

Лекция 4

Инструменты профилирования:

- nvprof
- Nsight Compute CLI
- nvvp
- Nsight Compute

nvprof и Nsight Compute CLI

```
__global__ void gInit(float* a, float* b){  
    int i=threadIdx.x+blockIdx.x*blockDim.x;  
    a[i]=(float)2*i;  
    b[i]=(float)(2*i+1);  
}
```

Тестовые ядра

```
__global__ void gSum(float* a, float *b){  
    int i=threadIdx.x+blockIdx.x*blockDim.x;  
    a[i]+=b[i];  
}
```



```
/Lecture3/Lab3-cuda-gdb # ncu --target-processes all ./lab3c
```

```
gInit(int *, int *), 2023-Feb-13 15:09:06, Context 1, Stream 7
```

```
Section: GPU Speed Of Light Throughput
```

```
-----  
DRAM Frequency cycle/nsecond          6.40  
SM Frequency    cycle/nsecond          1.29  
Elapsed Cycles          cycle          3,327  
Memory [%]              %              1.10  
DRAM Throughput         %              0.02  
Duration                usecond          2.56  
-----
```

WRN This kernel grid is too small to fill the available resources on this device, resulting in only 0.0 full waves across all SMs. Look at Launch Statistics for more details.

```
.....
```

```
/Lecture3/Lab3-cuda-gdb # ncu
```

```
--metrics gpu__time_duration.sum ./lab3c
```

```
gInit(int *, int *), 2023-Feb-13 18:42:52, Context 1, Stream 7
```

```
Section: Command line profiler metrics
```

```
-----
```

```
gpu  time duration.sum    usecond                29.50
```

```
-----
```

```
gSum(int *, int *), 2023-Feb-13 18:42:52, Context 1, Stream 7
```

```
Section: Command line profiler metrics
```

```
-----
```

```
gpu  time duration.sum    usecond                37.57
```

```
-----
```

```
/Lecture3/Lab3-cuda-gdb> nvprof --query-metrics  
===== Warning: Skipping profiling on device 0 since  
profiling is not supported on devices with compute  
capability 7.5 and higher.
```

Use NVIDIA Nsight Compute for GPU profiling and NVIDIA Nsight Systems for GPU tracing and CPU sampling.

Refer <https://developer.nvidia.com/tools-overview> for more details.

```
ip-011@linux-47dw:/home/malkov/WORKSHOP/PGP-2023>
nvprof --query-metrics | less
```

Available Metrics:	Name	Description
Device 0 (GeForce GTX 1050):		
inst_per_warp:	Average number of instructions executed by each warp	
warp_execution_efficiency:	Ratio of the average active threads per warp to the maximum number of threads per warp supported on a multiprocessor	
.....		
gld_transactions_per_request:	Average number of global memory load transactions performed for each global memory load.	
gst_transactions_per_request:	Average number of global memory store transactions performed for each global memory store	
.....		

```
ip-011@linux-47dw:/home/malkov/WORKSHOP/PGP-2023>
```

```
nvprof -m gst_throughput ./lab3c
```

Invocations	Metric Name	Metric Description	Min	Max	Avg
Device "GeForce GTX 1050 (0)"					
Kernel: gSum(int*, int*)					
1	gst_throughput	Global Store Throughput	40.582MB/s	40.582MB/s	40.582MB/s
Kernel: gInit(int*, int*)					
1	gst_throughput	Global Store Throughput	71.303MB/s	71.303MB/s	71.302MB/s

