretical Max	Pct-of-Peak
	Avg	Unit		
VALU FLOPs	0	GFLOP	23,936	0%
VALU IOPs	379	GIOP	23,936	2%
MFMA FLOPs (BF16)	0	GFLOP	95,744	0%
MFMA FLOPs (F16)	0	GFLOP	191,488	0%
MFMA FLOPs (F32)	0	GFLOP	47,872	0%
MFMA FLOPs (F64)	0	GFLOP	47,872	0%
MFMA IOPs (Int8)	0	GIOP	191,488	0%
Active CUs	75	CUs	110	68%
SALU Util	4	pct	100	4%
VALU Util	6	pct	100	6%
MFMA Util	0	pct	100	0%
VALU Active Threads/Wave	64	Threads	64	100%
IPC - Issue	1	Instr/cycle	5	20%
LDS BW	0	GB/sec	23,936	0%
LDS Bank Conflict		Conflicts/access	32	
Instr Cache Hit Rate	100	pct	100	100%
Instr Cache BW	217	GB/s	6,093	4%
Scalar L1D Cache Hit Rate	100	pct	100	100%
Scalar L1D Cache BW	217	GB/s	6,093	4%
Vector L1D Cache Hit Rate	50	pct	100	50%
Vector L1D Cache BW	1,733	GB/s	11,968	14%
L2 Cache Hit Rate	36	pct	100	36%
L2-Fabric Read BW	434	GB/s	1,638	26%
L2-Fabric Write BW	301	GB/s	1,638	18%



Key Insights from Omnipperf Analyzer

Initial assessment with kernel statistics

Initial Assessment

Instruction/data flow

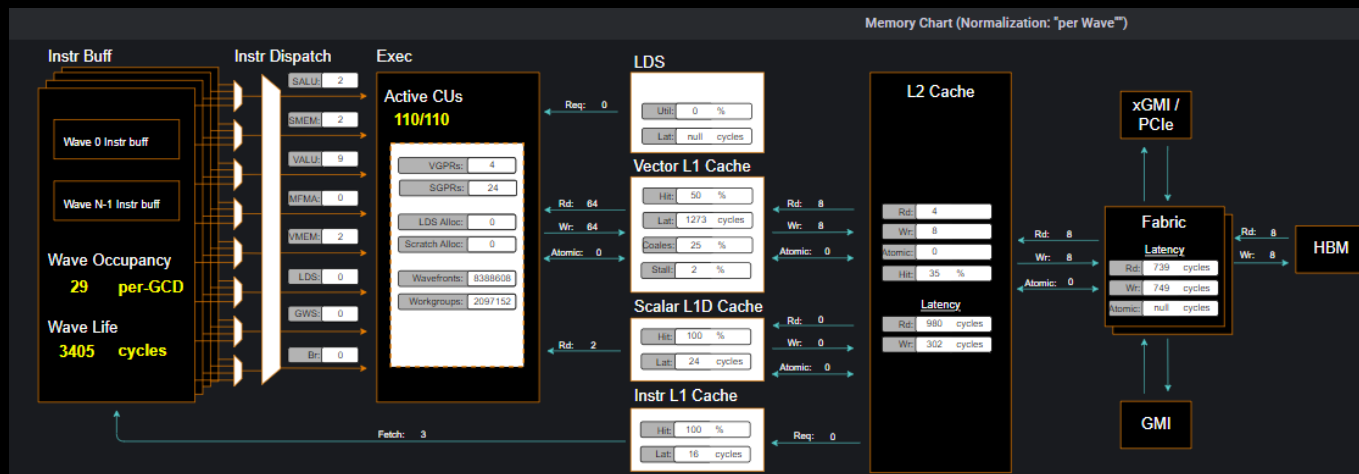
Speed-of-Light (SOL)

Omniperf tooling support

System SOL

Memory Chart

Kernel statistics



Metric	Avg	Unit	Theoretical Max	Pct-of-Peak
VALU FLOPs	0	GFLOP	23,936	0%
VALU IOPs	433	GIOP	23,936	2%
MFMA FLOPs (BF16)	0	GFLOP	95,744	0%
MFMA FLOPs (F16)	0	GFLOP	191,488	0%
MFMA FLOPs (F32)	0	GFLOP	47,872	0%
MFMA FLOPs (F64)	0	GFLOP	47,872	0%
MFMA IOPs (int8)	0	GIOP	191,488	0%
Active CUs	110	CUs	110	100%
SALU Util	3	pct	100	3%
VALU Util	8	pct	100	8%
MFMA Util	0	pct	100	0%
VALU Active Threads/Wave	64	Threads	64	100%
IPC - Issue	1	Instr/cycle	5	20%
LDS BW	0	GB/sec	23,936	0%
LDS Bank Conflict		Conflicts/access	32	
Instr Cache Hit Rate	100	pct	100	100%

Roofline: the first-step characterization of workload performance

Workload characterization

Compute bound

Memory bound

L1/L2 cache access

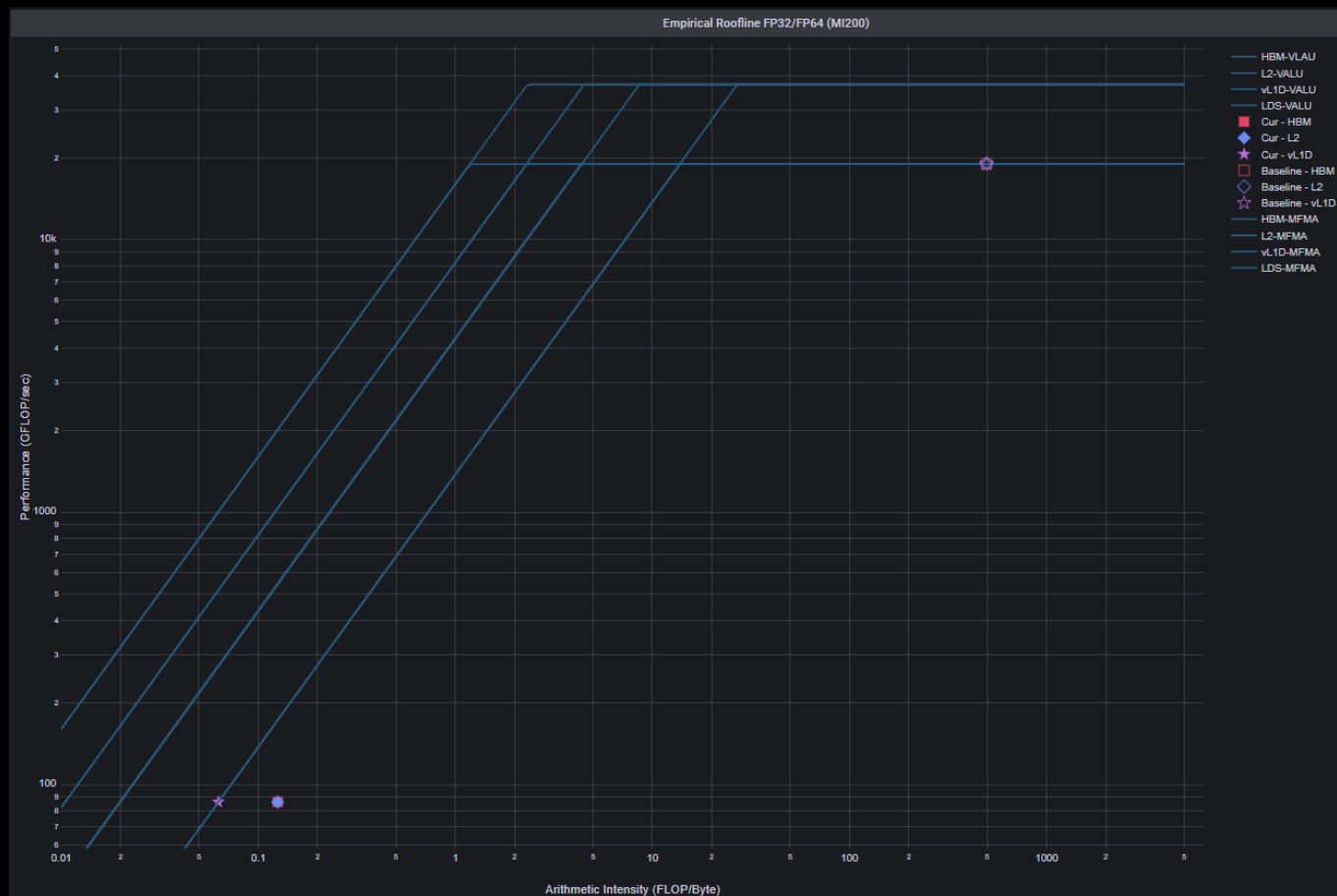
Performance margin

Omniperf tooling support

System SOL

Memory Chart

Kernel statistics



Top Kernels												
Name	Calls	Performance	HBM BW	Total Duration	Avg Duration	AI (Vector L1D Cache)	AI (L2 Cache)	AI (HBM)	Total FLOPs	VALU FLOPs	MFMA FLOPs (F16)	MFMA FLOPs (BF16)
void dot_kernel<doubl...	100	86.5 GFLOPS	689 GB/s	244 ms	2.44 ms	0.063	0.126	0.126	210,583,552	210,583,552	0	0
void triad_kernel<dou...	100	111 GFLOPS	1.33 TB/s	189 ms	1.89 ms	0.042	0.083	0.083	209,715,200	209,715,200	0	0
void add_kernel<doubl...	100	55.7 GFLOPS	1.34 TB/s	188 ms	1.88 ms	0.021	0.042	0.042	104,857,600	104,857,600	0	0
void copy_kernel<dou...	100	0 GFLOPS	1.37 TB/s	122 ms	1.22 ms	0	0	0	0	0	0	0
void mul_kernel<doubl...	100	86.1 GFLOPS	1.38 TB/s	122 ms	1.22 ms	0.031	0.063	0.063	104,857,600	104,857,600	0	0

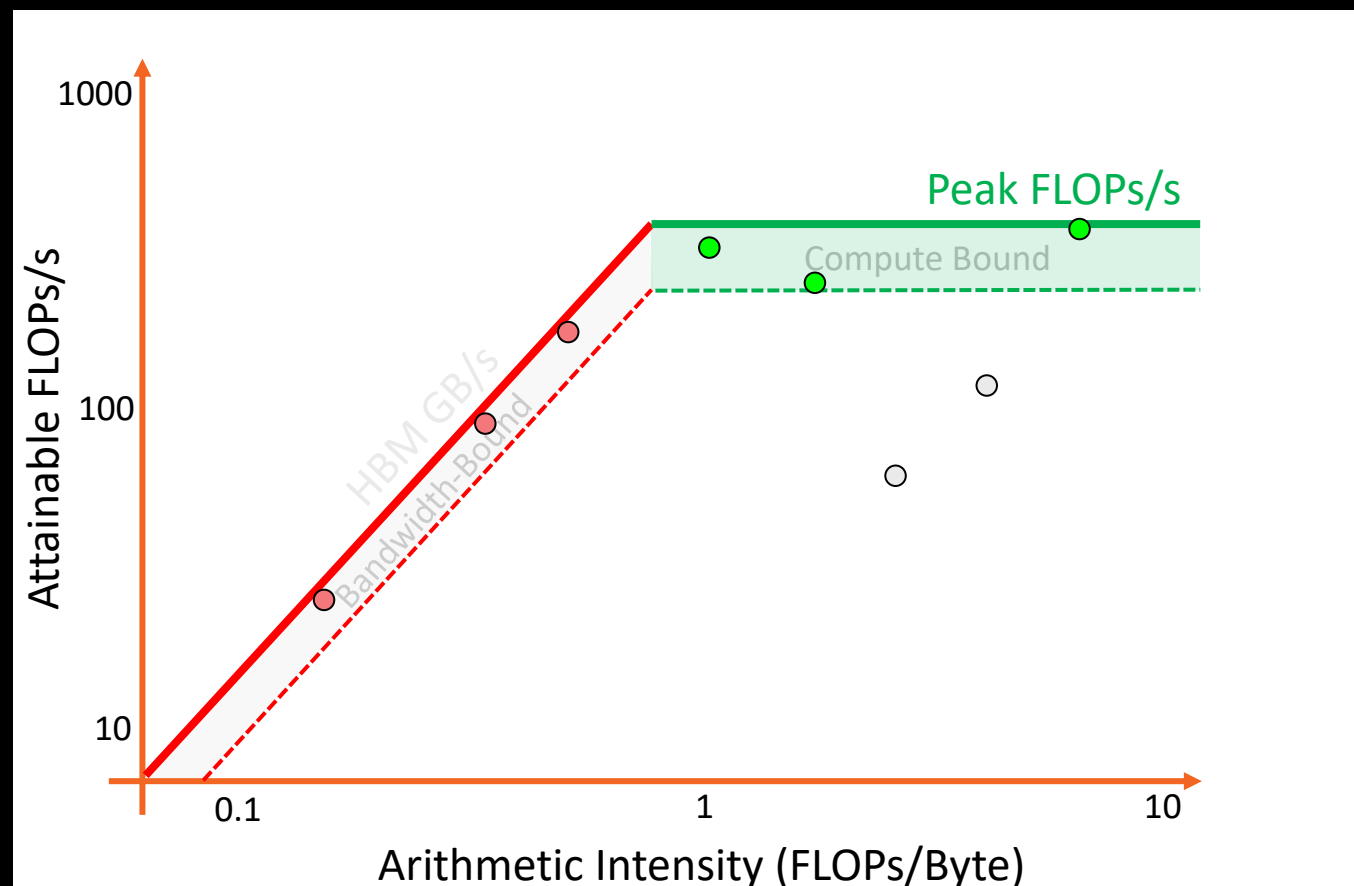


Background - What is a roofline?

