



Advanced Topics in CUDA

Cliff Woolley, NVIDIA
Developer Technology Group



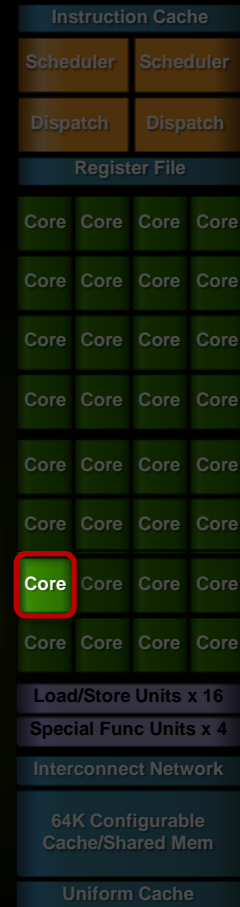
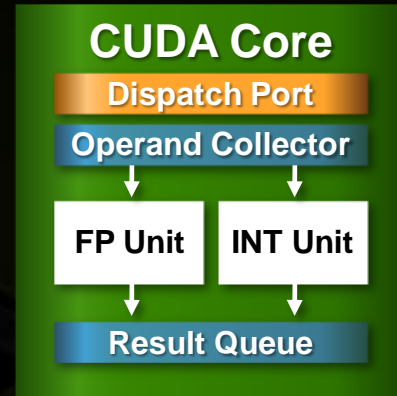


RECAP: SCHEDULING

GPU Architecture – Fermi: CUDA Core



- Floating point & Integer unit
 - IEEE 754-2008 floating-point standard
 - Fused multiply-add (FMA) instruction for both single and double precision
- Logic unit
- Move, compare unit
- Branch unit

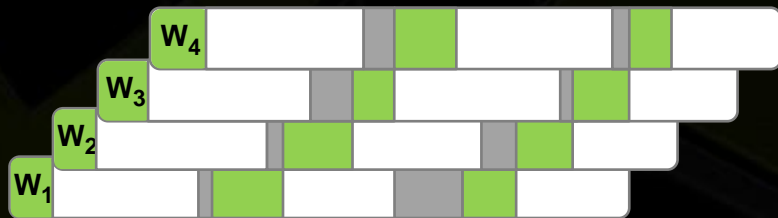


Low Latency or High Throughput?



- CPU architecture must **minimize latency** within each thread
- GPU architecture **hides latency** with computation from other thread warps

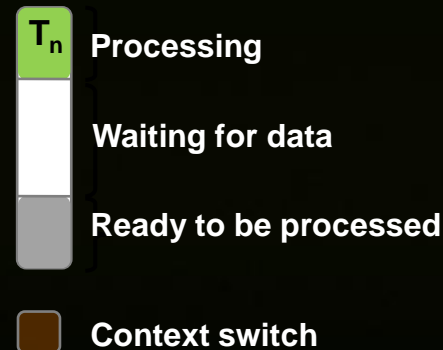
GPU Stream Multiprocessor – High Throughput Processor



