

Page: Details

Process: All

Launch: 693 - cachedRegsAssignmentKernel

Add Baseline

Apply Rules

Save as PDF

Current

693 - cachedRegsAssignmentKernel (5242880, 4, 1) Time: 9.47 second Cycles: 2,472,201,405 Regs: 49 GPU: NVIDIA GeForce GTX 1060 SM Frequency: 260.50 Mhz CC: 6.1 Process: [14884] k-means.exe

Baseline 4

643 - cachedRegsAssignmentKernel (5242880, 4, 1) Time: 1.74 second Cycles: 2,462,538,504 Regs: 49 GPU: NVIDIA GeForce GTX 1060 SM Frequency: 1.42 Ghz CC: 6.1 Process: [14884] k-means.exe

Baseline 5

653 - cachedRegsAssignmentKernel (5242880, 4, 1) Time: 1.70 second Cycles: 2,467,761,367 Regs: 49 GPU: NVIDIA GeForce GTX 1060 SM Frequency: 1.45 Ghz CC: 6.1 Process: [14884] k-means.exe

Baseline 6

Baseline 7

Baseline 8

Baseline 9

Baseline 10

GPU Speed of Light

SOL SM [%]	71.23	(-0.23%, z=-0.86)	Duration [second]	9.47	(+141.77%, z=+1.39)
SOL Memory [%]	44.48	(+0.29%, z=+1.29)	Elapsed Cycles [cycle]	2,47,22,01,405	(+0.07%, z=+0.31)
SOL TEX [%]	6.81	(-0.25%, z=-0.92)	SM Active Cycles [cycle]	2,46,63,41,650.50	(-0.25%, z=-0.92)
SOL L2 [%]	nan	(+0.00%, z=+0.00)	SM Frequency [Mhz]	260.50	(-76.43%, z=-1.29)
SOL FB [%]	2.81	(+63.07%, z=+1.31)	Memory Frequency [Ghz]	1.56	(-74.62%, z=-1.28)

GPU Utilization

Recommendations

Bottleneck

[Warning] Compute is more heavily utilized than Memory: Look at "Compute Workload Analysis" report section to see what the compute pipelines are spending their time doing. Also, consider whether any computation is redundant and could be reduced or moved to look-up tables.

Compute Workload Analysis

Executed Ipc Elapsed [inst/cycle]	3.27	(-0.25%, z=-0.92)	SM Busy [%]	71.23	(-0.23%, z=-0.86)
Executed Ipc Active [inst/cycle]	3.30	(+0.29%, z=+1.29)	Issue Slots Busy [%]	55.06	(+0.29%, z=+1.29)
Issued Ipc Active [inst/cycle]	3.30	(+0.29%, z=+1.29)	-	-	-

Memory Workload Analysis

Memory Throughput [dbyte/second]	1.86	(-48.01%, z=-1.29)	Mem Busy [%]	44.48	(+0.29%, z=+1.29)
L1 Hit Rate [%]	57.44	(+0.02%, z=+0.00)	Max Bandwidth [%]	44.48	(+0.29%, z=+1.29)
L2 Hit Rate [%]	74.36	(-10.65%, z=-1.28)	Mem Pipes Busy [%]	45.22	(+0.29%, z=+1.29)

Memory Chart

Shared Memory

	Instructions	Requests	% Peak	Bank Conflicts
Shared Load	10,74,13,50,400 (+0.00%, z=+0.00)	10,74,13,50,329 (-0.00%, z=+0.00)	44.03 (+0.29%, z=+1.29)	0 (+0.00%, z=+0.00)
Shared Store	11,01,00,480 (+0.00%, z=+0.00)	11,01,00,480 (+0.00%, z=+0.00)	0.45 (+0.29%, z=+0.00)	0 (+0.00%, z=+0.00)
Shared Atomic	0 (+0.00%, z=+0.00)	-	-	-
Total	10,85,14,50,880 (+0.00%, z=+0.00)	10,85,14,50,809 (-0.00%, z=+0.00)	44.48 (+0.29%, z=+1.29)	0 (+0.00%, z=+0.00)