Background – What is “Good” Performance?

Example:

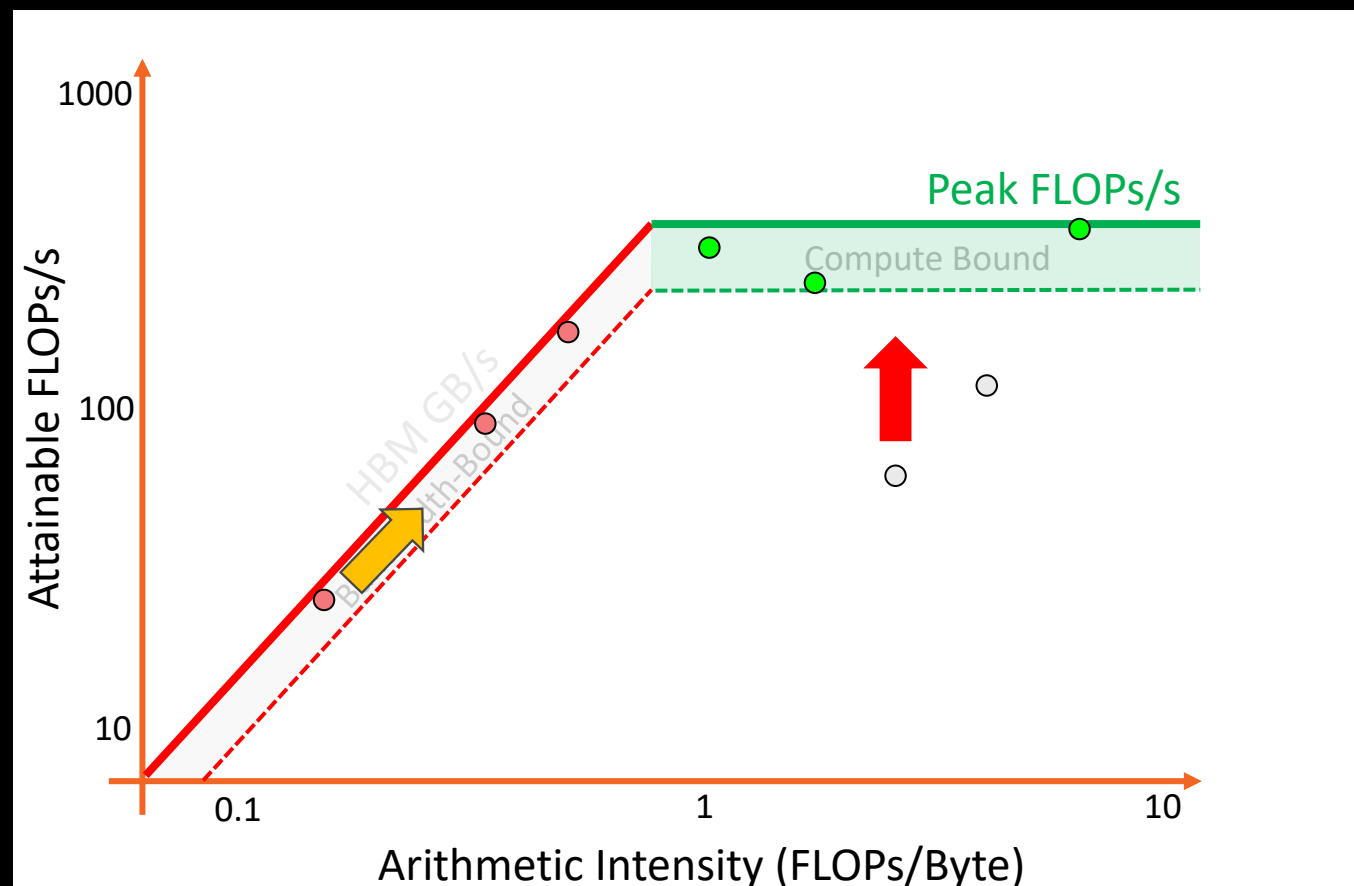
- We run a number of kernels and measure FLOPs/s
- Sort kernels by arithmetic intensity
- Compare performance relative to hardware capabilities
- Kernels near the roofline are making good use of computational resources
 - Kernels can have low performance (FLOPS/s), but make good use of BW



Background – What is “Good” Performance?

Example:

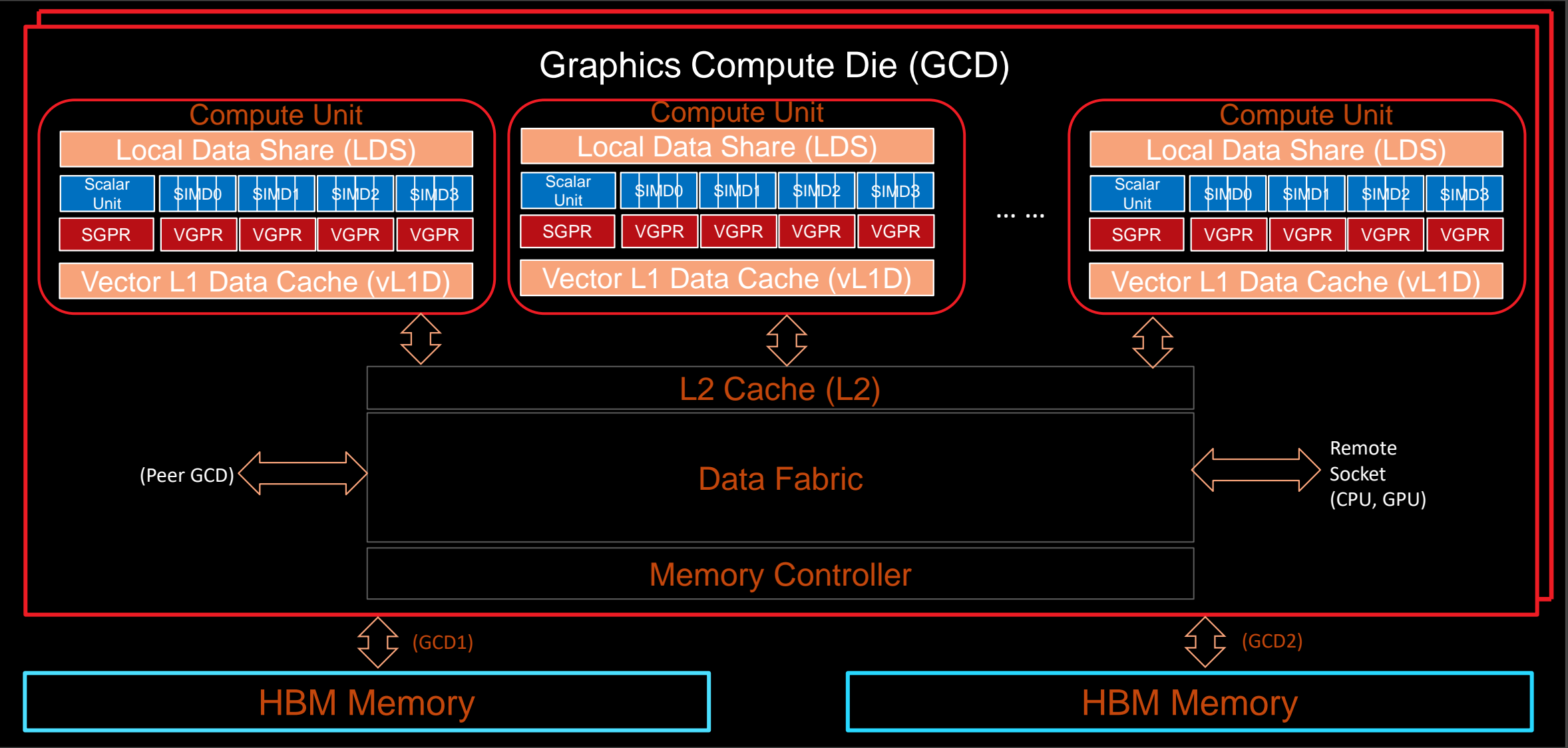
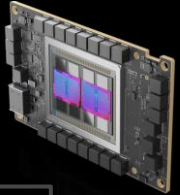
- We run a number of kernels and measure FLOPs/s
- Sort kernels by arithmetic intensity
- Compare performance relative to hardware capabilities
- Kernels near the roofline are making good use of computational resources
 - Kernels can have low performance (FLOPs/s), but make good use of BW
- Increase arithmetic intensity when bandwidth limited**
 - Reducing data movement increases AI
- Kernels not near the roofline *should** have optimizations that can be made to get closer to the roofline