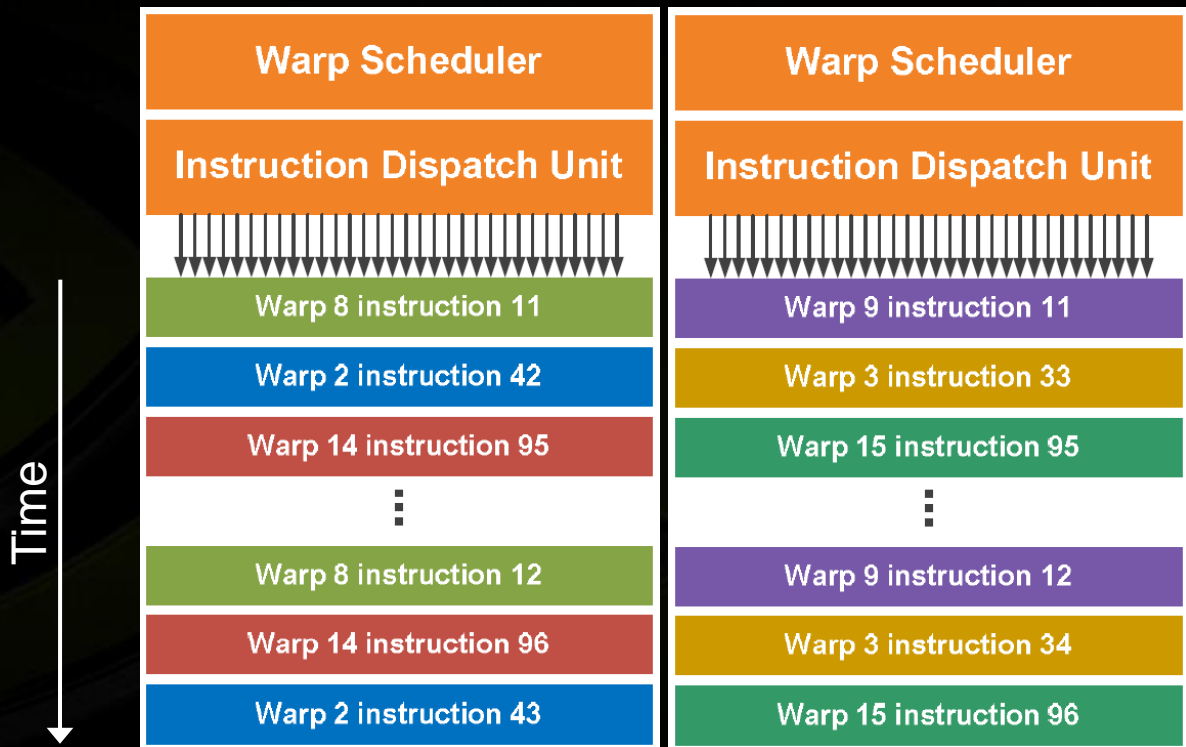
CPU core – Low Latency Processor



Computation Thread/Warp



GPU Latency Hiding: Warp Switching

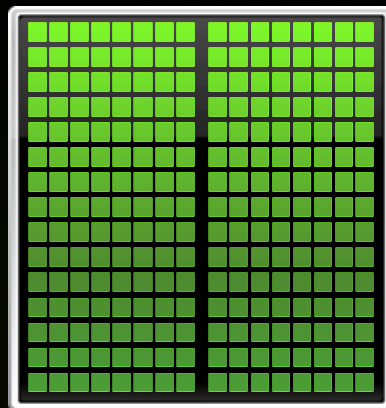


OVERLAPPING PROCESSING WITH DATA TRANSFERS

Available engines



CPU

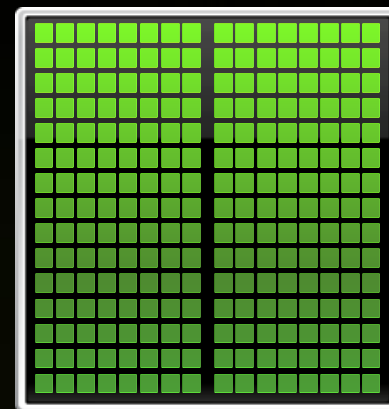
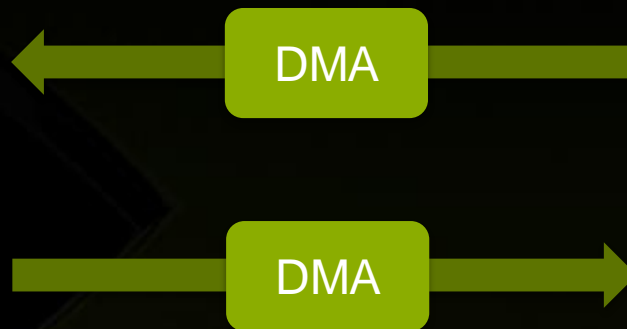