/Лекция4/lab4> **ncu --list-sections**

Identifier	Display Name	Enabled	Filename

<i>ComputeWorkloadAnalysis</i>	Compute Workload Analysis	yes	...2024.2.1/Sections/ ComputeWorkloadAnalysis.section
<i>InstructionStats</i>	Instruction Statistics	yes	...2024.2.1/Sections/ InstructionStatistics.section
<i>LaunchStats</i>	Launch Statistics	yes	...2024.2.1/Sections/ LaunchStatistics.section
<i>MemoryWorkloadAnalysis</i>	Memory Workload Analysis	yes	...
.....			

```
/Лекция4/lab4> ncu --section InstructionStats ./lab4c
```

```
glNit(float *, float *) (2, 1, 1)x(128, 1, 1), Context 1, Stream 7, Device 0, CC 7.5
```

```
Section: Instruction Statistics
```

```
-----  
Metric Name                      Metric Unit Metric Value
```

```
-----  
Avg. Executed Instructions Per Scheduler    inst    0,93  
Executed Instructions                      inst    112  
Avg. Issued Instructions Per Scheduler      inst    1,27  
Issued Instructions                       inst    152  
-----
```

```
gSum(float *, float *) (2, 1, 1)x(128, 1, 1), Context 1,  
Stream 7, Device 0, CC 7.5
```

```
Section: Instruction Statistics
```

```
-----
```

/Лекция4/lab4> **ncu --section ComputeWorkloadAnalysis ./lab4c**

gSum(float *, float *) (2, 1, 1)x(128, 1, 1), Context 1, Stream 7, Device 0, CC 7.5

Section: Compute Workload Analysis

Metric Name	Metric Unit	Metric Value
Executed Ipc Active	inst/cycle	0,04
Executed Ipc Elapsed	inst/cycle	0,00
Issue Slots Busy	%	1,35
Issued Ipc Active	inst/cycle	0,05
SM Busy	%	1,35

OPT Est. Local Speedup: 99.33%

All compute pipelines are under-utilized. Either this kernel is very small or it doesn't issue enough warps per scheduler. Check the Launch Statistics and Scheduler Statistics sections for further details.

```
/Лекция4/lab4> ncu --query-metrics > metrics.txt
```

```
Device NVIDIA GeForce RTX 2060 (TU104)
```

Metric Name	Metric Type	Metric Unit	Metric Description
dram__bytes	Counter	byte	# of bytes accessed in DRAM
dram__bytes_read	Counter	byte	# of bytes read from DRAM
dram__bytes_write	Counter	byte	# of bytes written to DRAM
.....			
smsp__average_inst_executed_pipe_lsu_per_warp	Ratio	inst/warp	average # of instructions executed by pipe lsu per warp
.....			

/Лекция4/lab4> ***ncu --metrics***

l1tex__t_bytes_pipe_lsu_mem_global_op_st.sum.per_second ./lab4c

glnit(float *, float *) (2, 1, 1)x(128, 1, 1), Context 1, Stream 7, Device 0, CC 7.5

Section: Command line profiler metrics

Metric Name	Metric Unit	Metric Value

l1tex__t_bytes_pipe_lsu_mem_global_op_st.sum.per_second	Mbyte/s	688,17

gSum(float *, float *) (2, 1, 1)x(128, 1, 1), Context 1, Stream 7, Device 0, CC 7.5

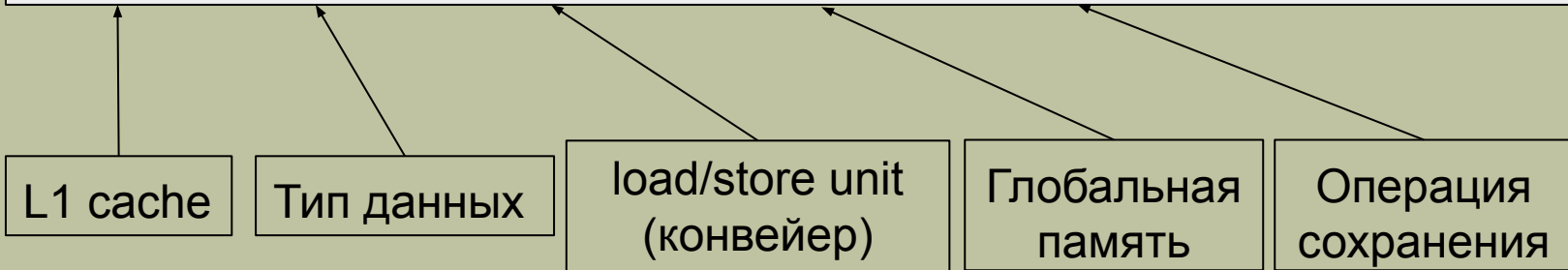
Section: Command line profiler metrics

Metric Name	Metric Unit	Metric Value

l1tex__t_bytes_pipe_lsu_mem_global_op_st.sum.per_second	Mbyte/s	347,83

Кодирование метрики ncu:

l1tex_t_bytes_pipe_lsu_mem_global_op_st.sum.per_second ./lab4c



```
ip-011@linux-47dw:/home/malkov/WORKSHOP/PGP-2023>
```

```
nvprof -m gld_throughput ./lab3c
```

Invocations	Metric Name	Metric Description	Min	Max	Avg
Device "GeForce GTX 1050 (0)"					
	Kernel: gInit(int*, int*)				
1	gld_throughput	Global Load Throughput	0.0B/s	0.0B/s	0.0B/s
	Kernel: gSum(int*, int*)				
1	gld_throughput	Global Load Throughput	87.694MB/s	87.694MB/s	87.694MB/s

```
/Lecture3/Lab3-cuda-gdb # ncu --metrics
```

```
lltex  t_bytes_pipe_lsu_mem_global_op_ld.sum.per_second  
./lab3c
```

```
gInit(int *, int *), 2023-Feb-13 15:25:41, Context 1, Stream 7
```

```
Section: Command line profiler metrics
```

```
-----
```

```
lltex  t bytes pipe lsu mem global op ld.sum.per_second  
byte/second                                0
```

```
-----
```

```
gSum(int *, int *), 2023-Feb-13 15:25:41, Context 1, Stream 7
```

```
Section: Command line profiler metrics
```

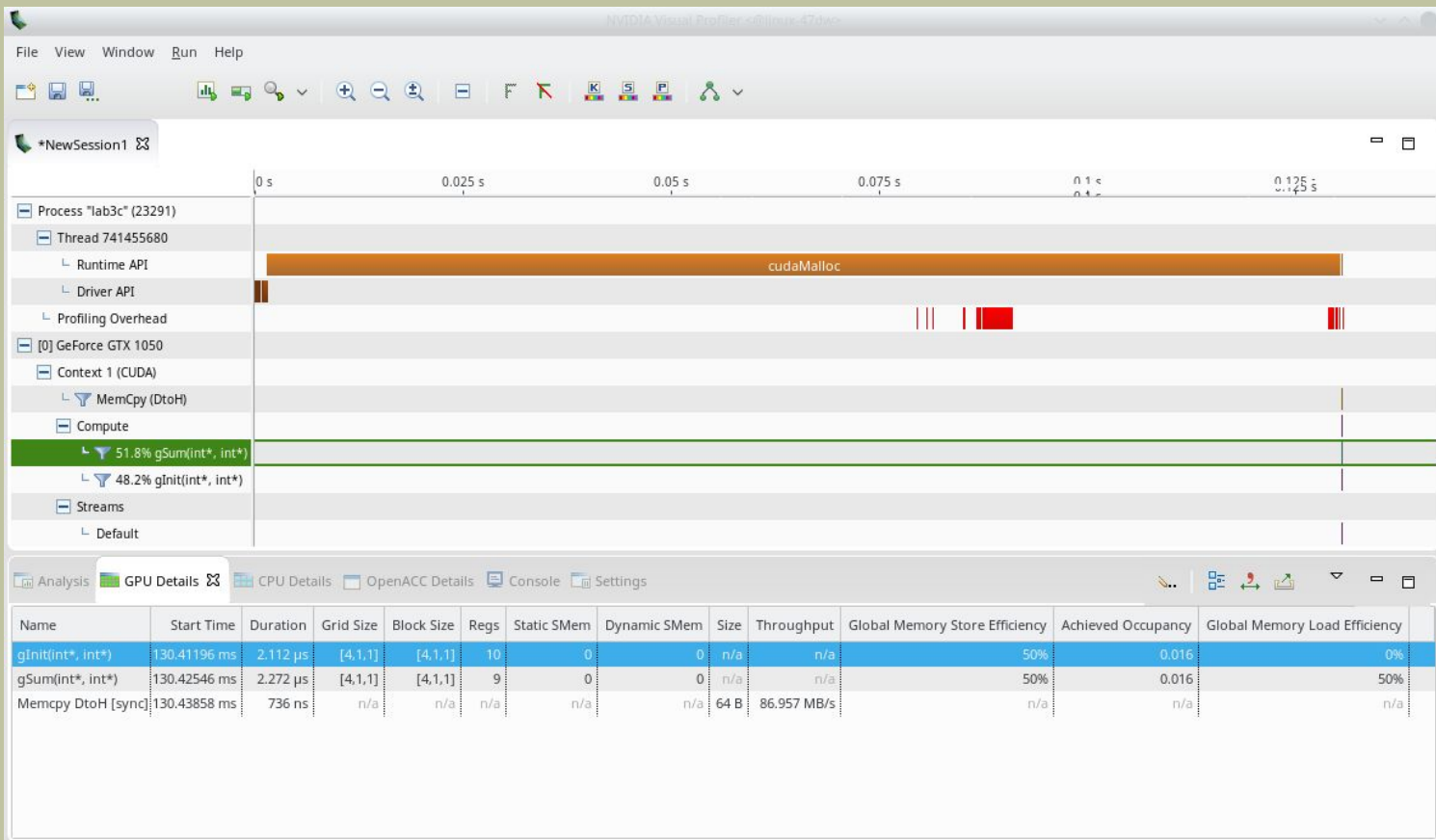
```
-----
```

```
lltex  t bytes pipe lsu mem global op ld.sum.per_second  
Mbyte/second                              82.47
```

```
-----
```