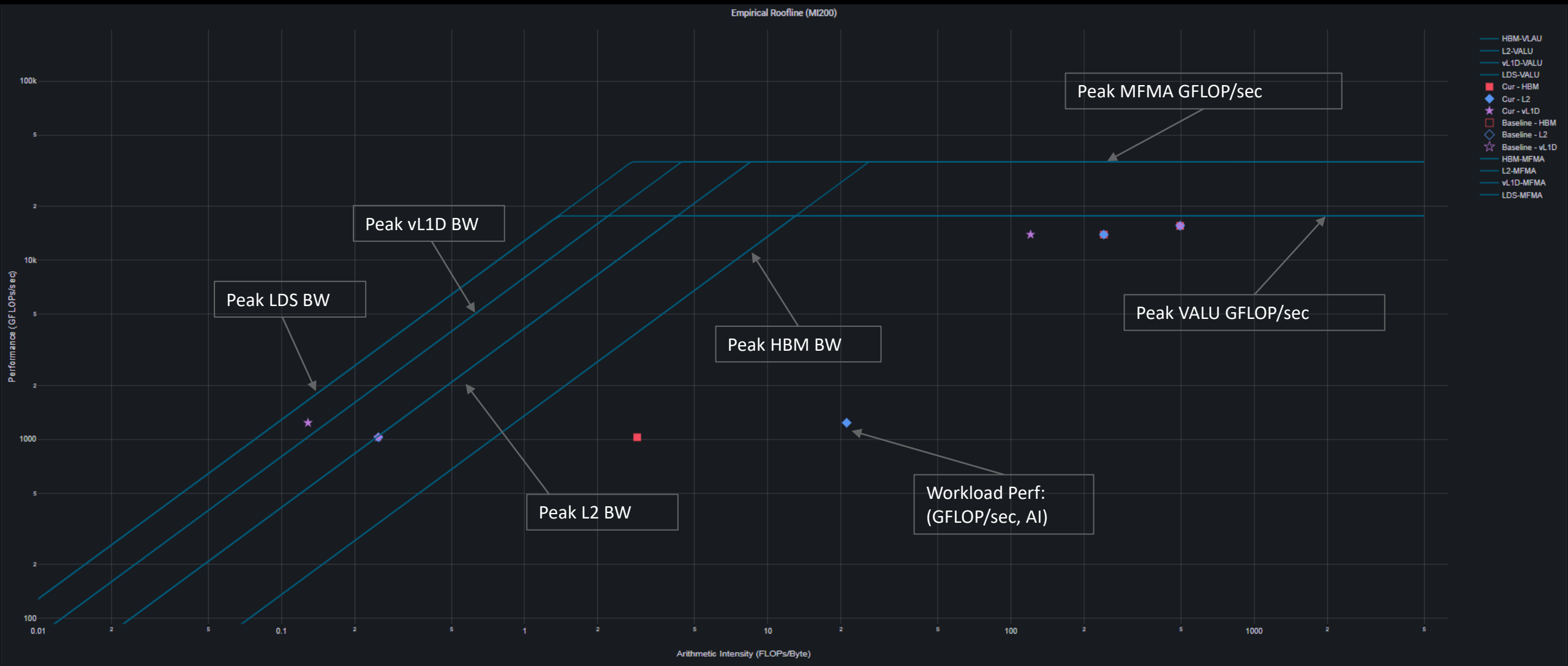


Roofline Calculations on AMD Instinct™ MI200 GPUs

Overview - AMD Instinct™ MI200 Architecture



Empirical Hierarchical Roofline on MI200 - Overview



Empirical Hierarchical Roofline on MI200 – Roofline Benchmarking

- Empirical Roofline Benchmarking
 - Measure achievable Peak FLOPS
 - VALU: F32, F64
 - MFMA: F16, BF16, F32, F64
 - Measure achievable Peak BW
 - LDS
 - Vector L1D Cache
 - L2 Cache
 - HBM
- Internally developed micro benchmark algorithms
 - Peak VALU FLOP: axpy
 - Peak MFMA FLOP: Matrix multiplication based on MFMA intrinsic
 - Peak LDS/vL1D/L2 BW: Pointer chasing
 - Peak HBM BW: Streaming copy

```

10:57:35 amd@node-bp126-014a'utils ±[master ✕]→ ./roofline
Total detected GPU devices: 2
GPU Device 0: Profiling...
99% [|||||]
HBM BW, GPU ID: 0, workgroupSize:256, workgroups:2097152, experiments:100, Total Bytes=8589934592, Duration=6.2 ms, Mean=1382.7 GB/sec, stdev=2.6 GB/s
99% [|||||]
L2 BW, GPU ID: 0, workgroupSize:256, workgroups:8192, experiments:100, Total Bytes=687194767360, Duration=157.3 ms, Mean=4321.3 GB/sec, stdev=59.1 GB/s
99% [|||||]
L1 BW, GPU ID: 0, workgroupSize:256, workgroups:16384, experiments:100, Total Bytes=26843545600, Duration=3.3 ms, Mean=8262.6 GB/sec, stdev=5.9 GB/s
99% [|||||]
LDS BW, GPU ID: 0, workgroupSize:256, workgroups:16384, experiments:100, Total Bytes=33554432000, Duration=1.8 ms, Mean=18780.4 GB/sec, stdev=33.0 GB/s
nSize:134217728, 268435456000
99% [|||||]
Peak FLOPs (FP32), GPU ID: 0, workgroupSize:256, workgroups:16384, experiments:100, Total FLOPS=274877906944, Duration=14.482 ms, Mean=18977.7 GFLOPs/sec, stdev=3.6 GFLOPs/s
99% [|||||]
Peak FLOPs (FP64), GPU ID: 0, workgroupSize:256, workgroups:16384, experiments:100, Total FLOPS=137438953472, Duration=7.5 ms, Mean=18336.156250.1 GFLOPs/sec, stdev=5.0 GFLOPs/s
99% [|||||]
Peak MFMA FLOPs (BF16), GPU ID: 0, workgroupSize:256, workgroups:16384, experiments:100, Total FLOPS=2147483648000, Duration=14.0 ms, Mean=153763.7 GFLOPs/sec, stdev=61.0 GFLOPs/s
99% [|||||]
Peak MFMA FLOPs (F16), GPU ID: 0, workgroupSize:256, workgroups:16384, experiments:100, Total FLOPS=2147483648000, Duration=14.5 ms, Mean=147890.9 GFLOPs/sec, stdev=32.2 GFLOPs/s
99% [|||||]
Peak MFMA FLOPs (F32), GPU ID: 0, workgroupSize:256, workgroups:16384, experiments:100, Total FLOPS=536870912000, Duration=14.4 ms, Mean=37200.4 GFLOPs/sec, stdev=9.3 GFLOPs/s
99% [|||||]
Peak MFMA FLOPs (F64), GPU ID: 0, workgroupSize:256, workgroups:16384, experiments:100, Total FLOPS=268435456000, Duration=7.3 ms, Mean=36978.4 GFLOPs/sec, stdev=10.0 GFLOPs/s

```

Empirical Hierarchical Roofline on MI200 - Arithmetic

$$\begin{aligned}
 \text{Total_FLOP} = & 64 * (\text{SQ_INSTS_VALU_ADD_F16} + \text{SQ_INSTS_VALU_MUL_F16} + \text{SQ_INSTS_VALU_TRANS_F16} + 2 * \text{SQ_INSTS_VALU_FMA_F16}) \\
 & + 64 * (\text{SQ_INSTS_VALU_ADD_F32} + \text{SQ_INSTS_VALU_MUL_F32} + \text{SQ_INSTS_VALU_TRANS_F32} + 2 * \text{SQ_INSTS_VALU_FMA_F32}) \\
 & + 64 * (\text{SQ_INSTS_VALU_ADD_F64} + \text{SQ_INSTS_VALU_MUL_F64} + \text{SQ_INSTS_VALU_TRANS_F64} + 2 * \text{SQ_INSTS_VALU_FMA_F64}) \\
 & + 512 * \text{SQ_INSTS_VALU_MFMA_MOPS_F16} \\
 & + 512 * \text{SQ_INSTS_VALU_MFMA_MOPS_BF16} \\
 & + 512 * \text{SQ_INSTS_VALU_MFMA_MOPS_F32} \\
 & + 512 * \text{SQ_INSTS_VALU_MFMA_MOPS_F64}
 \end{aligned}$$

$$\text{Total_IOP} = 64 * (\text{SQ_INSTS_VALU_INT32} + \text{SQ_INSTS_VALU_INT64})$$

$$AI_{LDS} = \frac{\text{TOTAL_FLOP}}{LDS_{BW}}$$

$$LDS_{BW} = 32 * 4 * (\text{SQ_LDS_IDX_ACTIVE} - \text{SQ_LDS_BANK_CONFLICT})$$

$$AI_{vL1D} = \frac{\text{TOTAL_FLOP}}{vL1D_{BW}}$$

$$vL1D_{BW} = 64 * \text{TCP_TOTAL_CACHE_ACCESSES_sum}$$

$$\begin{aligned}
 L2_{BW} = & 64 * \text{TCP_TCC_READ_REQ_sum} \\
 & + 64 * \text{TCP_TCC_WRITE_REQ_sum} \\
 & + 64 * (\text{TCP_TCC_ATOMIC_WITH_RET_REQ_sum} + \\
 & \text{TCP_TCC_ATOMIC_WITHOUT_RET_REQ_sum})
 \end{aligned}$$

$$AI_{L2} = \frac{\text{TOTAL_FLOP}}{L2_{BW}}$$

$$\begin{aligned}
 HBM_{BW} = & 32 * \text{TCC_EA_RDREQ_32B_sum} + 64 * (\text{TCC_EA_RDREQ_sum} - \\
 & \text{TCC_EA_RDREQ_32B_sum}) \\
 & + 32 * (\text{TCC_EA_WRREQ_sum} - \text{TCC_EA_WRREQ_64B_sum}) + 64 * \\
 & \text{TCC_EA_WRREQ_64B_sum}
 \end{aligned}$$

$$AI_{HBM} = \frac{\text{TOTAL_FLOP}}{HBM_{BW}}$$



* All calculations are subject to change

Empirical Hierarchical Roofline on MI200 - Manual Rocprof

- For those who like getting their hands dirty

- Generate input file

- See example roof-counters.txt →

- Run rocprof

```
foo@bar:~$ rocprof -i roof-counters.txt --timestamp on ./myCoolApp
```

- Analyze results

- Load *results.csv* output file in csv viewer of choice
 - Derive final metric values using equations on previous slide

- Profiling Overhead

- Requires one application replay for each pmc line

```
## roof-counters.txt

# FP32 FLOPs
pmc: SQ_INSTS_VALU_ADD_F32 SQ_INSTS_VALU_MUL_F32 SQ_INSTS_VALU_FMA_F32 SQ_INSTS_VALU_TRANS_F32

# HBM Bandwidth
pmc: TCC_EA_RDREQ_sum TCC_EA_RDREQ_32B_sum TCC_EA_WRREQ_sum TCC_EA_WRREQ_64B_sum

# LDS Bandwidth
pmc: SQ_LDS_IDX_ACTIVE SQ_LDS_BANK_CONFLICT

# L2 Bandwidth
pmc: TCP_TCC_READ_REQ_sum TCP_TCC_WRITE_REQ_sum TCP_TCC_ATOMIC_WITH_RET_REQ_sum
TCP_TCC_ATOMIC_WITHOUT_RET_REQ_sum

# vL1D Bandwidth
pmc: TCP_TOTAL_CACHE_ACSESSES_sum
```



Omniperf Performance Analyzer (cont..)

Subsystem performance analysis

Memory subsystems

L2 Cache

HBM access

LDS

vL1D

Omniperf tooling support

L2 Cache SOL

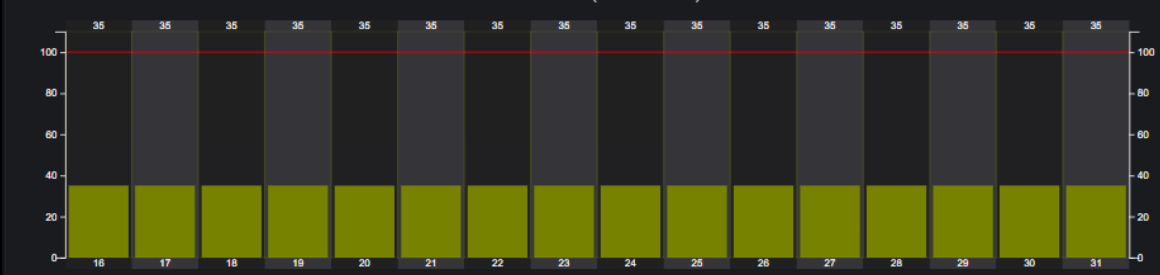
L2 fabric metrics

Per-channel statistics

Speed-of-Light: L2 Cache



Cache Hit Rate % (Channel 16 - 31)



L2 - Fabric Transactions

Metric	Avg	Min	Max	Unit
Read BW	693,148,700,953	664,565,016,054	695,197,543,698	Bytes per Sec
Write BW	692,659,558,092	664,096,634,666	694,705,946,653	Bytes per Sec
Read (32B)	0	0	0	Req per Sec
Read (Uncached 32B)	2,304,240	1,434,649	2,370,898	Req per Sec
Read (64B)	10,830,448,452	10,383,828,376	10,862,461,620	Req per Sec
HBM Read	10,830,362,679	10,383,764,324	10,862,381,992	Req per Sec
Write (32B)	0	0	0	Req per Sec
Write (Uncached 32B)	0	0	0	Req per Sec
Write (64B)	10,822,805,595	10,376,509,917	10,854,780,416	Req per Sec
HBM Write	10,822,801,389	10,376,488,102	10,854,762,613	Req per Sec
Read Latency	739	732	801	Cycles
Write Latency	749	737	784	Cycles
Atomic Latency				Cycles
Read Stall	3	2	3	pct
Write Stall	6	5	8	pct

L2 - Fabric Interface Stalls (Cycles "per Wave")



Shader compute components

Shader compute

Wavefront life

Instruction mix

Floating/
Integer Ops

Compute
pipeline

Instruction Mix



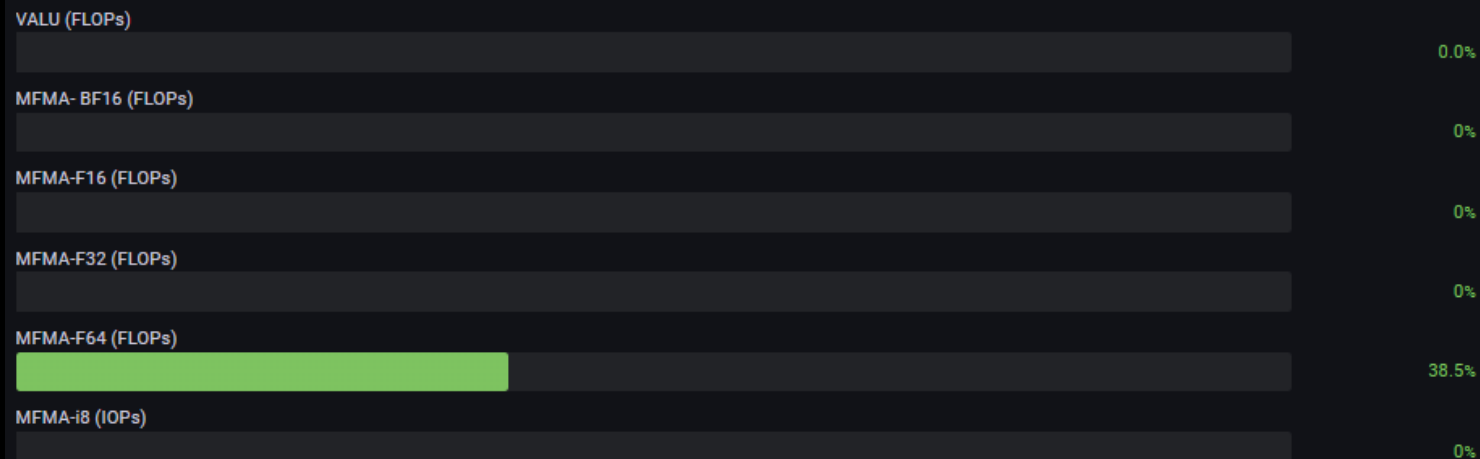
MFMA Arithmetic Instr Mix

MFMA Instr	Count
MFMA-i8	0
MFMA-F16	0
MFMA-BF16	0
MFMA-F32	0
MFMA-F64	995

Wavefront Runtime Stats

Metric	Avg	Min	Max	Unit
Kernel Time (Nanosec)	6,197,098	6,178,719	6,463,519	ns
Kernel Time (Cycles)	9,007,899	8,905,122	9,137,368	Cycle
Instr/wavefront	18	18	18	Instr/wavefro...
Wave Cycles	3,405	3,335	3,455	Cycles/wave
Dependency Wait Cycles	3,209	3,186	3,240	Cycles/wave
Issue Wait Cycles	165	112	193	Cycles/wave
Active Cycles	64	64	64	Cycles/wave
Wavefront Occupancy	3,198	3,166	3,210	Wavefronts

Speed-of-Light: Compute Pipeline



Omniperf profile – Roofline only

Profile with roofline:

```
$ omniperf profile -n roofline_case_app --roof-only -- <CMD> <ARGS>
```

Analyze the profiled workload:

```
$ omniperf analyze -p path/to/workloads/roofline_case_app/mi200 --gui
```

Open web page <http://IP:8050/>

When profile with --roof-only, a PDF with the roofline will be created. In order to see the name of the kernels, add the --kernel-names and a second PDF will be created with names for the kernel markers:

```
$ omniperf profile -n roofline_case_app --roof-only --kernel-names -- <CMD> <ARGS>
```

Empirical Roofline Analysis (FP32/FP64)