GPU

Available engines

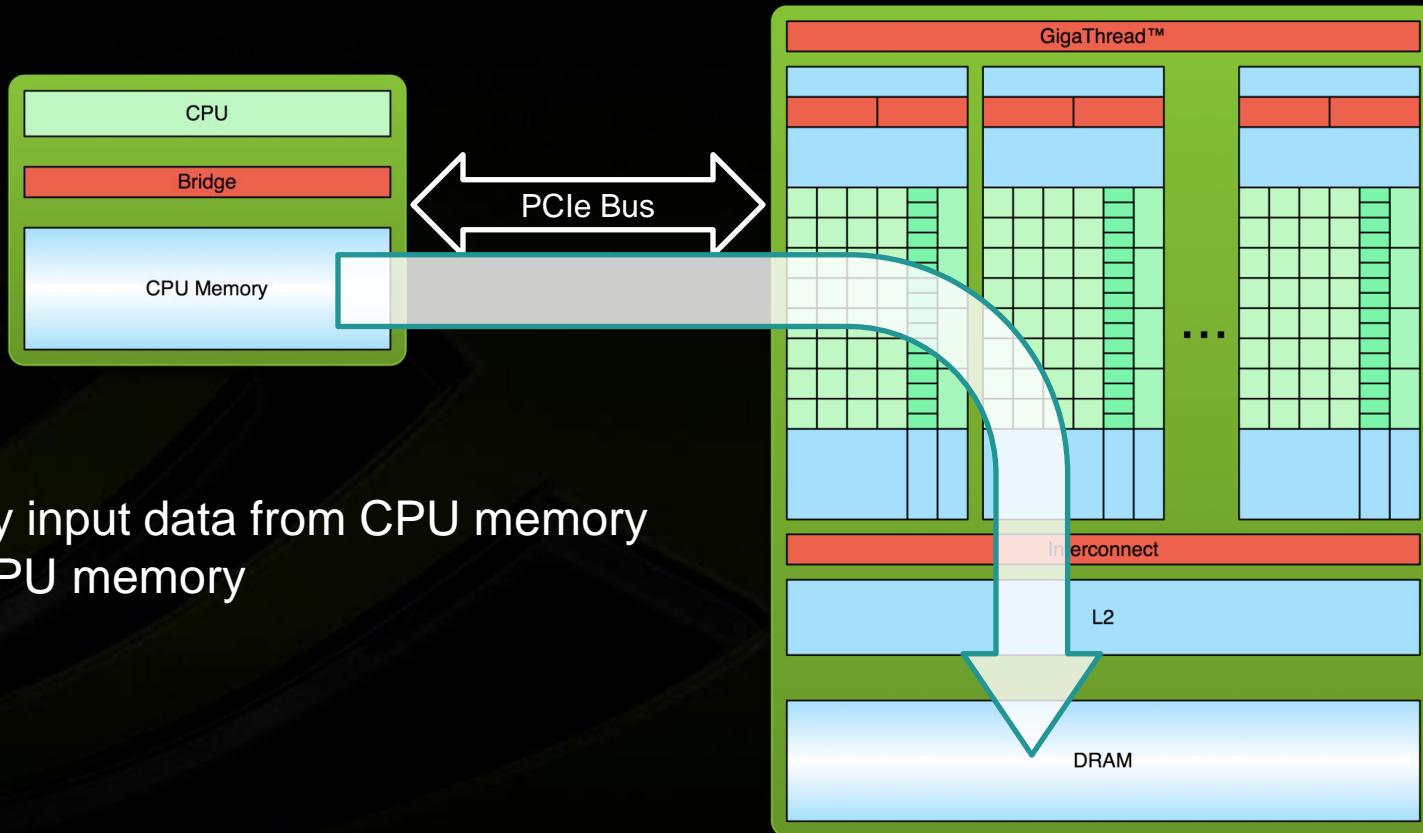


CPU



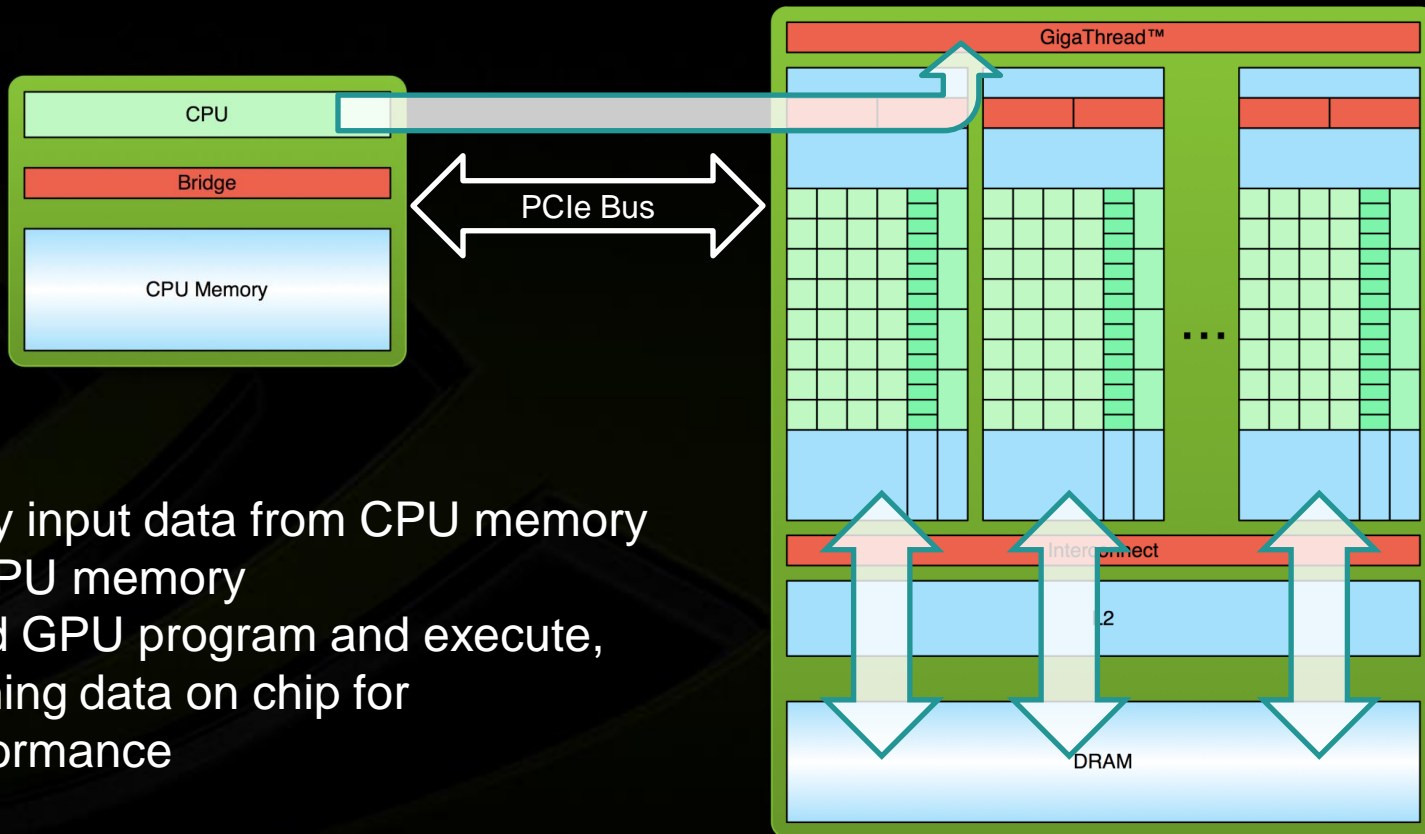
GPU

Standard Processing Flow



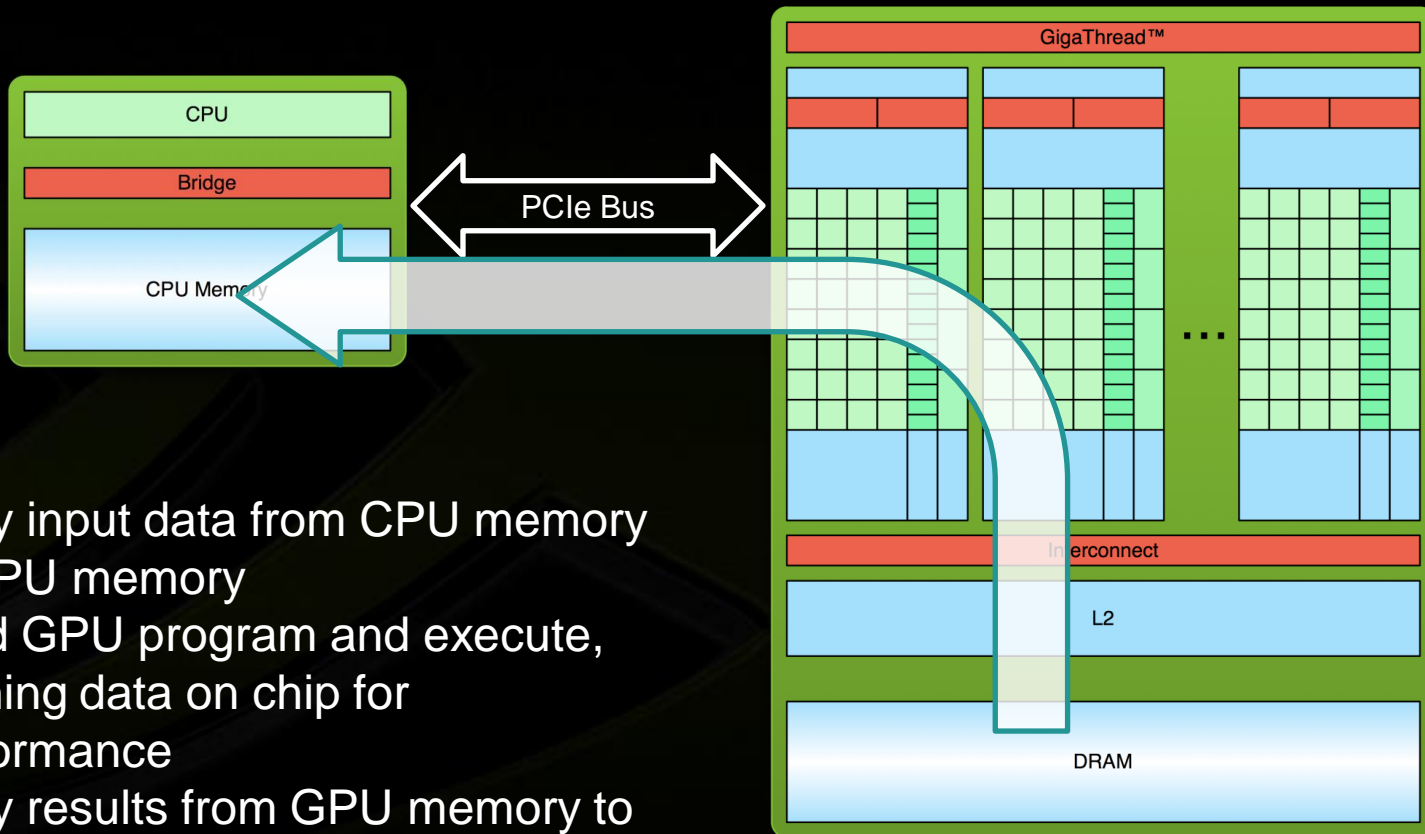
1. Copy input data from CPU memory to GPU memory

Standard Processing Flow