nvvp и Nsight Compute

```
ip-011@linux-47dw: /home/malkov/WORKSHOP/PGP-2023> nvvp ./lab3c
```



0.1 s0.125 s

Achieved Occupan

0.0160%

Metrics and Events <@linux-47dw>

Metrics and Events

Select metrics and events to be collected on individual devices

Device: [0] GeForce GTX 1050

MetricsEvents

☐ Device Memory Write Throughput

☐ Device Memory Write Transactions

☐ ECC Throughput

☐ ECC Transactions

☐ Global Load Throughput

☐ Global Load Transactions

☐ Global Load Transactions Per Request

☒ Global Memory Load Efficiency

☒ Global Memory Store Efficiency

☐ Global Store Throughput

☐ Global Store Transactions

☐ Global Store Transactions Per Request

☐ L2 Read Transactions

☐ L2 Write Transactions


Apply and Run

Cancel

OK

Open a dialog to configure metrics and events, and to run the application to collect

/Лекция4/lab4> ncu-ui &



NVIDIA Nsight Compute

2024.2.1.0 (build 34372528) (public-release)

File Connection Debug Profile Tools Window Help

Connect Disconnect Terminate Profile Kernel


Baselines Metric Details Launch Details


Project Explorer


Welcome


Search project...

Default Project


**Start Activity...**
Start a new Profile, Interactive Profile, Occupancy Calculator, or other activity.


**Open File...**
Open a previously saved file.


**New Project...**
Create a new project.

**Load Project...**
Load a previously saved project.

Continue

 lab3c-25.ncu-rep
22.02.2025 17:17

**NVIDIA Nsight Compute**
2024.2.1.0 (build 34372528) (public-release)

 NVIDIA Nsight Compute 2025.1.0 is now available. [Download now](#), or see [what's new](#).