First-Level (Unified) Cache

	Instructions	SM->TEX Requests	% Peak	Hit Rate	TEX->L2 Requests	% Peak	L2->TEX Returns	% Peak	TEX->SM Returns	% Peak
Global Load Cached	10,48,57,600 (+0.00%, z=+0.00)	41,94,30,400 (+0.00%, z=+0.00)	1.72 (+0.29%, z=+1.29)	57.44 (+0.02%, z=+0.28)	-	-	85,68,68,272 (-0.03%, z=-0.28)	1.74 (-0.28%, z=-0.94)	-	-
Global Load Uncached	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	-	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	-
Local Load Cached	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	41,94,30,400 (+0.00%, z=+0.00)	0.85 (-0.25%, z=-0.92)
Local Load Uncached	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	-	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	-
Surface Load	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	-	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	-
Texture Load	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	-	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	-
Global Store	6,55,360 (+0.00%, z=+0.00)	26,21,440 (+0.00%, z=+0.00)	0.01 (+0.29%, z=+0.00)	-	26,21,440 (+0.00%, z=+0.00)	0.01 (-0.25%, z=-0.00)	-	-	-	-
Local Store	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	-	-	-
Surface Store	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	-	-	-
Global Reduction	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	-	-	-
Surface Reduction	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	-	-	-
Global Atomic	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Global Atomic Cas	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Surface Atomic	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Surface Atomic Cas	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Loads	10,48,57,600 (+0.00%, z=+0.00)	41,94,30,400 (+0.00%, z=+0.00)	1.72 (+0.29%, z=+1.29)	57.44 (+0.02%, z=+0.28)	-	-	85,68,68,272 (-0.03%, z=-0.28)	1.74 (-0.28%, z=-0.94)	41,94,30,400 (+0.00%, z=+0.00)	0.85 (-0.25%, z=-0.92)
Stores	6,55,360 (+0.00%, z=+0.00)	26,21,440 (+0.00%, z=+0.00)	0.01 (+0.29%, z=+0.00)	-	26,21,440 (+0.00%, z=+0.00)	0.01 (-0.25%, z=-0.00)	-	-	-	-
Total	10,55,12,960 (+0.00%, z=+0.00)	42,20,51,840 (+0.00%, z=+0.00)	1.73 (+0.29%, z=+1.29)	57.44 (+0.02%, z=+0.28)	26,21,440 (+0.00%, z=+0.00)	0.01 (-0.25%, z=-0.00)	85,68,68,272 (-0.03%, z=-0.28)	1.74 (-0.28%, z=-0.94)	41,94,30,400 (+0.00%, z=+0.00)	0.85 (-0.25%, z=-0.92)

Second-Level (L2) Cache

	TEX->L2 Requests	% Peak	L2->TEX Returns	% Peak	Total Bytes	Total Throughput
Global Load Cached	-	-	85,68,68,272 (-0.03%, z=-0.28)	0 (-100.00%, z=-1.29)	27,41,97,84,704 (-0.03%, z=-0.28)	2,89,61,27,091,51 (-76.48%, z=-1.29)
Global Load Uncached	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Local Load Cached	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Local Load Uncached	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Surface Load	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Texture Load	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Global Store	26,21,440 (+0.00%, z=+0.00)	nan (+0.00%, z=+0.00)	-	-	8,38,86,080 (+0.00%, z=+0.00)	88,60,199,00 (-76.48%, z=-1.29)
Local Store	0 (+0.00%, z=+0.00)	-	-	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Surface Store	0 (+0.00%, z=+0.00)	nan (+0.00%, z=+0.00)	-	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Global Reduction	0 (+0.00%, z=+0.00)	nan (+0.00%, z=+0.00)	-	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Surface Reduction	0 (+0.00%, z=+0.00)	nan (+0.00%, z=+0.00)	-	-	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Global Atomic	0 (+0.00%, z=+0.00)	nan (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	nan (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Global Atomic Cas	0 (+0.00%, z=+0.00)	nan (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	nan (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Surface Atomic	0 (+0.00%, z=+0.00)	nan (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	nan (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Surface Atomic Cas	0 (+0.00%, z=+0.00)	nan (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	nan (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)	0 (+0.00%, z=+0.00)
Loads	-	-	85,68,68,272 (-0.03%, z=-0.28)	0 (-100.00%, z=-1.29)	27,41,97,84,704 (-0.03%, z=-0.28)	2,89,61,27,091,51 (-76.48%, z=-1.29)
Stores	26,21,440 (+0.00%, z=+0.00)	nan (+0.00%, z=+0.00)	-	-	8,38,86,080 (+0.00%, z=+0.00)	88,60,199,00 (-76.48%, z=-1.29)
Total	26,21,440 (+0.00%, z=+0.00)	nan (+0.00%, z=+0.00)	85,68,68,272 (-0.03%, z=-0.28)	nan (+0.00%, z=+0.00)	27,50,36,70,784 (-0.03%, z=-0.28)	2,90,49,87,290,50 (-76.48%, z=-1.29)

Device Memory (FB)