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance

Standard Processing Flow



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance
3. Copy results from GPU memory to CPU memory

Objective

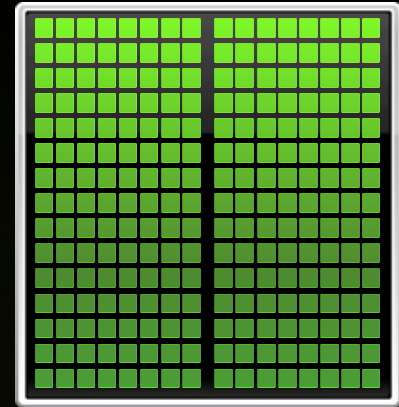


- Tasks not on the critical path should be hidden
 - i.e. overlapped with other tasks
- Critical path is frequently the transfer of result data from GPU to CPU
 - e.g. path data

Objective: Overlap Processing With Transfers



CPU

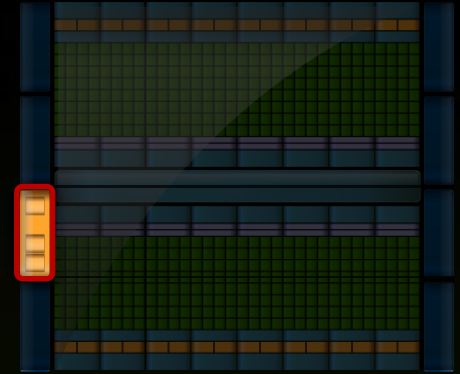
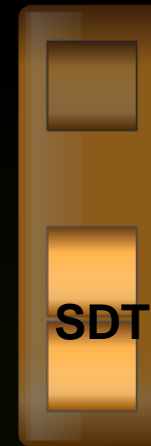