Roofline Analysis – Kokkos code

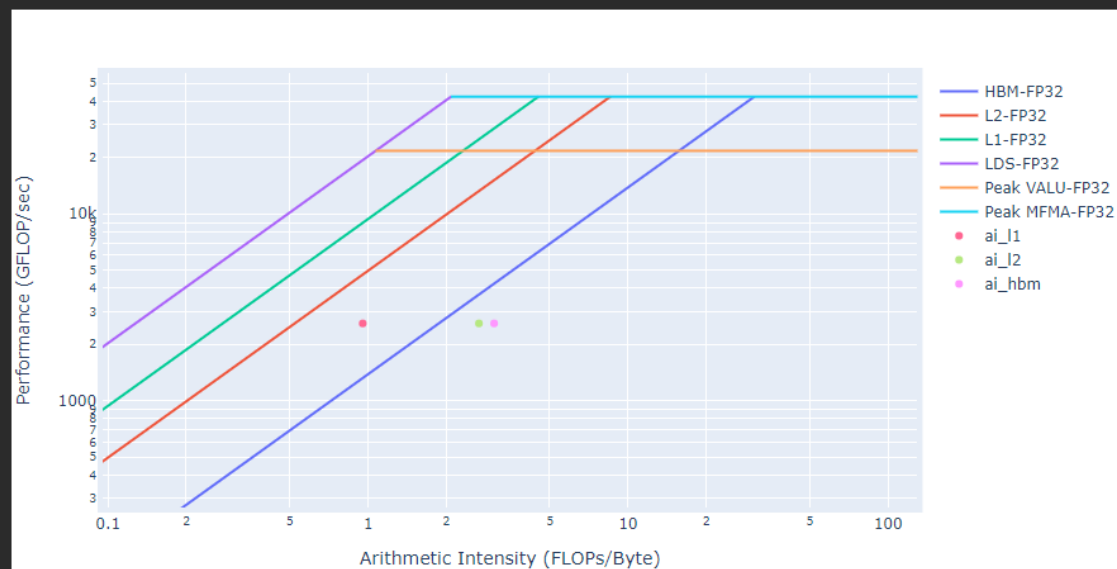
Menu ▾ NORMALIZATION: per Wave ▾ KERNELS: Fetch: 346 GCD: ALL ▾ DISPATCH FILTER: ALL ▾ Report Bug

```

void
Kokkos::Experimental::Impl::hip_parallel_launch_constant_memory<Kokkos::Impl::ParallelFor<idfix_for<Hydro::HlIdMHD<2>
()::{lambda(int, int, int)#1}>(std::_cxx11::basic_string<char, std::char_traits<char>, std::allocator<char>> const&, int const&,
int const&, int const&, int const&, int const&, int const&, Hydro::HlIdMHD<2>()::{lambda(int, int, int)#1}>::lambda(int
const&)#1), Kokkos::RangePolicy<Kokkos::Experimental::HIP>, Kokkos::Experimental::HIP>, 256u, 1u>() [clone .kd]

```

Empirical Roofline Analysis (FP32/FP64)



- Roofline: the first-step characterization of workload performance
 - Workload characterization
 - Compute bound
 - Memory bound
 - Performance margin
 - L1/L2 cache accesses
- Thorough SoC perf analysis for each subsystem to identify bottlenecks
 - HBM
 - L1/L2
 - LDS
 - Shader compute
 - Wavefront dispatch
- Omniperf tooling support
 - Roofline plot (float, integer)
 - Baseline roofline comparison
 - Kernel statistics

SPI Resource Allocation

- Dispatch Bound
 - Wavefront dispatching failure due to resources limitation
 - Wavefront slots
 - VGPR
 - SGPR
 - LDS allocation
 - Barriers
 - Etc.
 - Omniperf tooling support
 - Shader Processor Input (SPI) metrics

6.2 SPI Resource Allocation

Metric	Avg	Min	Max	Unit
Wave request Failed (CS)	613303.00	613303.00	613303.00	Cycles
CS Stall	356961.00	356961.00	356961.00	Cycles
CS Stall Rate	62.95	62.95	62.95	Pct
Scratch Stall	0.00	0.00	0.00	Cycles
Insufficient SIMD Waveslots	0.00	0.00	0.00	Simd
Insufficient SIMD VGPRs	16252333.00	16252333.00	16252333.00	Simd
Insufficient SIMD SGPRs	0.00	0.00	0.00	Simd
Insufficient CU LDS	0.00	0.00	0.00	Cu
Insufficient CU Barries	0.00	0.00	0.00	Cu
Insufficient Bulky Resource	0.00	0.00	0.00	Cu
Reach CU Threadgroups Limit	0.00	0.00	0.00	Cycles
Reach CU Wave Limit	0.00	0.00	0.00	Cycles
VGPR Writes	4.00	4.00	4.00	Cycles/wave
SGPR Writes	5.00	5.00	5.00	Cycles/wave



What if Grafana and web GUI crashes when loading performance data?
(real case)

When profiling produces too large data...

- We had an application that the realistic case was dispatching 6.7 million calls to kernels
- Executing Omniperf without any options, it would take up to 36 hours to finish while single non instrumented execution takes less than 1 hour.
- HW counters add overhead
- We had totally around 9 GB of profiling data from 1 MPI process
- Uploading the data to a Grafana server was crashing Grafana server and we had to reboot the service
- Using standalone GUI was never finishing loading the data

- Omniperf profile has an option called `-k` where you define which specific kernel to profile. You can define the id 0-9 of the top 10 kernels.
- This creates profiling data **only** for the selected kernel
- This way you can split the profiling data to 10 executions, one per kernel:
 - You can use different resources to do the experiments in parallel (remember there can be performance variation between different GPUs)
 - You can visualize each kernel

Profile with roofline for a specific kernel:

```
$ srun -N 1 -n 1 --ntasks-per-node=1 --gpus=1 --hint=nomultithread omniperf profile -n kernel_roof  
-k kernel_name --roof-only -- ./binary args
```



Example – DAXPY with a loop in the kernel

DAXPY – with a loop in the kernel

```
#include <hip/hip_runtime.h>

__constant__ double a = 1.0f;

__global__
void daxpy (int n, double const* x, int incx, double* y, int incy)
{
    int i = blockDim.x * blockIdx.x + threadIdx.x;
    if (i < n)
        for(int ll=0;ll<20;ll++) {
            y[i] = a*x[i] + y[i];
        }
}

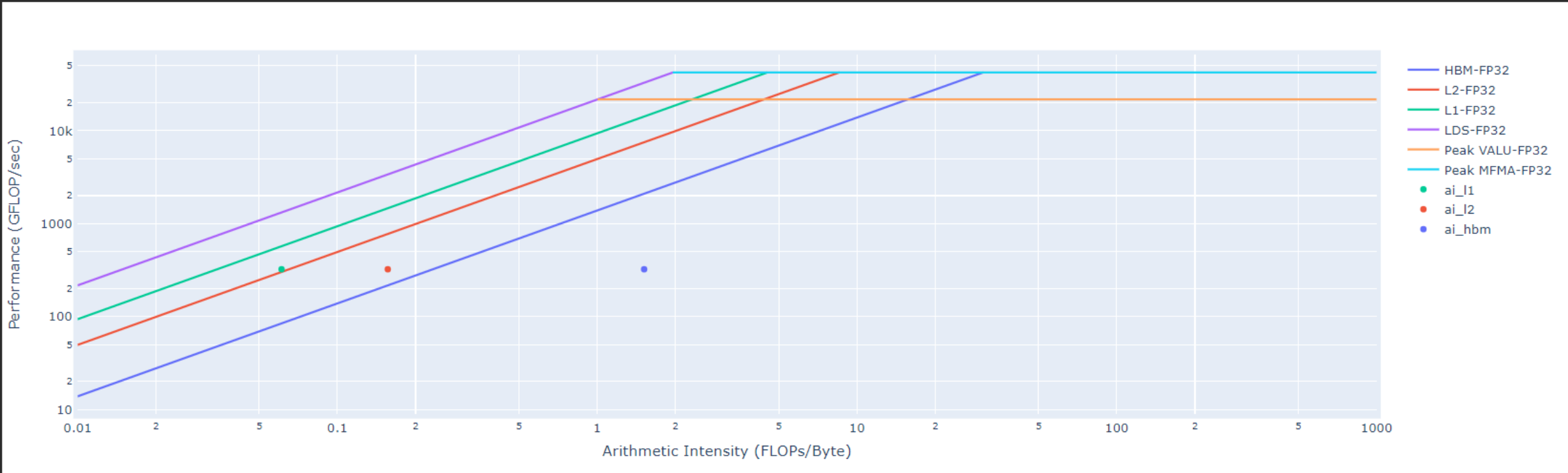
int main()
{
    int n = 1<<24;
    std::size_t size = sizeof(double)*n;

    double* d_x;
    double *d_y;
    hipMalloc(&d_x, size);
    hipMalloc(&d_y, size);

    int num_groups = (n+255)/256;
    int group_size = 256;
    daxpy<<<num_groups, group_size>>>(n, d_x, 1, d_y, 1);
    hipDeviceSynchronize();
}
```

Roofline

Empirical Roofline Analysis (FP32/FP64)



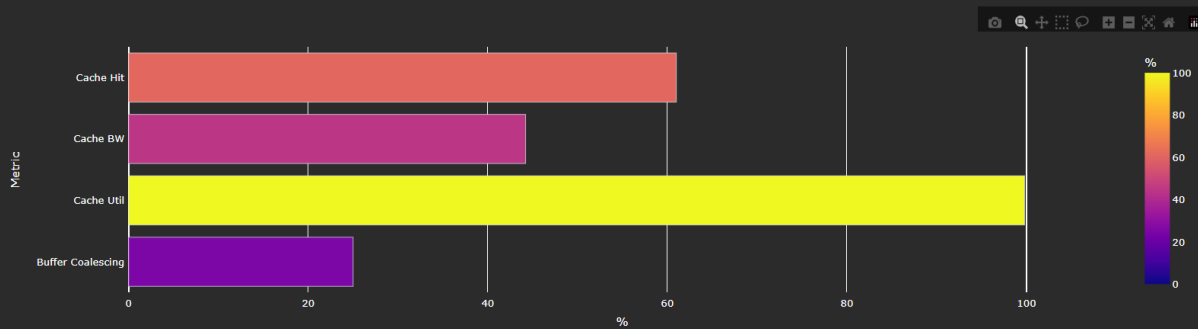
• Performance: almost 330 GFLOPs

Kernel execution time and L1D Cache Accesses

KernelName	Count	Sum(ns)	Mean(ns)	Median(ns)	Pct
daxpy(int, double const*, int, double*, int) [clone .kd]	1.00	2024491.00	2024491.00	2024491.00	100.00

16. Vector L1 Data Cache

16.1 Speed-of-Light



16.2 L1D Cache Stalls

Metric	Mean	Min	Max	Unit
Stalled on L2 Data	73.69	73.69	73.69	Pct
Stalled on L2 Req	19.47	19.47	19.47	Pct
Tag RAM Stall (Read)	0.00	0.00	0.00	Pct
Tag RAM Stall (Write)	0.00	0.00	0.00	Pct
Tag RAM Stall (Atomic)	0.00	0.00	0.00	Pct

16.3 L1D Cache Accesses

Metric	Avg	Min	Max	Unit
Total Req	2624.00	2624.00	2624.00	Req per wave
Read Req	1344.00	1344.00	1344.00	Req per wave
Write Req	1280.00	1280.00	1280.00	Req per wave
Atomic Req	0.00	0.00	0.00	Req per wave
Cache BW	5291.66	5291.66	5291.66	Gb/s
Cache Accesses	656.00	656.00	656.00	Req per wave
Cache Hits	400.16	400.16	400.16	Req per wave
Cache Hit Rate	61.00	61.00	61.00	Pct

DAXPY – with a loop in the kernel - Optimized

```
#include <hip/hip_runtime.h>

__constant__ double a = 1.0f;

__global__
void daxpy (int n, double const* __restrict__ x, int incx, double* __restrict__ y, int incy)
{
    int i = blockDim.x * blockIdx.x + threadIdx.x;
    if (i < n)
        for(int ll=0;ll<20;ll++) {
            y[i] = a*x[i] + y[i];
        }
}

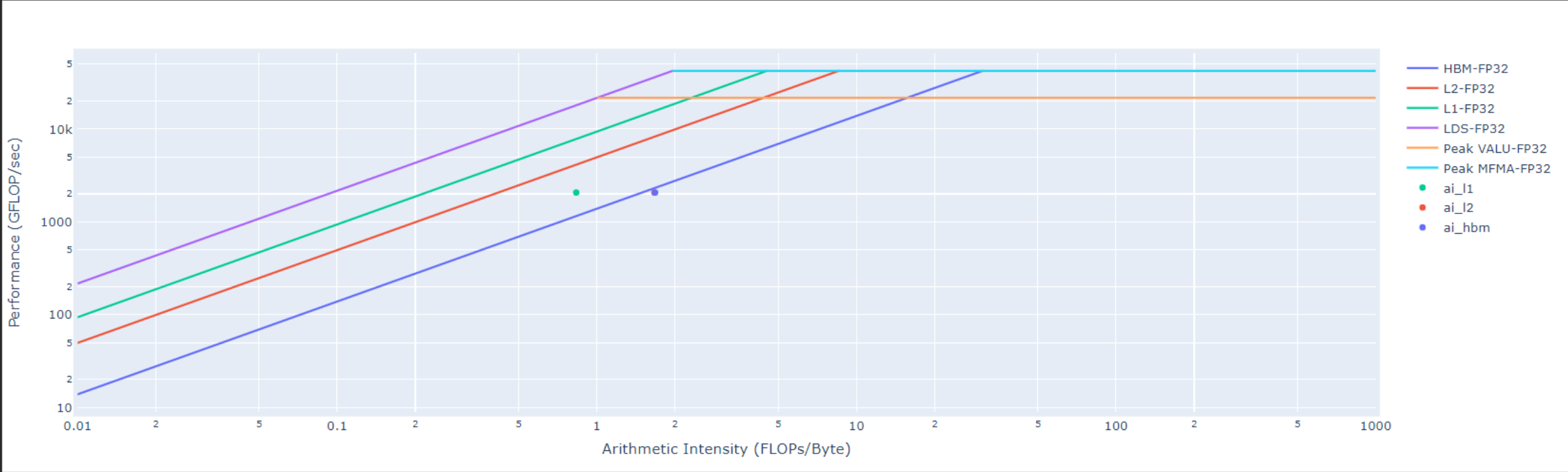
int main()
{
    int n = 1<<24;
    std::size_t size = sizeof(double)*n;

    double* d_x;
    double *d_y;
    hipMalloc(&d_x, size);
    hipMalloc(&d_y, size);

    int num_groups = (n+255)/256;
    int group_size = 256;
    daxpy<<<num_groups, group_size>>>(n, d_x, 1, d_y, 1);
    hipDeviceSynchronize();
}
```