What's New (4/5)
Source Comparison
Added support for SASS view, Source Markers. Improved diff visualization by adding empty lines on other side of inserted/deleted lines.
Released in 2024.1.0

View: SASS

Report	Result	Time	GPU	SM Frequency	CC
sobelDouble	632 - Sobel (64, 64, 1)x(16, 16, 1)	627.87 usecond	0 - NVIDIA RTX A4500	1.05 cycle/usecond	8.6

Source: Sobel

Navigation: Instructions Executed

# Address	Source	Live arp Stall Sampling	Registers (All Samples)
12 00007f99 c3264cb8	IMAD R13, R13, R13, c[0x0]	15	< 0.61%
13 00007f99 c3264cc9	IAD03 R12, R13, -R13, R2	16	< 0.61%
14 00007f99 c3264cd0	IMAD IAD0 R3, R13, R41, R13	17	< 0.61%
15 00007f99 c3264ce9	ISETP.GT.AND R0, PT, R13, c[0]	17	< 0.61%
16 00007f99 c3264cf8	LOPS.LUT R0, R12, R3, R2, R4	18	< 0.61%
17 00007f99 c3264d00	IMAD.WIDE R0, R4, R13, c[0x0]	19	< 0.61%
18 00007f99 c3264d18	ISETP.GE.AND R0, PT, R13, c[0]	19	< 0.61%
19 00007f99 c3264d28	ISETP.LT.OR P3, PT, R0, R2	19	< 0.61%
20 00007f99 c3264d38	LOPS.LUT R0, R13, R3, R2, R4	19	< 0.61%
21 00007f99 c3264d40	ISETP.GT.OR P3, PT, R13, c[0]	19	< 0.61%
22 00007f99 c3264d50	ISETP.LT.OR R5, PT, R0, R2	19	< 0.61%
23 00007f99 c3264d60	ISETP.GT.OR R0, PT, R13, c[0]	18	0.62%
24 00007f99 c3264d70	LDO.E.UR R0, [R0, R4+R0]	18	0.62%
25 00007f99 c3264d80	LDO.E.UR R14, [R0, R4]	18	0.62%
26 00007f99 c3264d90	IAD03 R13, R0, R41, R2	19	< 0.61%
27 00007f99 c3264da0	CS2R R0, R02	21	< 0.61%
28 00007f99 c3264db0	LOPS.LUT R0, R13, R13, R2, R4	22	< 0.61%
29 00007f99 c3264dc9	LOPS.LUT R13, R13, R3, R2, R4	22	< 0.61%
30 00007f99 c3264de0	ISETP.GE.AND P1, PT, R13, c[0]	22	0.62%
31 00007f99 c3264df0	ISETP.LT.OR P4, PT, R0, R2	22	< 0.61%
32 00007f99 c3264e00	ISETP.LT.OR R2, R0, R2, R4	22	< 0.61%

Report	Result	Time	GPU	SM Frequency	CC
sobelFloat	632 - Sobel (64, 64, 1)x(16, 16, 1)	31.55 usecond	0 - NVIDIA RTX A4500	937.59 cycle/usecond	8.6

Source: Sobel

Navigation: Instructions Executed

# Address	Source	Live arp Stall Sampling	Registers (All Samples)
12 00007f26 c3264cb8	IMAD R13, R0, R13, c[0x0]	13	< 0.61%
13 00007f26 c3264cc9	IAD03 R2, R0, -R13, R2	14	< 0.61%
14 00007f26 c3264cd0	IMAD IAD0 R0, R0, R41, R11	15	0.99%
15 00007f26 c3264ce9	ISETP.GT.AND P4, PT, R0, c[0]	15	0.13%
16 00007f26 c3264cf8	LOPS.LUT R0, R2, R3, R2, R4	16	0.59%
17 00007f26 c3264d00	IMAD.WIDE R0, R4, R13, c[0x0]	17	0.67%
18 00007f26 c3264d18	ISETP.GE.AND P2, PT, R0, c[0]	17	0.67%
19 00007f26 c3264d28	ISETP.LT.OR R0, PT, R0, R2	17	0.53%
20 00007f26 c3264d38	LOPS.LUT R0, R0, R3, R2, R4	17	0.31%
21 00007f26 c3264d40	ISETP.GT.OR R0, PT, R0, c[0]	17	0.26%
22 00007f26 c3264d50	ISETP.LT.OR P3, PT, R0, R2	17	3.37%
23 00007f26 c3264d60	ISETP.GT.OR P3, PT, R0, c[0]	16	0.79%
24 00007f26 c3264d70	LDO.E.UR R0, [R0, R4+R0]	16	0.75%
25 00007f26 c3264d80	LDO.E.UR R14, [R0, R4]	16	0.51%
26 00007f26 c3264d90	IAD03 R12, R0, R41, R2	17	0.97%
27 00007f26 c3264da0	LOPS.LUT R0, R0, R3, R2, R4	18	0.31%
28 00007f26 c3264db0	LOPS.LUT R0, R12, R3, R2, R4	19	0.07%
29 00007f26 c3264dc9	ISETP.GE.AND P1, PT, R12, c[0]	19	0.67%
30 00007f26 c3264de0	ISETP.LT.OR P0, PT, R0, R2	19	2.12%
31 00007f26 c3264df0	ISETP.LT.OR R2, R0, R2, R4	19	0.31%