	L2<->FB Sectors	% Peak	Bytes	Throughput
Load	28,31,30,679 (+69.62%, z=+1.29)	2.55 (+60.37%, z=+1.30)	9,06,01,81,728 (+69.62%, z=+1.29)	95,69,52,727.37 (-49.26%, z=-1.30)
Store	2,80,14,438 (+105.18%, z=+1.31)	0.26 (+95.02%, z=+1.28)	92,84,62,016 (+105.19%, z=+1.31)	9,80,65,831.92 (-31.68%, z=-0.57)
Total	31,21,45,117 (+72.40%, z=+1.30)	2.81 (+63.07%, z=+1.31)	9,98,86,43,744 (+72.40%, z=+1.30)	1,05,50,18,559.29 (-48.01%, z=-1.29)

Scheduler Statistics

Active Warps Per Scheduler [warp/cycle]	4.00	(+0.00%, z=+0.00)	Instructions Per Active Issue Slot [inst/issue]	1.19	(-0.00%, z=+0.00)
Eligible Warps Per Scheduler [warp/cycle]	1.17	(-0.31%, z=+1.23)	No Eligible [%]	38.49	(-0.76%, z=-1.25)
Issued Warp Per Scheduler [issue/cycle]	0.70	(+0.37%, z=+1.25)	One or More Eligible [%]	69.54	(+0.37%, z=+1.25)

Warps Per Scheduler

Recommendations

CPI Stall 'Wait'

[Warning] Every scheduler is capable of issuing two instructions per cycle, but for this kernel each scheduler only issues an instruction every 1.4 cycles. This might leave hardware resources underutilized and may lead to less optimal performance. Out of the maximum of 16 warps per scheduler, this kernel allocates an average of 4.00 active warps per scheduler, but only an average of 1.17 warps were eligible per cycle. Eligible warps are the subset of active warps that are ready to issue their next instruction. Every cycle with no eligible warp results in no instruction being issued and the issue slot remains unused. To increase the number of eligible warps either increase the number of active warps or reduce the time the active warps are stalled.

Warp State Statistics

Warp Cycles Per Issued Instruction [cycle/inst]	4.27	(-0.48%, z=-1.24)	Avg. Active Threads Per Warp [thread/inst]	32	(+0.00%, z=+0.00)
Warp Cycles Per Issue Active [cycle/issue]	5.07	(-0.48%, z=-1.24)	Avg. Not Predicated Off Threads Per Warp [thread/inst]	31.83	(+0.00%, z=+0.00)
Warp Cycles Per Executed Instruction [cycle/inst]	4.27	(-0.48%, z=-1.24)	-	-	-

Warp State (All Cycles)

Recommendations

CPI Stall 'Wait'

[Warning] On average each warp of this kernel spends 1.7 cycles being stalled on a fixed latency execution dependency. This represents about 32.6% of the total average of 5.1 cycles between issuing two instructions. Typically, this stall reason should be very low and only shows up as a top contributor in already highly optimized kernels. If possible, try to further increase the number of active warps to hide the corresponding instruction latencies.

Instruction Statistics

Executed Instructions [inst]	88,59,79,75,848	(-0.00%, z=-0.39)	Avg. Executed Instructions Per Scheduler [inst]	2,01,49,49,376	(-0.00%, z=+0.00)
Issued Instructions [inst]	88,59,86,63,059	(+0.00%, z=+1.23)	Avg. Issued Instructions Per Scheduler [inst]	2,01,49,66,576.47	(+0.00%, z=+1.23)

Executed Instruction Mix

Launch Statistics

Grid Size	81,928	(+0.00%, z=+0.00)	Registers Per Thread [register/thread]	49	(+0.00%, z=+0.00)
Block Size	256	(+0.00%, z=+0.00)	Static Shared Memory Per Block [byte/block]	0	(+0.00%, z=+0.00)
Theoretical Active Warps Per SM [warp/cycle]	2,09,71,528	(+0.00%, z=+0.00)	Dynamic Shared Memory Per Block [byte/block]	36.86	(+0.00%, z=+0.00)
Waves Per SM	4,896	(+0.00%, z=+0.00)	Shared Memory Configuration Size [byte]	73.73	(+0.00%, z=+0.00)

Block Durations

Warp Durations

Occupancy

Theoretical Occupancy [%]	25	(+0.00%, z=+0.00)	Block Limit Registers [block]	4	(+0.00%, z=+0.00)
Theoretical Active Warps Per SM [warp/cycle]	25	(+0.00%, z=+0.00)	Block Limit Shared Mem [block]	2	(+0.00%, z=+0.00)
Achieved Occupancy [%]	24.99	(+0.00%, z=+0.00)	Block Limit Warps [block]	8	(+0.00%, z=+0.00)
Achieved Active Warps Per SM [warp/cycle]	16.00	(+0.00%, z=+0.00)	Block Limit SM [block]	32	(+0.00%, z=+0.00)

Impact of Varying Register Count Per Thread

Impact of Varying Register Block Size

Impact of Varying Shared Memory Usage Per Block