GPU

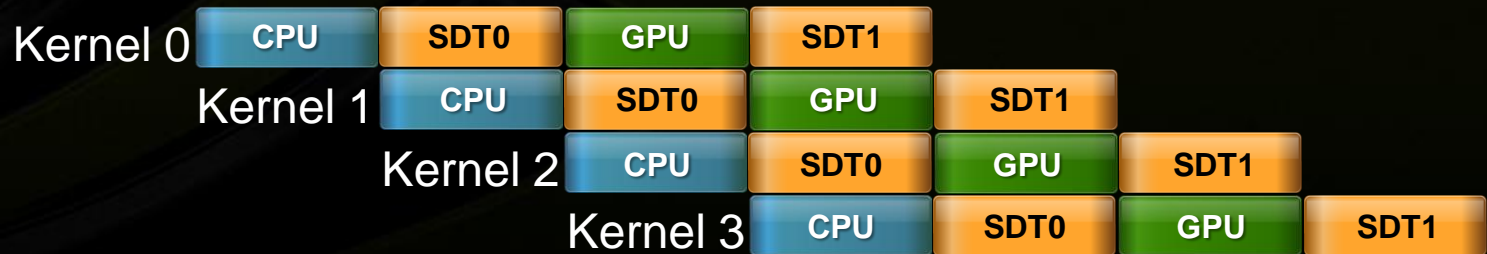
Overlapping Processing With Transfers



- **Dual DMA engines**
 - Simultaneous CPU→GPU and GPU→CPU data transfer
 - Fully overlapped with CPU and GPU processing time



- **Activity Snapshot:**



Phase A: No Pipelining



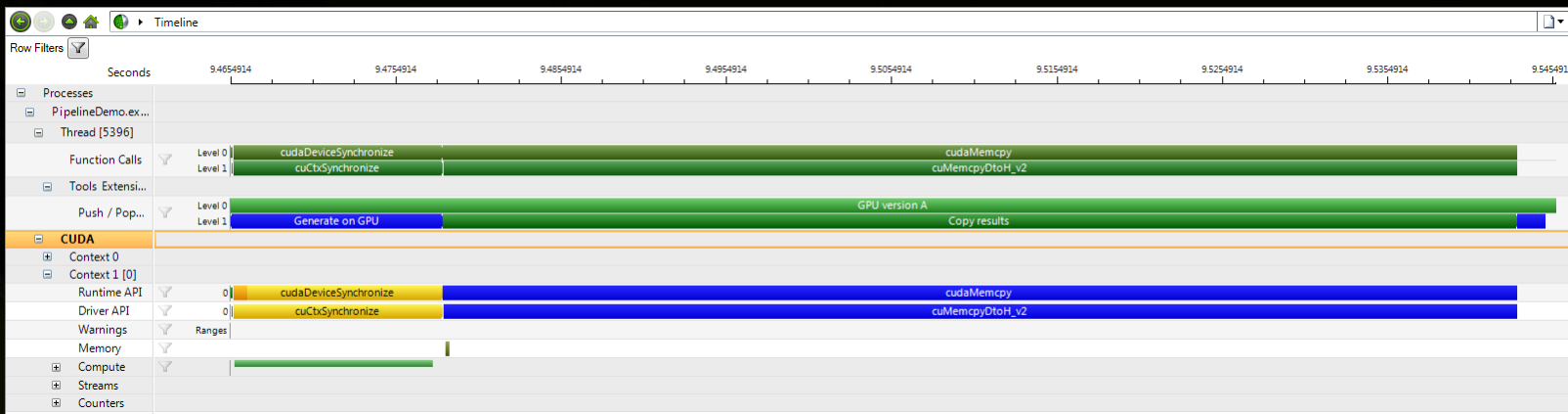
Kernel

Transfer

Post process

- No overlap
 - Kernel runs on GPU
 - Transfer data across PCIe from device to host
 - Post process on CPU