Target Platform

- Linux (aarch64 sbasa)
- Linux (x86_64)
- Windows

Connection: localhost

Launch

Attach

Application Executable: /IOP/EDUCATION/PGP-2025/Лекции/Лекция4/lab4/lab4c

Working Directory: \$(ApplicationDir)

Command Line Arguments:

Environment:

Activity

- Profile
- Interactive Profile
- Occupancy Calculator
- System Trace

Profile an application using the command line profiler. All GPU workloads are serialized. Note: Attach is not supported for this activity.

Supported APIs: CUDA, OptiX

Common

Filter

Metrics

PM Sampling

Warp Sampling

Other

Output File: lab4c-repl

Force Overwrite: Yes

Target Processes: All

Replay Mode: Kernel

Application Replay Match: Grid

Application Replay Buffer: File

Application Replay Mode: Strict

Graph Profiling: Node

Command Line: /opt/nvidia/nsight-compute/2024.2.1/target/linux-desktop-nlibc 2 11 3-x64/ncu --config-file off --export "/

Cancel

Reset Activity

Launch

	Result	Size	Time	Cycles	GPU	SM Frequency	Process	Attributes
Current	542 - glnit	(2, 1, 1)x(128, 1, 1)	1,70 us	2,283	0 - NVIDIA GeForce RTX 2060	1,34 Ghz	[9459] lab4c	

Summary	Details	Source	Context	Comments	Raw	Session
---------	---------	--------	---------	----------	-----	---------

Compare	Tools	View	Export	
---------	-------	------	--------	--

This table shows all results in the report. Use the column headers to sort the results in this report. Double-click a result to see detailed metrics. Double-click on demangled names to rename it.

ID	Estimated Speedup	Function Name	Demangled Name	Duration (4384)	Runtime Improvement (4091.73)	Compute Throughput	Memory Throughput	# Registers
0	93.33	glnit	glnit(float *, float..	1,70	1,58	0,07	1,36	
1	93.33	gSum	gSum(float *, float..	2,69	2,51	0,08	1,52	

The following performance optimization opportunities were discovered for this result. Follow the rule links to see more context on the Details page.

Note: Speedup estimates provide upper bounds for the optimization potential of a kernel assuming its overall algorithmic structure is kept unchanged.

Small Grid

Est. Speedup: 93.33%

The grid for this launch is configured to execute only 2 blocks, which is less than the GPU's 30 multiprocessors. This can underutilize some multiprocessors. If you do not intend to execute this kernel concurrently with other workloads, consider reducing the block size to have at least one block per multiprocessor or increase the size of the grid to fully utilize the available hardware resources. See the [Hardware Model](#) description for more details on launch configurations.

Achieved Occupancy

Est. Speedup: 87.08%

The difference between calculated theoretical (100.0%) and measured achieved occupancy (12.9%) can be the result of warp scheduling overheads or workload imbalances during the kernel execution. Load imbalances can occur between warps within a block as well as across blocks of the same kernel. See the [CUDA Best Practices Guide](#) for more details on optimizing occupancy.

