



First Experiences with Maxwell GPUs

October 8, 2014 | Andrew V. Adinetz



Agenda

- New in Maxwell
- Shared-memory atomics
 histograms, filtering
- PANDA Triplet Finder



Agenda

- New in Maxwell
- Shared-memory atomics
 histograms, filtering
- PANDA Triplet Finder

New in Maxwell



- Functionally same as Kepler (K20)
- Dynamic parallelism is mainstream
- Architectural improvements
 - new SM design (SMM)
 - native shared-memory 32-bit atomics
 - other improvements
- CC 5.x (5.0 for GM 207, 5.2 for GM 204)
 - nvcc -arch=sm_50 ...

Kepler SM (SMX)





- Scheduler not tied to cores
- Double issue for max utilization

Maxwell SM (SMM)



SMM



- Simplified design
 - power-of-two, quadrant-based
 - scheduler tied to cores
- Better utilization
 - single issue sufficient
 - lower instruction latency
- Efficiency
 - <10% difference from SMX</p>
 - ~50% SMX chip area



Other Maxwell Features



- More thread blocks/SM
 - 32 (vs. 16), full occupancy with 64-thread TBs
- Shared memory
 - 64 KiB/SM in GM 207, 96 KiB/SM in GM 204
 - 4-byte mode only
- L1 cache
 - read-only (same as texture cache)
 - not sharable with shared memory
- Larger L2 cache
 - 2 MiB vs. 1.5 MiB on Kepler

Parameters of Test GPUs



	Kepler K20X	Maxwell GTX 750 Ti
#SMs	14	5
#CUDA cores	2688	640
Frequency, MHz	732	1110
Peak SP GFlop/s	3935	1305.6
Memory BW, GiB/s	169.5	68.9

All results are preliminary



Agenda

New in Maxwell

Shared-memory atomics

- histograms, filtering
- PANDA Triplet Finder



Histogram

Common operation in signal and image processing

```
on CPU:
for(i = 0; i < n; i++)
  histo[data[i]]++;
on GPU:
 global void histo k
(int *histo, uchar *data, int n) {
  int i = threadIdx.x +
    blockIdx.x * blockDim.x;
  if(i >= n)
    return;
  atomicAdd(&histo[data[i]], 1);
```





Global atomics are the bottleneck



Shared-Memory Histogram

Avoid global atomics bottleneck

```
#define NCLASSES 256
#define BS 256
#define PER THREAD 32
  global void histo k(int *histo, const unsigned* data, int n) {
 // init per-block histogram
    shared int lhisto[NCLASSES];
  for(int i = threadIdx.x; i < NCLASSES; i += BS)</pre>
    lhisto[i] = 0;
    syncthreads();
  // compute per-block histogram
  int istart = blockIdx.x * (BS * PER THREAD) + threadIdx.x;
  int iend = min(istart + BS * PER THREAD, n);
  for(int i = istart; i < iend; i += BS) {</pre>
    union { unsigned char c[sizeof(unsigned)]; unsigned i; } el;
    el.i = data[i];
    for (int j = 0; j < size of (unsigned); j++)
      atomicAdd(&lhisto[el.c[j]], 1); // shared-memory atomic
    syncthreads();
  // accumulate histogram to global storage
  for(int i = threadIdx.x; i < NCLASSES; i += BS)</pre>
    atomicAdd(&histo[i], lhisto[i]); // global atomics
} // histo kernel
```

Histogram256 Performance



- Maxwell clearly wins (against sever-grade GPUs!)
- 60% memory bandwidth achieved

October 8, 2014

Performance per SM





Higher performance expected with larger GPUs (more SMs)

October 8, 2014

NVidia Application Lab Workshop 2014



Filtering

Copying elements based on predicate (e.g., > 0)

1	-4	2	-3	3	-2	5	0	7			3	1	2	7	5
---	----	---	----	---	----	---	---	---	--	--	---	---	---	---	---

```
on CPU: nres = 0;
for(i = 0; i < n; i++)
    if(src[i] > 0)
        tgt[nres++] = src[i];
```

```
on GPU:
```

```
global____void filter_k(int *tgt, int *nres, int *src, int n) {
    int i = threadIdx.x + blockIdx.x * blockDim.x;
    if(i >= n)
        return;
    if(src[i] > 0)
        tgt[atomicAdd(nres, 1)] = src[i];
}
```

Global atomics are the bottleneck

October 8, 2014

Shared vs Global Atomics





Qualitatively different behavior on Maxwell

Warp-Aggregated Atomics



http://devblogs.nvidia.com/parallelforall/cuda-pro-tip-optimized-filtering-warp-aggregated-atomics/

Higher shared-memory performance with more SMs expected

LICH



Agenda

- New in Maxwell
- Shared-memory atomics
 histograms, filtering
- PANDA Triplet Finder

PANDA Experiment

- state-of-the-art triggerless HEP experiment
- starts in 2019
- latest GPUs are of interest (e.g., Maxwell)
 - single precision OK







p a n)d a



Triplet Finder Maxwell Tests

- Optimal bunch size
- Bunching strategies
 - dynamic parallelism vs. host streams
- Specific optimizations
 - SoA vs. AoS
 - all-tube testing vs. sector-rows
 - all-skewlet testing vs. binning
 - optimal #connections
- Maxwell vs. Kepler comparison

Optimal Bunch Size





- Kepler: 1 (dynamic parallelism) and 2 μs
- more similar behavior between approaches
- possibly due to smaller GPUs

Number of Connections





- more similar behavior between #connections
- possibly due to smaller GPUs
- 24 connections used (as for Kepler)

Bunching Strategy





- more similar behavior between approaches
- possibly due to smaller GPUs



SoA vs AoS

array-of-structures (AoS)

struct fvec3 {

- float x;
 float y;
 float z;
- } a[N];

xyzxyz . . . xyz

better language support

structure-of-arrays (SoA)

float ax[N];
float ay[N];
float az[N];



- better coalescing
- better performance



Automating SoA

• R. Strzodka, <u>http://people.mpi-inf.mpg.de/~strzodka/projects/layout/</u>





~ 20% improvement with SoA, similar to Kepler



Tube Testing



1.5x improvement, less than Kepler (1.7x)





- >30% time taken by skewlet testing
- like sector-rows, but with skewlets
- bin assignment required, different statistics

Skewlet Binning Performance





- similar performance
- Kepler: binning 2x slower



Maxwell vs. Kepler



- Kepler faster when saturated
- Maxwell faster on smaller #hits (< 40K)
 - higher frequency?



Maxwell vs. Kepler (per-SM)

Adjusted for frequency



SMM vs. SMX performance difference within 10%
 ~20% error range for Maxwell

Runtime Breakdown





- similar runtimes before track evaluation
- skewlet testing dominates even more for Maxwell

Conclusions



- Very promising first experiences
 - histogram faster than on K20X
 - PANDA: good performance of some stages
- Larger Maxwell GPU wanted
 - smaller GPUs: not all effects are clear
 - = GM 204 or server-grade

Questions?

• NVLab at FZJ: <u>http://fz-juelich.de/ias/jsc/nvlab</u>



JÜLICH APPLICATION LAB NVIDIA