Phase A: No Pipelining

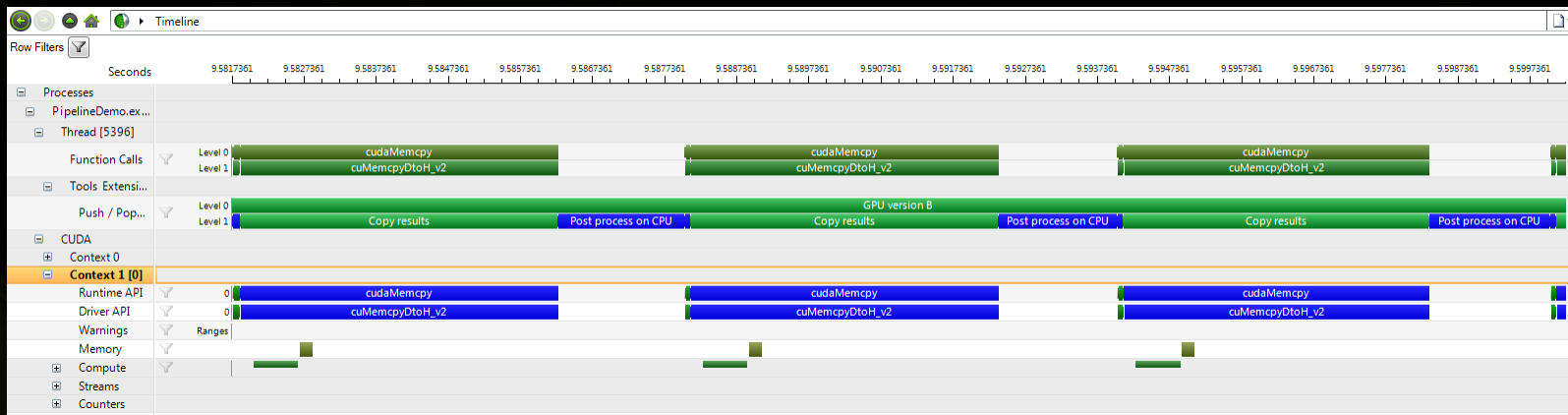


Phase B: Batching

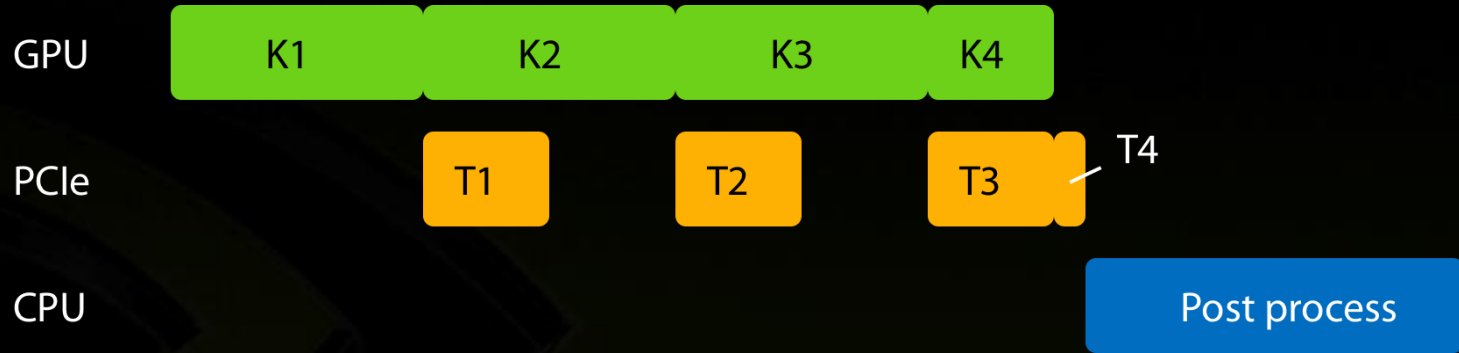


- No overlap
 - Preparation for next phase...