Roofline - Optimized

Empirical Roofline Analysis (FP32/FP64)



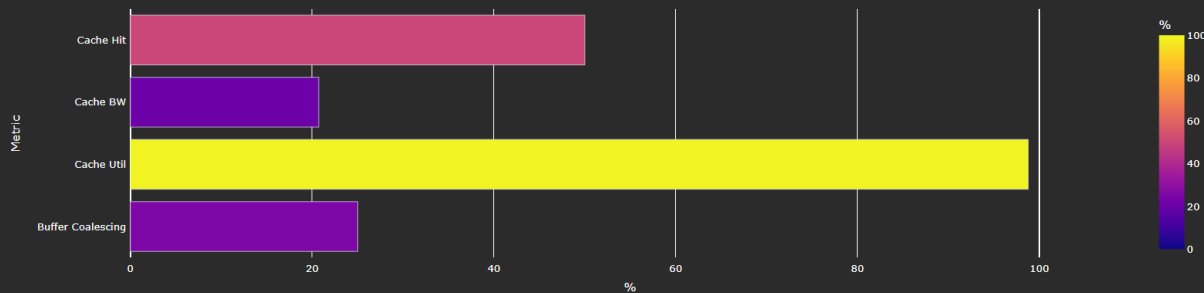
• Performance: almost 2 TFLOPs

Kernel execution time and L1D Cache Accesses - Optimized

KernelName	Count	Sum(ns)	Mean(ns)	Median(ns)	Pct
daxpy(int, double const*, int, double*, int) [clone .kd]	1.00	323522.00	323522.00	323522.00	100.00

6.2 times faster!

16.1 Speed-of-Light



16.2 L1D Cache Stalls

Metric	Mean	Min	Max	Unit
Stalled on L2 Data	79.08	79.08	79.08	Pct
Stalled on L2 Req	15.17	15.17	15.17	Pct
Tag RAM Stall (Read)	0.00	0.00	0.00	Pct
Tag RAM Stall (Write)	0.00	0.00	0.00	Pct
Tag RAM Stall (Atomic)	0.00	0.00	0.00	Pct

16.3 L1D Cache Accesses

Metric	Avg	Min	Max	Unit
Total Req	192.00	192.00	192.00	Req per wave
Read Req	128.00	128.00	128.00	Req per wave
Write Req	64.00	64.00	64.00	Req per wave
Atomic Req	0.00	0.00	0.00	Req per wave
Cache BW	2480.60	2480.60	2480.60	Gb/s
Cache Accesses	48.00	48.00	48.00	Req per wave
Cache Hits	24.00	24.00	24.00	Req per wave
Cache Hit Rate	50.00	50.00	50.00	Pct
Invalidate	0.00	0.00	0.00	Req per wave

Summary

- Omniperf is a tool that collects many counters automatically
- It can create roofline analysis to understand how efficient are your kernels
- It displays a lot of metrics regarding your kernels, however, it is required to know more about your kernel
- It does not have learning curve to start running it, but requies knowledge for the analysis
- It supports Grafana, standalone GUI, and CLI
- Includes several features such as:
 - System Speed-of-Light Panel
 - Memory Chart Analysis Panel
 - Vector L1D Cache Panel
 - Shader Processing Input (SPI) Panel

Questions?

DISCLAIMERS

The information presented in this document is for informational purposes only and may contain technical inaccuracies, omissions, and typographical errors. The information contained herein is subject to change and may be rendered inaccurate for many reasons, including but not limited to product and roadmap changes, component and motherboard version changes, new model and/or product releases, product differences between differing manufacturers, software changes, BIOS flashes, firmware upgrades, or the like. Any computer system has risks of security vulnerabilities that cannot be completely prevented or mitigated. AMD assumes no obligation to update or otherwise correct or revise this information. However, AMD reserves the right to revise this information and to make changes from time to time to the content hereof without obligation of AMD to notify any person of such revisions or changes.

THIS INFORMATION IS PROVIDED 'AS IS.' AMD MAKES NO REPRESENTATIONS OR WARRANTIES WITH RESPECT TO THE CONTENTS HEREOF AND ASSUMES NO RESPONSIBILITY FOR ANY INACCURACIES, ERRORS, OR OMISSIONS THAT MAY APPEAR IN THIS INFORMATION. AMD SPECIFICALLY DISCLAIMS ANY IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY, OR FITNESS FOR ANY PARTICULAR PURPOSE. IN NO EVENT WILL AMD BE LIABLE TO ANY PERSON FOR ANY RELIANCE, DIRECT, INDIRECT, SPECIAL, OR OTHER CONSEQUENTIAL DAMAGES ARISING FROM THE USE OF ANY INFORMATION CONTAINED HEREIN, EVEN IF AMD IS EXPRESSLY ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

Third-party content is licensed to you directly by the third party that owns the content and is not licensed to you by AMD. ALL LINKED THIRD-PARTY CONTENT IS PROVIDED "AS IS" WITHOUT A WARRANTY OF ANY KIND. USE OF SUCH THIRD-PARTY CONTENT IS DONE AT YOUR SOLE DISCRETION AND UNDER NO CIRCUMSTANCES WILL AMD BE LIABLE TO YOU FOR ANY THIRD-PARTY CONTENT. YOU ASSUME ALL RISK AND ARE SOLELY RESPONSIBLE FOR ANY DAMAGES THAT MAY ARISE FROM YOUR USE OF THIRD-PARTY CONTENT.

© 2023 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo, AMD CDNA, AMD ROCm, AMD Instinct, and combinations thereof are trademarks of Advanced Micro Devices, Inc. in the United States and/or other jurisdictions. Other names are for informational purposes only and may be trademarks of their respective owners.

ATTRIBUTIONS

Docker and the Docker logo are trademarks or registered trademarks of Docker, Inc.

Git and the Git logo are either registered trademarks or trademarks of Software Freedom Conservancy, Inc., corporate home of the Git Project, in the United States and/or other countries.

MongoDB is a registered trademark of MongoDB, Inc.

Python is a registered trademark of the PSF.

AMD 

Guided Exercises

1. Launch Parameters
2. LDS Occupancy Limiter
3. VGPR Occupancy Limiter
4. Strided Data Access Pattern
5. Algorithmic Optimizations
6. Daxpy example

Guided Exercises: Logistics/Preamble

- To accommodate the virtual setting and attendees with varied access to Omniperf:
 - I'll read through the slides without waiting for everyone to finish working through each exercise
 - If you have access to a system with Omniperf, clone the repo and start working through the exercises:
 - `git clone https://github.com/amd/HPCTrainingExamples/`
 - The READMEs contain all of what I'm saying and include platform-specific instructions for this training in the top-level directory
- We have used a publicly available release candidate of Omniperf to generate output for these slides:
 - <https://github.com/AMDResearch/omniperf/releases/tag/v1.1.0-PR1>
 - Behavior may differ if using a different version of Omniperf (e.g. 1.0.10)
 - Generally, building stable releases is the best practice
- The numbers shown in the READMEs and these slides were generated using MI210 accelerators
- Implementations in these exercises are **not** fully-optimized kernels

Guided Exercises: Representative Optimization Tasks

- The Exercises are roughly in order of ease of development effort and performance impact:
 - Exercise 1: Verify Reasonable Launch Parameters
 - Exercise 2: Attempt to Cache Data in Shared Memory
 - Exercise 3: Determining a Source of Unexpected Resource Usage
 - Exercise 4: Verifying Efficient Data Access Patterns
 - Exercise 5: Analyzing an Algorithmic Change
- The underlying code is kept simple to emphasize the optimization techniques
- These slides are intended as a “Cheat Sheet” starting point providing:
 - Omnipperf commands to filter through output for common optimization concerns
 - Some optimization direction given certain Omnipperf output

Guided Exercises: Optimizing a yAx Kernel

- We'll be looking at a relatively simple kernel that solves the same problem in each exercise, yAx
 - yAx is a vector-matrix-vector product that can be implemented in serial as:

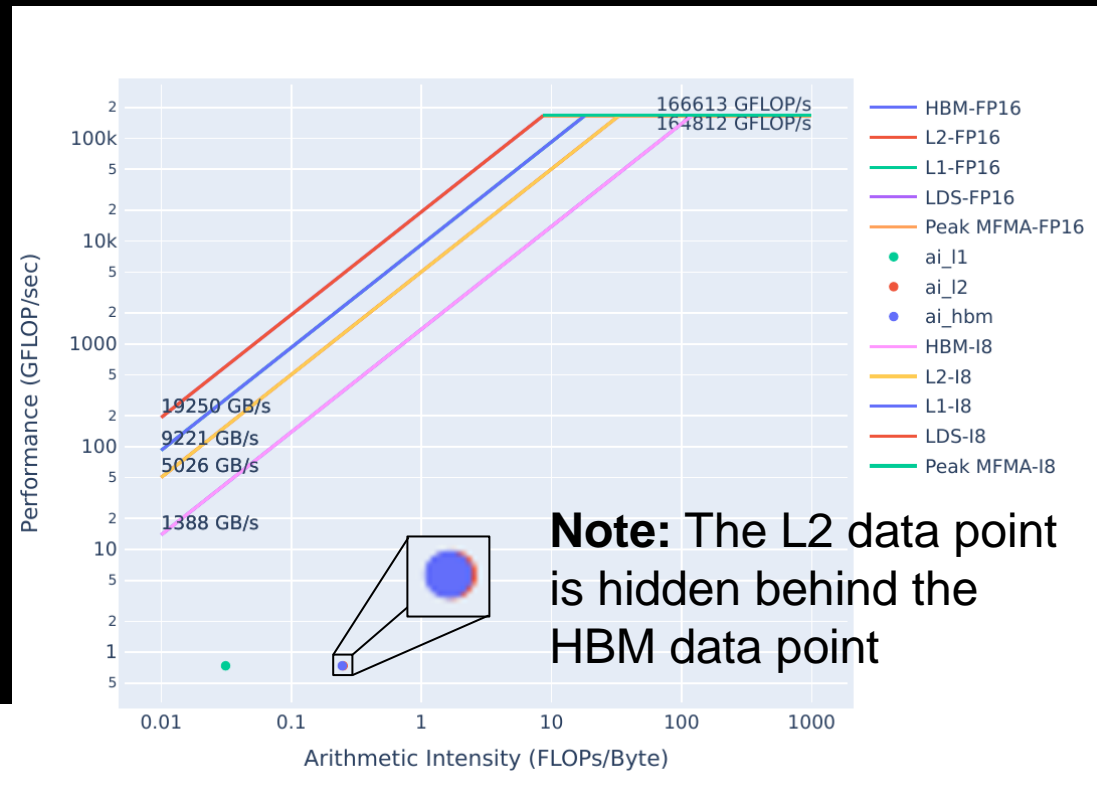
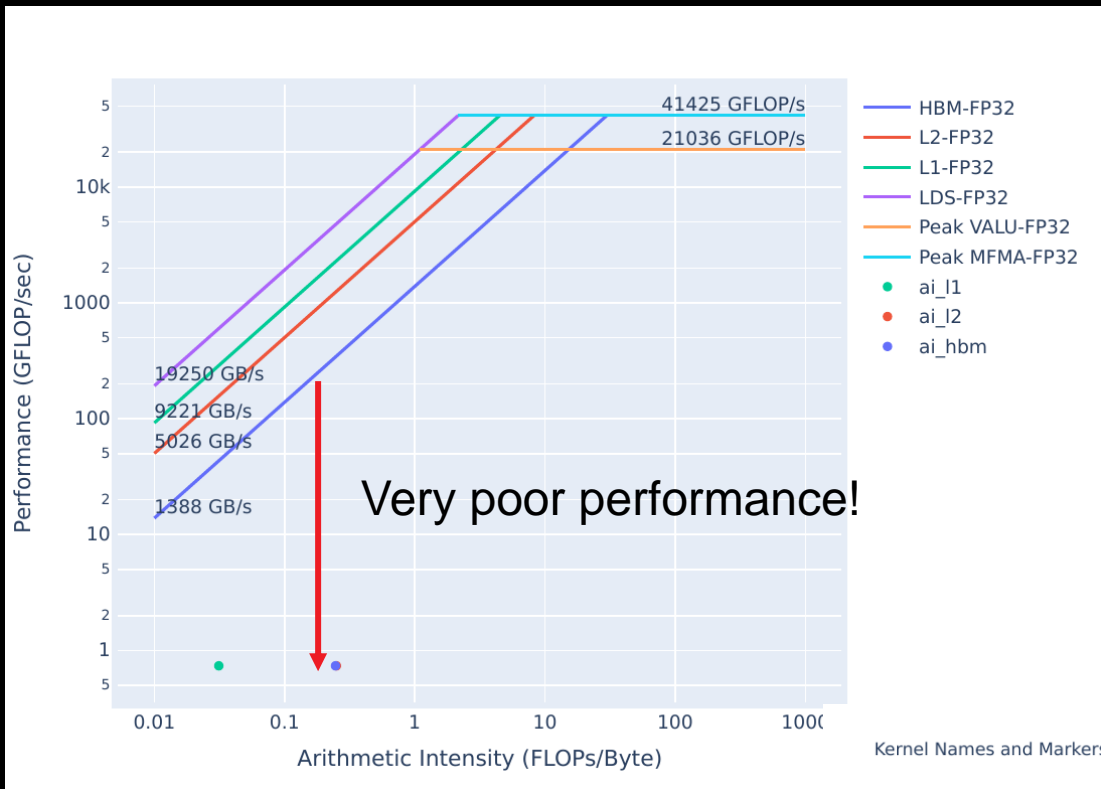
```
double result = 0.0;
for (int i = 0; i < n; i++){
    double temp = 0.0;
    for (int j = 0; j < m; j++){
        temp += A[i*m + j] * x[j];
    }
    result += y[i] * temp;
}
```

