



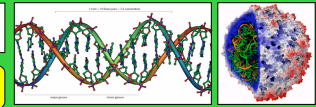
Accelerating Applications with Inter-Block GPU Communication via Fast Barrier Synchronization

SyNeRG
http://synergy.cs.vt.edu/

Wu Feng and Shucai Xiao

{wfeng, shucai}@vt.edu

Virginia Tech



1. Motivation and Background

GPGPU: General-Purpose Computation for GPU

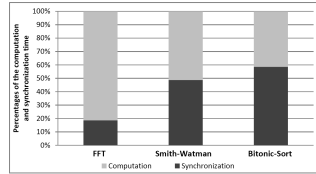
- The need to move from *data (or task) parallel computation* to *general-purpose computation*?

Current State of the Art in GPU Computing

- Focus on accelerating computation rather than inter-block GPU communication (via barrier synchronization).

- Inter-block GPU communication for more general-purpose computation consumes a substantial portion of total execution time.

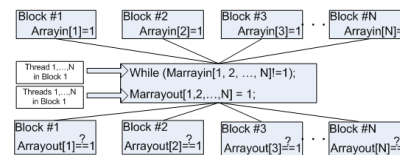
- Current Approach**
Barrier synchronization across multiprocessors via kernel launch.



3. Lock-Free Approach

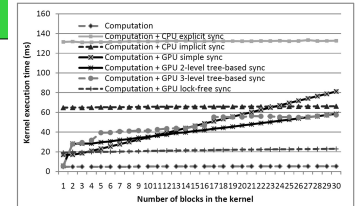
Communication via GPU Lock-Free Synchronization

- Implementation:** Two arrays control inter-block GPU communication. Controlling threads in "Block #1" synchronize via `__syncthreads()`.
- Properties**
 - No need for atomic add
 - All operations executed in parallel
 - Constant synchronization time

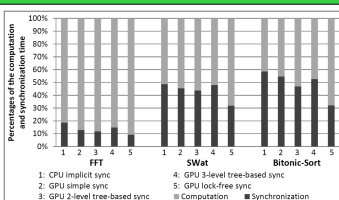


4. Performance Verification of GPU Communication via Barrier Synchronization

- CPU implicit synchronization takes *12 times longer* than the computation time.
- GPU-based barrier synchronization takes less time than both CPU implicit and explicit synchronization.
 - The execution time of GPU simple synchronization (GSS) varies linearly with the number of blocks.
 - The execution time of GPU lock-free synchronization is constant.
- The empirical results match our performance models.



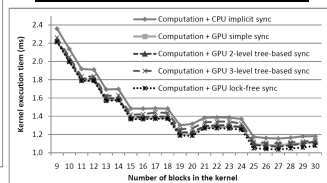
5. Performance Evaluation of GPU Communication: FFT, Smith-Watman, Bitonic Sort



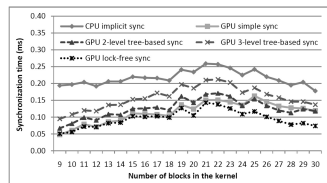
Percentage of Computation Time vs. Synchronization Time

- Bitonic sort experiences a 38% improvement, Smith-Watman 24%, and FFT 8%.
- The percentage of time spent in synchronization is smallest for FFT and largest for bitonic sort.
- The variances in synchronization time are consistent with our performance models.

Fast Fourier Transformation

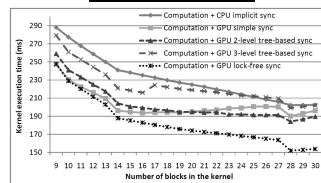


Execution Time (ms) vs. Block Number

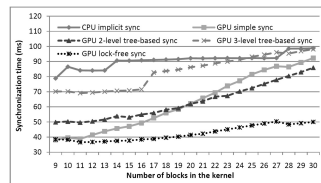


Synchronization Time (ms) vs. Block Number

Smith-Watman

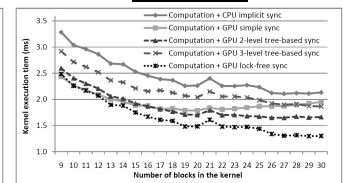


Execution Time (ms) vs. Block Number

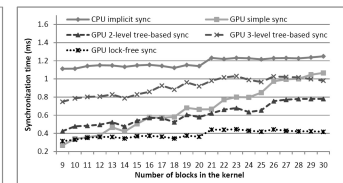


Synchronization Time (ms) vs. Block Number

Bitonic Sort



Execution Time (ms) vs. Block Number



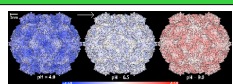
Synchronization Time (ms) vs. Block Number

6. Related Work: Data-Parallel Apps on the GPU

N-Body Molecular Dynamics: 5000x speed-up

CPU Version: 1,334 minutes.

GPU Version: 0.26 minutes = 16 seconds.



7. Acknowledgements



VirginiaTech
Invent the Future