Imc Miss Stalls

Est. Speedup: 74.66%

On average, each warp of this kernel spends 30.3 cycles being stalled waiting for an immediate constant cache (IMC) miss. A read from constant memory costs one memory read from device memory only on a cache miss; otherwise, it just costs one read from the constant cache. Immediate constants are encoded into the SASS instruction as 'c[bank] [offset]'. Accesses to different addresses by threads within a warp are serialized, thus the cost scales linearly with the number of unique addresses read by all threads within a warp. As such, the constant cache is best when threads in the same warp access only a few distinct locations. If all threads of a warp access the same location, then constant memory can be as fast as a register access. This stall type represents about 74.7% of the total average of 40.6 cycles between issuing two instructions.

Welcome ×lab4c-rep.ncu-rep ×

Current

546 - gSum

▽

▽

(2, 1, 1)x(128, 1, 1)

2,69 us

3.575

0 - NVIDIA GeForce RTX 2060

1,32 Ghz

[9459] lab4c

⚙

Summary

Details

Source

Context

Comments

Raw

Session

Compare

Tools

View

Export

⋮

GPU Speed Of Light Throughput

GPU Throughput Chart

🗨

High-level overview of the throughput for compute and memory resources of the GPU. For each unit, the throughput reports the achieved percentage of utilization with respect to the theoretical maximum. Breakdowns show the throughput for each individual sub-metric of Compute and Memory to clearly identify the highest contributor. High-level overview of the utilization for compute and memory resources of the GPU presented as a roofline chart.

Compute (SM) Throughput [%]	0,08	Duration [us]	2,69
Memory Throughput [%]	1,52	Elapsed Cycles [cycle]	3.575
L1/TEX Cache Throughput [%]	7,45	SM Active Cycles [cycle]	80,53
L2 Cache Throughput [%]	1,52	SM Frequency [Ghz]	1,32
DRAM Throughput [%]	0,60	DRAM Frequency [Ghz]	7,30

🔍 Small Grid

This kernel grid is too small to fill the available resources on this device, resulting in only 0.0 full waves across all SMs. Look at [Launch Statistics](#) for more details.

🗨

🔍 Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved close to 0% of this device's fp32 peak performance and 0% of its fp64 peak performance. See the [Kernel Profiling Guide](#) for more details on roofline analysis.

PM Sampling

🗨

Timeline view of PM metrics sampled periodically over the workload duration. Data is collected across multiple passes. Use this section to understand how workload behavior changes over its runtime.

Maximum Sampling Interval [cycle]	20.000	# Pass Groups	1
Maximum Buffer Size [Kbytes]	64	Dropped Samples [sample]	0

Compute Workload Analysis

🗨

Detailed analysis of the compute resources of the streaming multiprocessors (SM), including the achieved instructions per clock (IPC) and the utilization of each available pipeline. Pipelines with very high utilization might limit the overall performance.

Executed Ipc Elapsed [inst/cycle]	0,00	SM Busy [%]	1,32
Executed Ipc Active [inst/cycle]	0,04	Issue Slots Busy [%]	1,32
Issued Ipc Active [inst/cycle]	0,05		

Metric Details

×

Search metrics in current report or for a chip

🔍

sm__throughput.avg.pct_of_peak_s

⏪

⏩

Name	sm__throughput.avg.pct...
Unit	%
Value	0.06984459577440195
Report	lab4c-rep.ncu-rep
Chip	TU104

Additional Information

Description:

SM throughput assuming ideal load balancing across SMSPs (This throughput metric represents the percent of the peak sustained rate achieved during elapsed cycles across

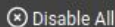
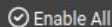
Knowledgebase Entry:

sm: The Streaming Multiprocessor handles execution of a kernel as groups of 32 threads, called warps. Warps are further grouped into cooperative thread arrays (CTA), called

Suffix	Value
.avg	0.069844595774401...
.sum	0.069844595774401...
.min	0
.max	1.0476689366160292