- Where:
 - A is a 1-D array of size n*m
 - x is an array of size m
 - y is an array of size n

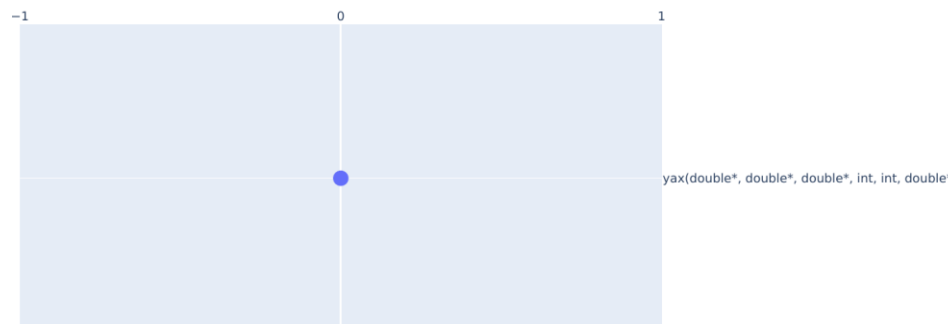
Exercise 1: First Things First, Generate a Roofline

- Run this command to generate roofline plots and a legend for each kernel (in PDF form):
 - `omniperf profile -n problem_roof_only --roof-only --kernel-names -- ./problem.exe`
 - The files will appear in the `./workloads/problem_roof_only/mi200` folder.
 - `--roof-only` generates PDF roofline plots, and does **not** generate any non-roofline profiling data
 - `--kernel-names` generates a PDF showing which kernel names correspond to which icons in the roofline
- Rooflines are a useful tool in determining which kernels are good optimization targets
 - They are only one perspective of performance: runtime of the kernel cannot be inferred from the roofline
- Generated PDF roofline plots can have overlapping data points but should still be instructive
 - There are fixes to this, but they may be difficult to setup for different cluster installations
 - Generating the PDF plots from the command line interface should always work
- Complete sets of Roofline plots and commands can be found in the READMEs for each exercise

Exercise 1: Problem Roofline Plots



Kernel legend



Exercise 1: Prep to use Omnipperf to Find Kernel Launch Parameters

- Launch parameters are given at the time of the kernel launch, as in lines 49 and 54:
 - `yax<<<grid,block>>>(y,A,x,n,m,result);`
 - Where grid and block are the kernel yax's launch parameters
 - In problem, `grid = (4,1,1)`, and `block = (64,1,1)`
 - In solution, `grid = (2048,1,1)`, and `block = (64,1,1)`
- Sometimes the launch parameters for a given kernel can be obfuscated
- Omnipperf can easily show launch parameter information regardless of the code
 - You just need the dispatch ID
- To generate profiling data, use the commands:
 - `omnipperf profile -n problem --no-roof -- ./problem.exe`
 - `omnipperf profile -n solution --no-roof -- ./solution.exe`
 - `--no-roof` saves time by not generating roofline data – profile commands can take a while
- **Real benchmarks can take prohibitively long to profile** – use smaller representative problems if possible

Exercise 1: CLI Omnipperf Comparisons are Easy

```
omnipperf analyze -p workloads/problem/mi200 -p workloads/solution/mi200 --dispatch 1 --metric 7.1.0 7.1.1 7.1.2
```

```
-----
Analyze
-----
```

```
-----
0. Top Stat
```

	KernelName	Count	Count	Sum(ns)	Sum(ns)	Mean(ns)	Mean(ns)	Median(ns)	Median(ns)	Pct	Pct
0	yax(double*, double*, double*, int, int, double*)	1.00	1.0 (0.0%)	754934306.50	69702016.5 (-90.77%)	754934306.50	69702016.5 (-90.77%)	754934306.50	69702016.5 (-90.77%)	100.00	100.0 (0.0%)

10.8x speedup

```
-----
7. Wavefront
```

```
7.1 Wavefront Launch Stats
```

Index	Metric	Avg	Avg	Min	Min	Max	Max	Unit
7.1.0	Grid Size	256.00	131072.0 (51100.0%)	256.00	131072.0 (51100.0%)	256.00	131072.0 (51100.0%)	Work items
7.1.1	Workgroup Size	64.00	64.0 (0.0%)	64.00	64.0 (0.0%)	64.00	64.0 (0.0%)	Work items
7.1.2	Total Wavefronts	4.00	2048.0 (51100.0%)	4.00	2048.0 (51100.0%)	4.00	2048.0 (51100.0%)	Wavefronts

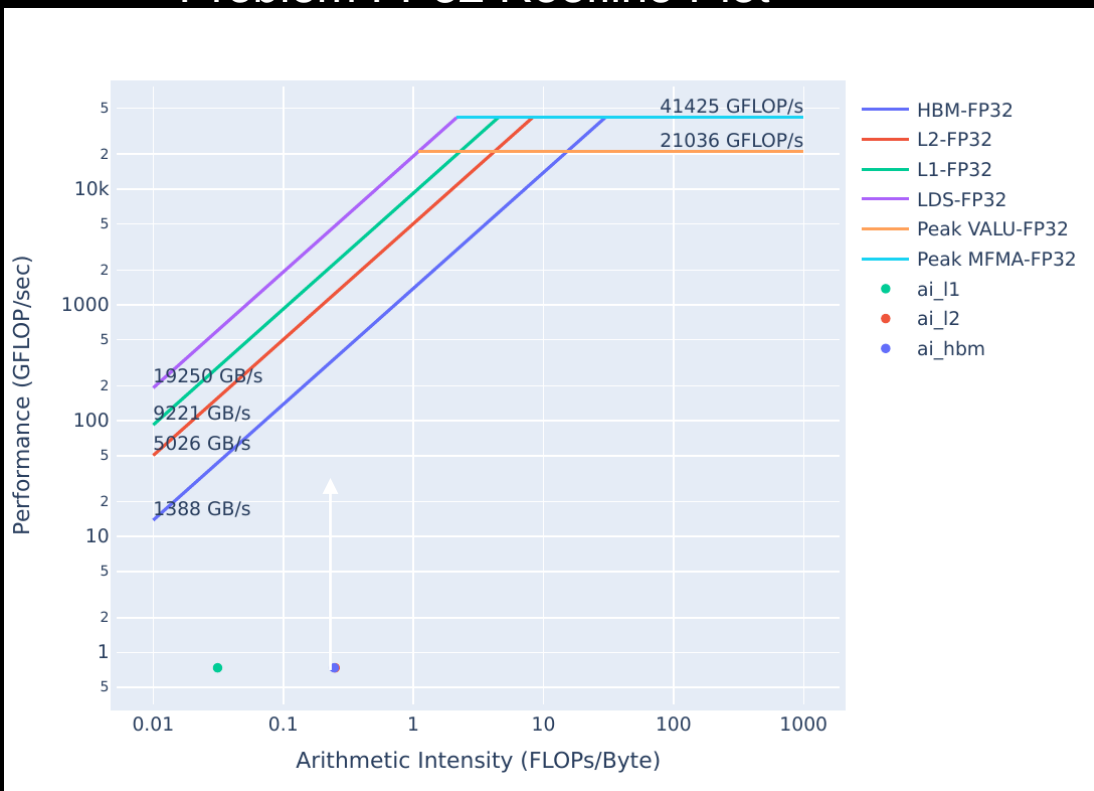
Increased launched wavefronts, which increases
Grid Size

In general, it is difficult to pre-determine optimal launch bounds, so some experimentation is likely necessary

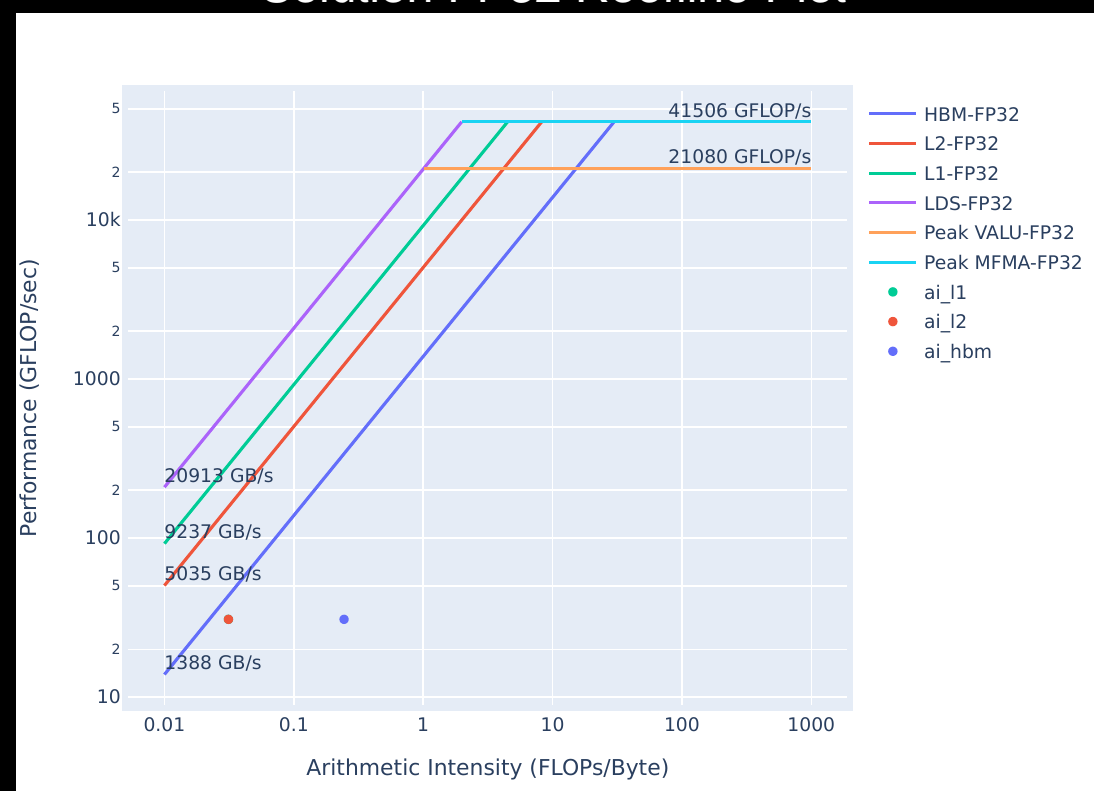
These slides always put `problem` as the baseline, and `solution` as the comparative

Exercise 1: Comparing Problem and Solution Roofline Plots

Problem FP32 Roofline Plot



Solution FP32 Roofline Plot



Generally, moving **up** and to the **right** is good.

