

GPU Speed of Light

High-level overview of the utilization for compute and memory resources of the GPU. For each unit, the Speed of Light (SoL) reports the achieved percentage of utilization with respect to the theoretical maximum.

SoL SM [%]	48.64 (+745.34%)	Duration [seconds]	2.23 (-94.15%)
SoL Memory [%]	47.21 (-35.24%)	Elapsed cycles [cycle]	3,885,874 (-94.15%)
SoL Tex [%]	47.75 (-44.48%)	SM Active Cycles [cycle]	3,842,115.75 (-94.15%)
SoL L1 [%]	5.47 (-28.45%)	SM Frequency [cycles/second]	1.75 (-49.18%)
SoL PE [%]	9.08 (+1,382.38%)	Memory Frequency [cycles/second]	1.78 (+1.73%)

GPU Utilization

Recommendations

Bottleneck The kernel is utilizing greater than 80.0% of the available compute or memory performance of the device. To further improve performance, work will likely need to be shifted from the most utilized to another unit.

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 Inst. Elapsed [Inst/cycle]	1.21 (+480.35%)	SM Busy [%]	35.16 (+87.87%)
Executed Inst. Active [Inst/cycle]	1.22 (+480.15%)	Issue slots Busy [%]	38.56 (+88.13%)
Issued ops Active [Inst/cycle]	1.22 (+480.15%)	-	-

Pipe Utilization

Memory Workload Analysis

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 (Mem Bandwidth), or by reaching the maximum throughput of issuing memory instructions (Mem Pipes Busy). Detailed chart of the memory units. Detailed tables with data for each memory unit.

Memory Throughput [bytes/second]

L1 Hit Rate [%]

L2 Hit Rate [%]

Memory Chart

The diagram illustrates the memory architecture and data flow. On the left, a vertical green bar represents the **Kernel**. It connects to several memory units: **Global**, **Local**, **Texture**, **Surface**, and **Shared**. Each unit shows its memory throughput and hit rate. These units connect to a central **Unified Cache**, which in turn connects to the **L2 Cache**. The **L2 Cache** then connects to **System Memory** and **Device Memory**. The **Shared** memory unit also connects to **Shared Memory**. The diagram uses arrows to show the direction of data flow and includes numerical values for throughput and hit rates.

Unit	Throughput [bytes/second]	Hit Rate [%]
Global	8.91 M	94.8%
Local	0.00	100%
Texture	0.00	100%
Surface	0.00	100%
Shared	82.89 M	985.7%
Unified Cache	524.29 M	98%
L2 Cache	64.02 M	85.69%
System Memory	0.00 B	100%
Device Memory	64.00 M	100%
Shared Memory	33.00	100%

Shared Memory

Scheduler Statistics

Summary of the activity of the scheduler issuing instructions. Each scheduler maintains a pool of warps that it can issue instructions for. The upper bound of warps in the pool (Theoretical Warps) is limited by the launch configuration. On every cycle each scheduler checks the state of the allocated warps in the pool (Active Warps). Active warps that are not stalled (Eligible Warps) are ready to issue their next instruction. From the set of eligible warps the scheduler selects a single warp from which to issue one or more instructions (Issued Warps). On cycles with no eligible warps, the issue slot is idled and no instruction is issued. Many warps idled near the end indicates poor steering logic.

Active Warps Per Scheduler [warp]	7.78 (+/-10%)	Instructions Per Active Issue Slot [inst/cycle]	1 (+/-100%)
Eligible Warps Per Scheduler [warp]	1.04 (+/-69.1%)	No Eligible [0]	69.79 (+/-24.1%)
Issued Warps Per Scheduler [warp]	0.31 (+/-89.1%)	One or More Eligible [1]	70.51 (+/-89.1%)

Warps Per Scheduler

Category	Value
Theoretical Warps Per Scheduler	~10.5
Active Warps Per Scheduler	~7.8
Eligible Warps Per Scheduler	~1.0
Issued Warps Per Scheduler	~0.3

Category	Red Series	Blue Series
L0P	~45	~75
L0P3	~95	~65

Category	Value
SQL	100
SP	100
SQL	100

The graph shows the weekly occupancy of the 1000-bed ward at the Royal Victoria Infirmary from 1990 to 2000. The occupancy starts at approximately 1400 beds in 1990, rises to a peak of about 1500 beds in 1992, and then fluctuates between 1400 and 1500 beds until 1994. After 1994, the occupancy generally trends upwards, reaching approximately 1550 beds by 2000.