Metric Selection

Metric Sections/Rules

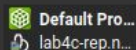


Enter filter

Name	Priority	Description	Sets	Metrics	Filename	State
▶ GPU Speed Of Light Throughput (3)	10	High-level overview of the throughput for compute and memory resource...	basic,detailed,f...	(53) arch:50:70:dram__cycles_elap...	SpeedOfLight.s...	Stock
▶ GPU Speed Of Light Roofline Chart (1)	11	High-level overview of the utilization for compute and memory resources...	detailed,full,roo...	(62) arch:50:70:dram_bytes.sum...	SpeedOfLight_R...	Stock
GPU Speed Of Light Hierarchical Roofline C...	12	High-level overview of the utilization for compute and memory resources...	roofline	(98) arch:50:70:dram_bytes.sum....	SpeedOfLight_...	Stock
GPU Speed Of Light Hierarchical Roofline C...	12	High-level overview of the utilization for compute and memory resources...	roofline	(98) arch:50:70:dram_bytes.sum....	SpeedOfLight_...	Stock
GPU Speed Of Light Hierarchical Roofline C...	12	High-level overview of the utilization for compute and memory resources...	roofline	(98) arch:50:70:dram_bytes.sum...	SpeedOfLight...	Stock

Metrics: Enter metrics, e.g. metric1,metric2

Search project...



Default Pro...

lab4c-rep.n...

Identifier: "MemoryWorkloadAnalysis"

DisplayName: "Memory Workload Analysis"

Description: "Detailed analysis of the memory resources of the GPU. Memory can become a limiting factor for the overall kernel performance when fully utilizing the involved hardware units (Mem Busy), exhausting the available communication bandwidth between those units (Max Bandwidth), or by reaching the maximum throughput of issuing memory instructions (Mem Pipes Busy)."

Order: 30

Sets {

Identifier: "detailed"

}

Sets {

Identifier: "full"

}

Header {

Metrics {

Label: "Memory Throughput"

Name: "dram_bytes.sum.per_second"

Filter {

MaxArch: CC_70

}

Options {

Name: "dram_bytes.sum.per_second"

Filter {

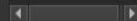
MinArch: CC_75

MaxArch: CC_86

}

Options {

Name: "dram_bytes.sum.per_second"

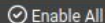


Metric Selection

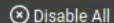
Metric Sections/Rules



Reload



Enable All



Disable All



Restore

Enter filter

Name	Priority	Description	Sets	Metrics	Filename	State
<input checked="" type="checkbox"/> Memory Workload Analysis	30	Detailed analysis of the memory resources of the GPU. Memory can bec...	detailed,full	(22) arch:50:70:dram_bytes.sum...	MemoryWorklo...	Stock
▶ <input type="checkbox"/> Memory Workload Analysis Chart (2)	31	Detailed chart of the memory units.	detailed,full	(38) arch:50:70:its_t_sectors_srcu...	MemoryWorklo...	Stock
▶ <input type="checkbox"/> Memory Workload Analysis Tables (2)	32	Detailed tables with data for each memory unit.	full	(44) arch:80:86:group.memory_l2...	MemoryWorklo...	Stock
▶ <input type="checkbox"/> Scheduler Statistics (1)	40	Summary of the activity of the schedulers issuing instructions. Each sch...	full	(25) smsp_issue_active.avg.pct_o...	SchedulerStat...	Stock
▶ <input type="checkbox"/> Warn State Statistics (2)	50	Analysis of the states in which all warns spent cycles during the kernel e...	full	(27) arch:90:90:smsn_averane_w...	WarnStateStati...	Stock

Metrics: Enter metrics, e.g. metric1,metric2

Metric Details

Search metrics in current report or for a chip

sm__throughput.avg.pct_of_peak_s...

Name	sm__throughput.avg.pct...
Unit	%
Value	0.06984459577440195
Report	lab4c-rep.ncu-rep
Chip	TU104

Additional Information

Description:

SM throughput assuming ideal load
Knowing how many SMs
sm: The Streaming Multiprocessor

Suffix	Value
.avg	0.069844595774401...
.sum	0.069844595774401...
.min	0
.max	1.0476689366160292