Exercise 1: It's Easy to Check Launch Parameters with Omniperf

- Use this omniperf command to check launch parameters:
 - `omniperf analyze -p workloads/problem/mi200 --dispatch 1 --metric 7.1.0 7.1.1 7.1.2`
 - Shows the launch parameters of the kernel with dispatch ID 1
 - `--metric` filters the output to **only** show these launch parameters
- Good launch parameters are essential to a performant GPU kernel
 - Determining which parameters give the best performance usually requires experimenting
- It can be difficult to track down where launch parameters are set in code
- Omniperf can easily show the launch parameters of a kernel
 - Need the dispatch ID or index given by `--list-kernels`
 - `--list-kernels` index can be passed to `-k` as in:
 - `omniperf analyze -p workloads/problem/mi200 -k 0 --metric 7.1.0 7.1.1 7.1.2`
- **Note:**
 - These metric numbers are for Omniperf 1.0.10

Exercise 2: Diagnosing a Shared Memory Occupancy Limiter

- Using LDS (Local Data Store – Shared Memory) to cache re-used data can be an effective optimization strategy
- Using **too much** LDS can restrict occupancy however, and reduce performance
- Line 12 in `problem.cpp` shows the allocation of LDS:
 - `__shared__ double tmp[fully_allocate_lds];`
- There are two solutions:
 - `solution-no-lds` removes the LDS allocation, and thus the occupancy limiter
 - `solution` reduces the size of the LDS allocation, removes occupancy limiter, and is faster than `solution-no-lds`
 - This is the solution used to generate the Omniperf output in the next slide
- Omniperf makes it easy to determine if LDS allocations restrict occupancy, as before profile with:
 - `omniperf profile -n problem --no-roof -- ./problem.exe`
 - `omniperf profile -n solution --no-roof -- ./solution.exe`

Exercise 2: LDS Occupancy Limiter – Relevant Omniperf Output

```
omniperf analyze -p workloads/problem/mi200 -p workloads/solution/mi200 --dispatch 1 --metric 2.1.26 6.2.7
```

Analyze

0. Top Stat

	KernelName	Count	Count	Sum(ns)	Sum(ns)	Mean(ns)	Mean(ns)	Median(ns)	Median(ns)	Pct	Pct
0	yax(double*, double*, double*, int, int, double*)	1.00	1.0 (0.0%)	175427205.00	50366185.0 (-71.29%)	175427205.00	50366185.0 (-71.29%)	175427205.00	50366185.0 (-71.29%)	100.00	100.0 (0.0%)

3.4x speedup

2. System Speed-of-Light

2.1 Speed-of-Light

Index	Metric	Value	Value	Unit	Peak	Peak	PoP	PoP
2.1.26	Wave Occupancy	102.70	487.32 (374.51%)	Wavefronts	3328.00	3328.0 (0.0%)	3.09	14.64 (373.88%)

+ ~11% Occupancy (overall)

6. Shader Processor Input (SPI)

6.2 SPI Resource Allocation

Index	Metric	Avg	Avg	Min	Min	Max	Max	Unit
6.2.7	Insufficient CU LDS	6015745446.00	0.0 (-100.0%)	6015745446.00	0.0 (-100.0%)	6015745446.00	0.0 (-100.0%)	Cu

Sharp decrease in SPI stat

Exercise 2: Use SPI Stats to Determine if LDS Limits Occupancy

- Occupancy limiters can negatively impact performance
- Workgroup manager (SPI) stats in Omniperf indicate whether a kernel resource limits occupancy
- You can get the SPI stat for LDS for a single kernel with:
 - `omniperf analyze -p workloads/problem/mi200 --dispatch 1 --metric 2.1.26 6.2.7`

Note:

- In current Omniperf release 1.0.10, the SPI “insufficient resource” stats are a count of cycles, meaning:
 - Large numbers (on the order of over 1 million) are expected if a field is not zero
 - The magnitude of these fields **does not** necessarily indicate how severely occupancy is impacted
 - If two fields are nonzero, the larger number indicates that resource is limiting occupancy more
- In a coming release, these “insufficient resource” fields are changing to percentages:
 - Large numbers will no longer be expected, but the other points will still hold

Exercise 3: Diagnosing a Register Occupancy Limiter

- Seemingly innocuous function calls inside kernels can lead to unexpected performance characteristics
 - In this case an assert on line 15 causes occupancy to be limited by register usage
 - The solution simply removes the assert
- The types of registers on AMD GPUs are:
 - **VGPRs (Vector General Purpose Registers):** registers that can hold distinct values for each thread in the wavefront
 - **SGPRs (Scalar General Purpose Registers):** uniform across a wavefront. If possible, using these is preferable
 - **AGPRs (Accumulation vector General Purpose Registers):** special-purpose registers for MFMA (Matrix Fused Multiply-Add) operations, or low-cost register spills
- Using too many of one of these register types can impact occupancy and negatively impact performance
- We use the same profile commands to get the profiling data:
 - `omniperf profile -n problem --no-roof -- ./problem.exe`
 - `omniperf profile -n solution --no-roof -- ./solution.exe`

Exercise 3: Register Occupancy Limiter – Relevant Omniperf Output

```
omniperf analyze -p workloads/problem/mi200 -p workloads/solution/mi200 --dispatch 1 --metric 2.1.26 6.2.5 7.1.5 7.1.6 7.1.7
```

0. Top Stat

	KernelName	Count	Count	Sum(ns)	Sum(ns)	Mean(ns)	Mean(ns)	Median(ns)	Median(ns)	Pct	Pct
0	yax(double*, double*, double*, int, int, double*)	1.00	1.0 (0.0%)	76983902.00	69815871.0 (-9.31%)	76983902.00	69815871.0 (-9.31%)	76983902.00	69815871.0 (-9.31%)	100.00	100.0 (0.0%)

Minor speedup

2. System Speed-of-Light

2.1 Speed-of-Light

Index	Metric	Value	Value	Unit	Peak	Peak	PoP	PoP
2.1.26	Wave Occupancy	438.00	444.1 (1.39%)	Wavefronts	3328.00	3328.0 (0.0%)	13.16	13.34 (1.4%)

Small increase in occupancy

6. Shader Processor Input (SPI)

6.2 SPI Resource Allocation

Index	Metric	Avg	Avg	Min	Min	Max	Max	Unit
6.2.5	Insufficient SIMD VGPRs	13733460.00	0.0 (-100.0%)	13733460.00	0.0 (-100.0%)	13733460.00	0.0 (-100.0%)	Simd

Large decrease in SPI stat

7. Wavefront

7.1 Wavefront Launch Stats

Index	Metric	Avg	Avg	Min	Min	Max	Max	Unit
7.1.5	VGPRs	92.00	32.0 (-65.22%)	92.00	32.0 (-65.22%)	92.00	32.0 (-65.22%)	Registers
7.1.6	AGPRs	132.00	0.0 (-100.0%)	132.00	0.0 (-100.0%)	132.00	0.0 (-100.0%)	Registers
7.1.7	SGPRs	48.00	96.0 (100.0%)	48.00	96.0 (100.0%)	48.00	96.0 (100.0%)	Registers

Able to use:
Fewer VGPRs,
No AGPRs,
more SGPRs

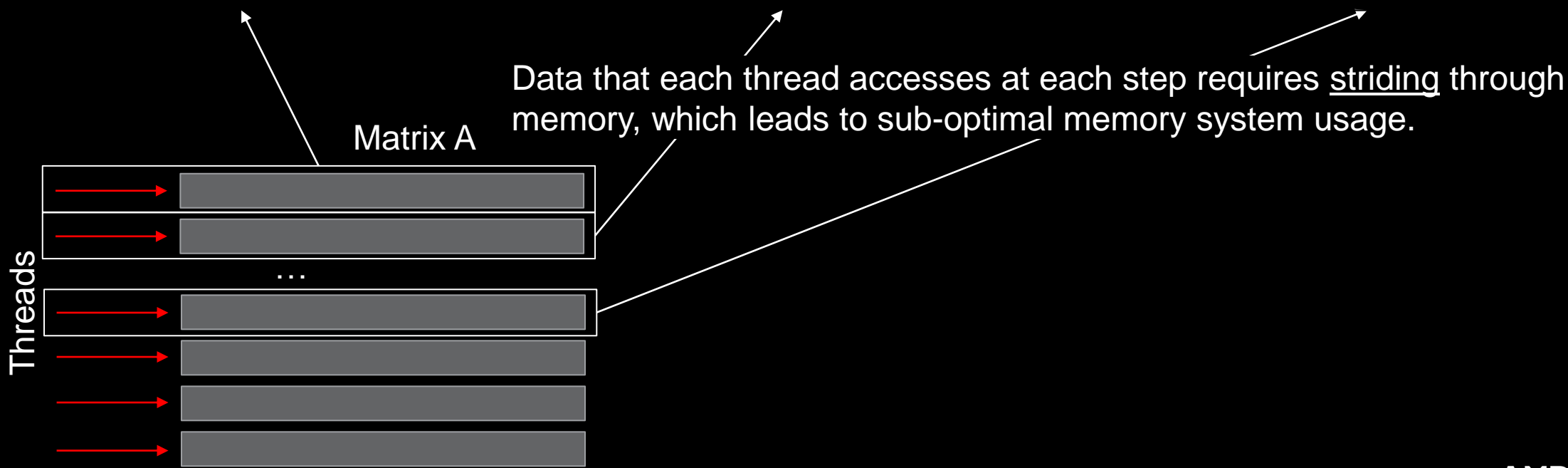
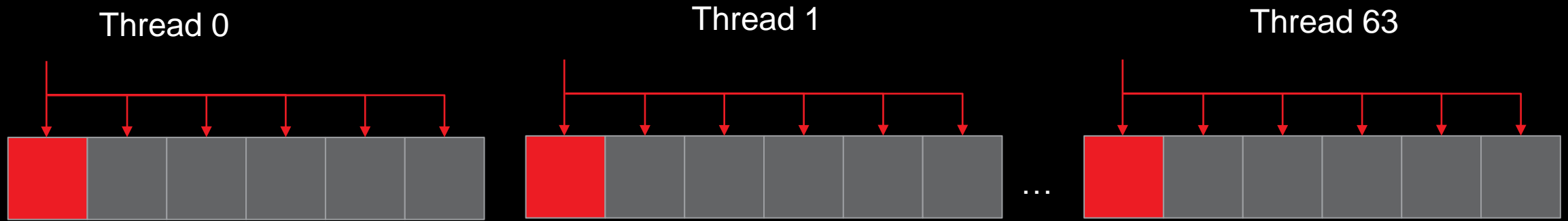
Exercise 3: Register Occupancy Limiter - Takeaways

- Seemingly innocuous function calls inside kernels can lead to unexpected performance characteristics
 - Asserts, and even excessive use of math functions in kernels can degrade performance
- In this case the occupancy limit was very minor, despite a large number in the SPI stat
- AGPR usage in the absence of MFMA (Matrix Fused Multiply Add) instructions can indicate degraded performance.
 - Spilling registers to AGPRs, due to running out of VGPRs
- To determine if any SPI “insufficient resource” stats are nonzero, you can do:
 - `omniperf analyze -p workloads/problem/mi200 --dispatch 1 --metric 6.2`
 - **Note:** This will report more than just all “insufficient resource” fields

Exercise 4: Data Access Patterns are Important to Performance

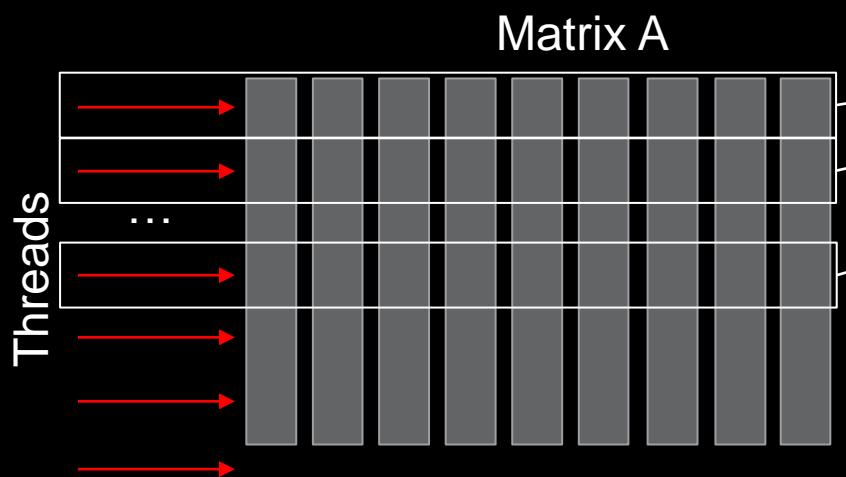
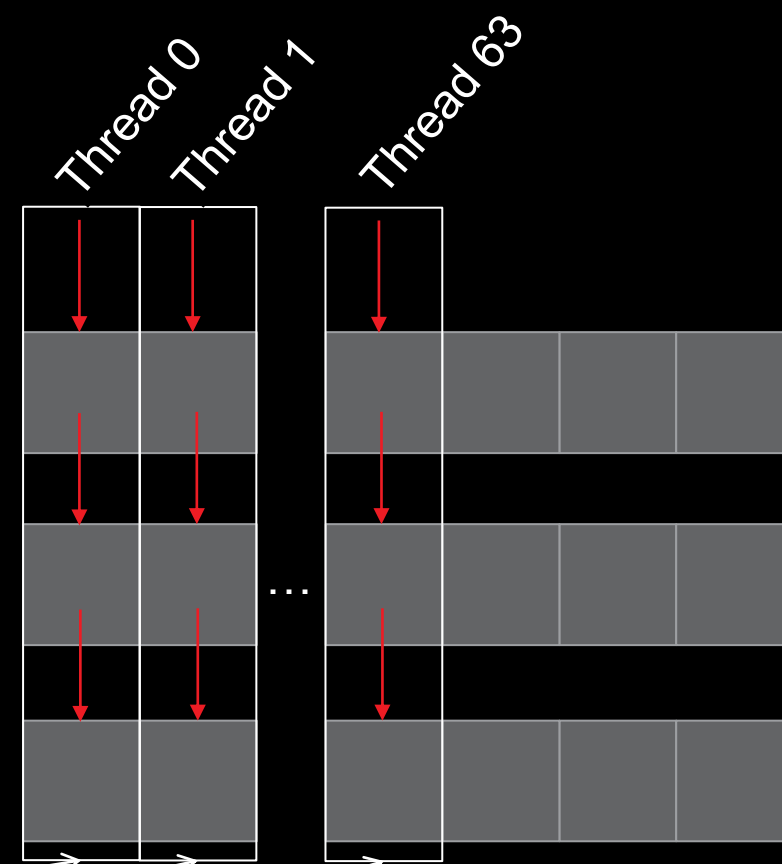
- The way in which threads access memory has a big impact on performance
- “Striding” in global memory has adverse effects on kernel performance, especially on GPUs.
 - “Strided data access patterns” lead to poor utilization of cache memory systems
- These access patterns can be difficult to spot in the code
 - They are valid methods of indexing data
- Using Omniperf can quickly show if a kernel’s data access is adversarial to the caches

Exercise 4: What is a “Strided Data Access Pattern”?



Exercise 4: Strided Data Access Patterns

Increasing the **locality** of data accesses of nearby threads allows for more efficient memory usage



Note: This is the same computation as before, only data layout has changed.

Exercise 4: Using Omnipperf to Diagnose a Strided Data Access Pattern

- This exercise's setup makes it very easy to change the data access pattern
 - Generally, these optimizations can have nontrivial development overhead
 - Re-conceptualizing the data structure can be difficult
- All the solution does is re-work the indexing scheme to better use caches
 - No required change to underlying data, because all the values in y, A, and x are set to 1
- To get started run:
 - `omnipperf profile -n problem --no-roof -- ./problem.exe`
 - `omnipperf profile -n solution --no-roof -- ./solution.exe`

Exercise 4: Strided Data Access Pattern – Relevant Omnipperf Output

```
omnipperf analyze -p workloads/problem/mi200 -p workloads/solution/mi200 --dispatch 1 --metric 16.1 17.1
```

0. Top Stat

	KernelName	Count	Count	Sum(ns)	Sum(ns)	Mean(ns)	Mean(ns)	Median(ns)	Median(ns)	Pct	Pct
0	yax(double*, double*, double*, int, int, double*)	1.00	1.0 (0.0%)	69875592.00	12469690.5 (-82.15%)	69875592.00	12469690.5 (-82.15%)	69875592.00	12469690.5 (-82.15%)	100.00	100.0 (0.0%)

5.6x speedup

16. Vector L1 Data Cache

16.1 Speed-of-Light

Index	Metric	Value	Value	Unit
16.1.0	Buffer Coalescing	25.00	25.0 (0.0%)	Pct of peak
16.1.1	Cache Util	87.80	98.08 (11.7%)	Pct of peak
16.1.2	Cache BW	8.69	12.18 (40.19%)	Pct of peak
16.1.3	Cache Hit	0.00	49.98 (inf%)	Pct of peak

+ ~50% in L1 hit

17. L2 Cache

17.1 Speed-of-Light

Index	Metric	Value	Value	Unit
17.1.0	L2 Util	98.74	98.39 (-0.36%)	Pct
17.1.1	Cache Hit	93.45	0.52 (-99.44%)	Pct
17.1.2	L2-EA Rd BW	125.69	688.98 (448.16%)	Gb/s
17.1.3	L2-EA Wr BW	0.00	0.0 (inf%)	Gb/s

L2 Cache Hit
decreases sharply,
Read BW from HBM
increases by ~5x

The solution better uses the L1, but our L2 hit rate has degraded, which points to a deficiency in our algorithm

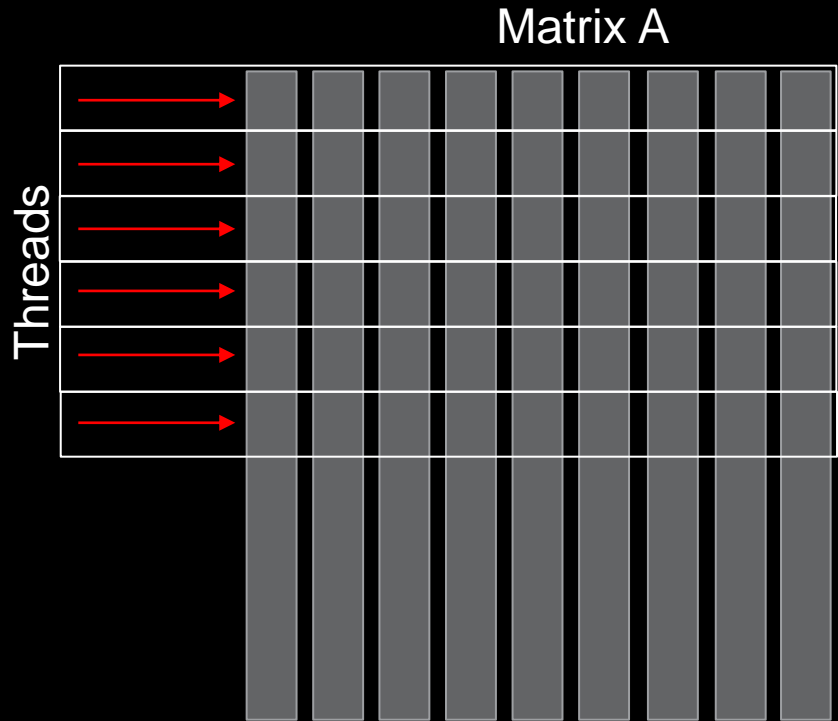
Exercise 4: Omnipperf Speed-of-Light Cache Access Statistics

- This Omnipperf command will show high-level details about L1 and L2 cache accesses:
 - `omnipperf analyze -p workloads/problem/mi200 --dispatch 1 --metric 16.1 17.1`
- Ensuring better data locality will generally provide better performance
- In this case, we start hitting in the L1 cache, rather than having to go out to L2 for everything
- **Note:** In a real code, optimizations of this type likely have much more development overhead
 - Need to change how the data structure is indexed everywhere

Exercise 5: Algorithmic Optimizations

- These types of optimizations are the most difficult to execute
 - Generally, it is difficult to determine if the runtime of one algorithm will be faster than another
- We start with the solution from last exercise as our problem
 - Speed-of-light cache statistics showed that we had ~0% hit rate in the L2, could it be better?
- Our initial algorithm is naïve in terms of parallelization:
 - Each thread computes the sum of a row
- Exposing more parallelism is possible and should get us more performance in this case

Exercise 5: Algorithmic Optimizations



In our current algorithm, each thread computes the sum of a single row

Exercise 5: Algorithmic Optimizations



In a more efficient implementation, wavefronts have multiple threads sum up the rows in parallel, using shared memory to reduce partial sums

Note: The original data layout allows the wavefronts to avoid striding memory

Exercise 5: Using Omnipperf to Evaluate an Algorithmic Optimization

- The strided data access pattern issue is everywhere
 - This solution gets about 2x faster when the data layout is switched to optimize locality
- Though the solution shows a **29x speedup** from the problem, cache speed-of-light stats aren't convincing
 - The rooflines for these problems do not tell the full performance story either
- Running the solution shows it is much faster, but does it use the caches more efficiently?
- To get started, run:
 - `omnipperf profile -n problem --no-roof -- ./problem.exe`
 - `omnipperf profile -n solution --no-roof -- ./solution.exe`

Exercise 5: Sometimes the Full Story is in the Details

```
omniperf analyze -p workloads/problem/mi200 -p workloads/solution/mi200 --dispatch 1 --metric 16.3 17.2 17.3
0. Top Stat
```

	KernelName	Count	Count	Sum(ns)	Sum(ns)	Mean(ns)	Mean(ns)	Median(ns)	Median(ns)	Pct	Pct
0	yax(double*, double*, double*, int, int, double*)	1.00	1.0 (0.0%)	12443928.00	408316.0 (-96.72%)	12443928.00	408316.0 (-96.72%)	12443928.00	408316.0 (-96.72%)	100.00	100.0 (0.0%)

16. Vector L1 Data Cache

16.3 L1D Cache Accesses

~29x faster

Index	Metric	Avg	Avg	Min	Min	Max	Max	Unit
16.3.0	Total Req	524368.00	16448.0 (-96.86%)	524368.00	16448.0 (-96.86%)	524368.00	16448.0 (-96.86%)	Req per wave
...								
16.3.5	Cache Accesses	131140.00	4097.0 (-96.88%)	131140.00	4097.0 (-96.88%)	131140.00	4097.0 (-96.88%)	Req per wave
16.3.6	Cache Hits	65538.00	2864.0 (-95.63%)	65538.00	2864.0 (-95.63%)	65538.00	2864.0 (-95.63%)	Req per wave
16.3.7	Cache Hit Rate	49.98	69.9 (39.87%)	49.98	69.9 (39.87%)	49.98	69.9 (39.87%)	Pct

- ~32x

- ~32x

+ ~40%

17. L2 Cache

17.2 L2 - Fabric Transactions

Index	Metric	Avg	Avg	Min	Min	Max	Max	Unit
17.2.0	Read BW	4194916.56	65688.69 (-98.43%)	4194916.56	65688.69 (-98.43%)	4194916.56	65688.69 (-98.43%)	Bytes per wave

- ~64x

17.3 L2 Cache Accesses

Index	Metric	Avg	Avg	Min	Min	Max	Max	Unit
17.3.0	Req	32945.33	617.41 (-98.13%)	32945.33	617.41 (-98.13%)	32945.33	617.41 (-98.13%)	Req per wave
...								
17.3.6	Hits	171.28	104.03 (-39.27%)	171.28	104.03 (-39.27%)	171.28	104.03 (-39.27%)	Hits per wave
17.3.7	Misses	32774.06	513.38 (-98.43%)	32774.06	513.38 (-98.43%)	32774.06	513.38 (-98.43%)	Misses per wave
17.3.8	Cache Hit	0.52	16.85 (3140.15%)	0.52	16.85 (3140.15%)	0.52	16.85 (3140.15%)	Pct

- ~53x

- ~64x

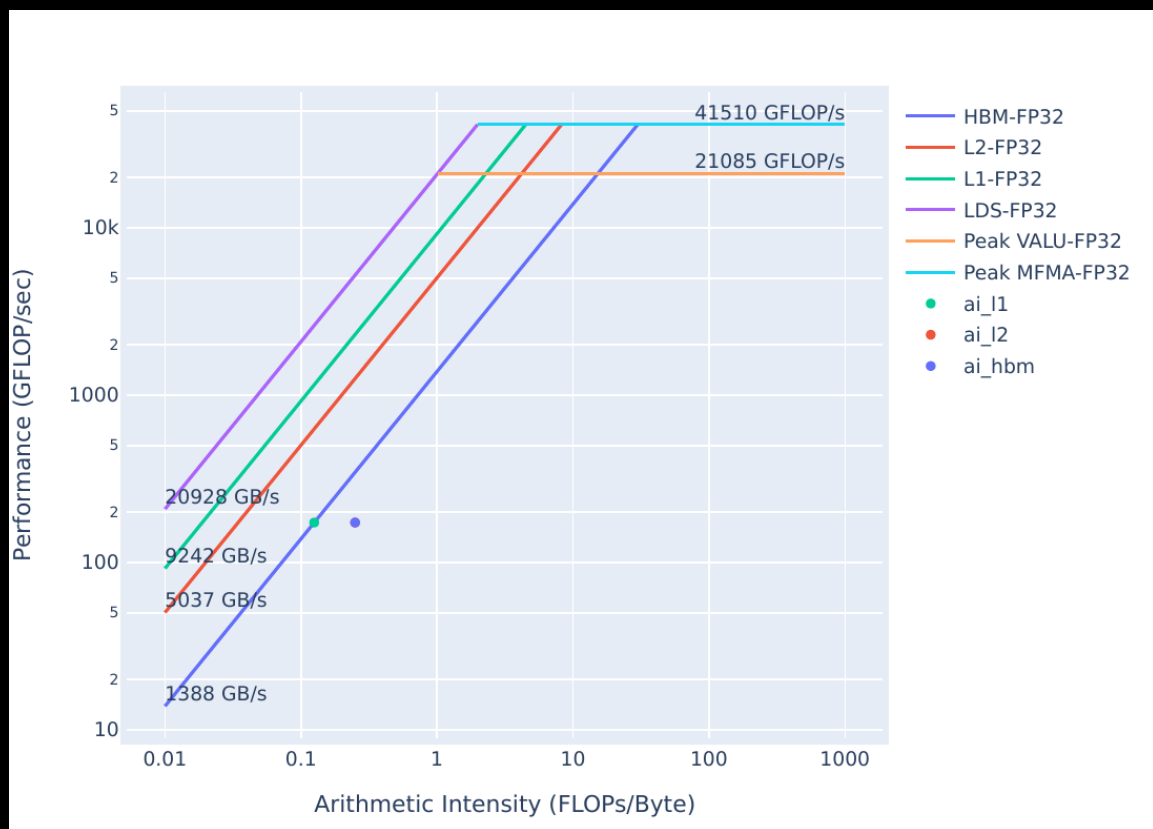
Large relative gain, + ~16% overall

Cache hit rates alone do not give a convincing reason for our performance increase

Exercise 5: It Can Be Hard to Compare Rooflines Between Algorithms

- `omniperf profile -n problem_roof_only --roof-only --kernel-names -- ./problem.exe`
- `omniperf profile -n solution_roof_only --roof-only --kernel-names -- ./solution.exe`

Problem FP32 Roofline



Solution FP32 Roofline



problem is closer to being HBM bandwidth bound: It needs to request much more data from HBM than the optimized version

Exercise 5: Omnipperf Detailed Cache Statistics - Takeaways

- To get detailed cache statistics (including data movement) for kernel with dispatch ID 1:
 - `omnipperf analyze -p workloads/problem/mi200 --dispatch 1 --metric 16.2 16.3 17.2 17.3`
 - **Note:** The slide omitted some Omnipperf output from this metric filtering
- Algorithmic optimizations can be powerful, but are usually time-intensive to design and implement
- It can be difficult to understand the performance differences between algorithms
 - Rooflines can be misleading
 - Assuming correctness is verified, timings don't lie
 - Detailed profiling data can help shed light on the *why* of performance differences