Phase B: Batching

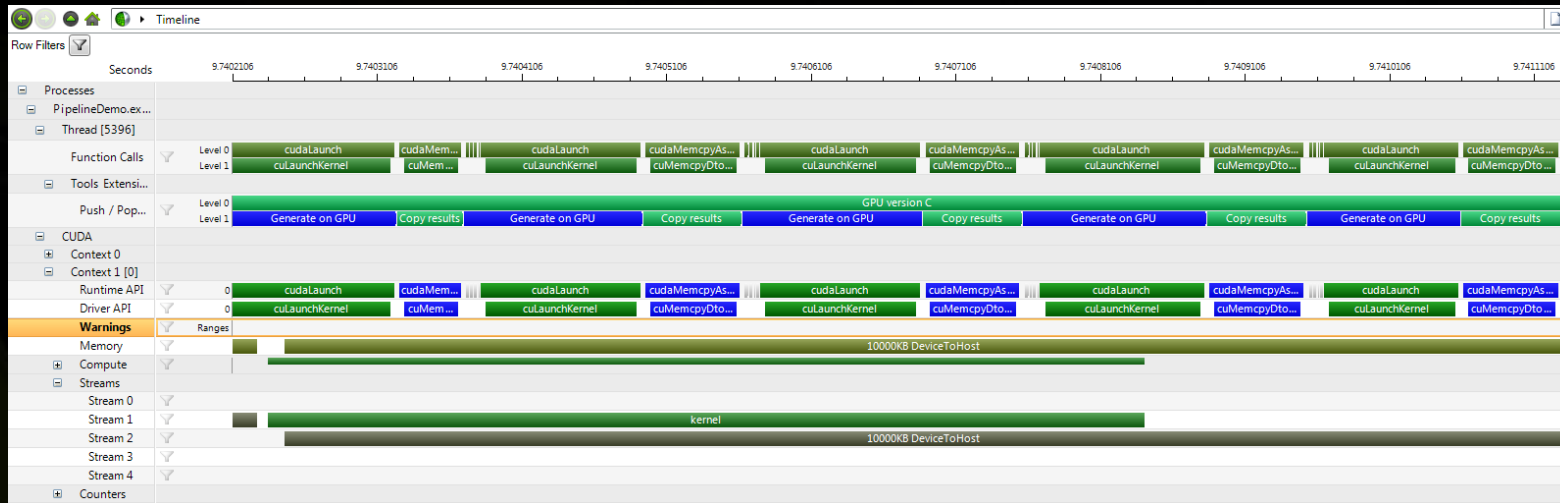


Phase C: Overlap Kernel and Transfer

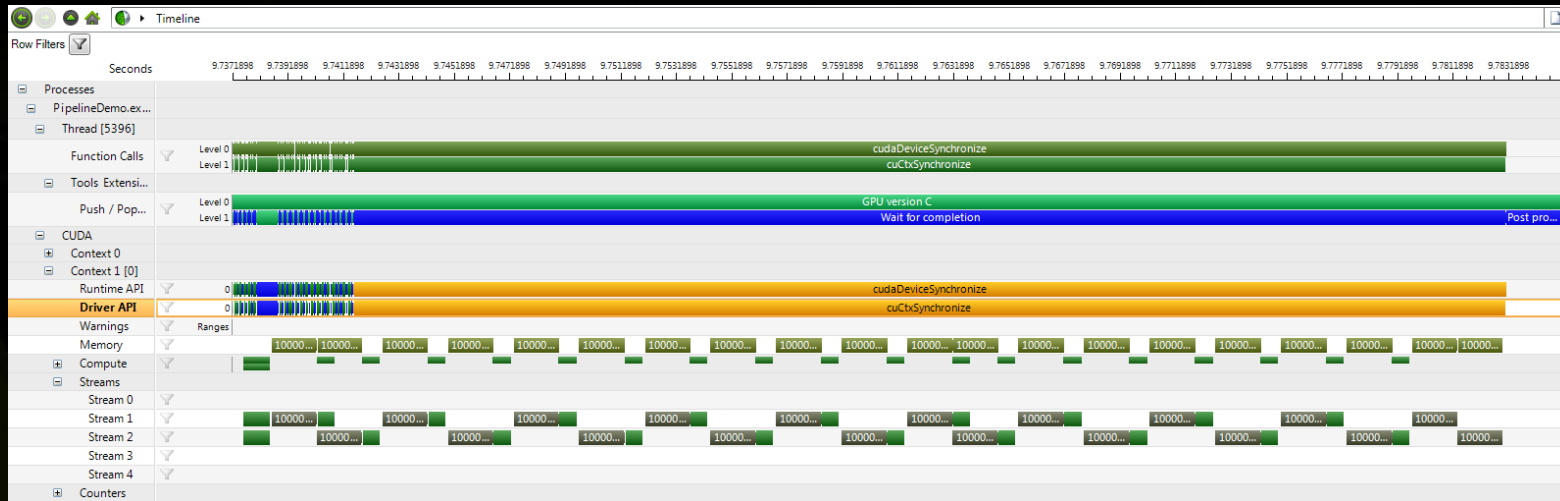


- Overlap GPU compute with PCIe transfer
- Use streams to specify dependencies
 - K1→T1, K2→T2 etc.

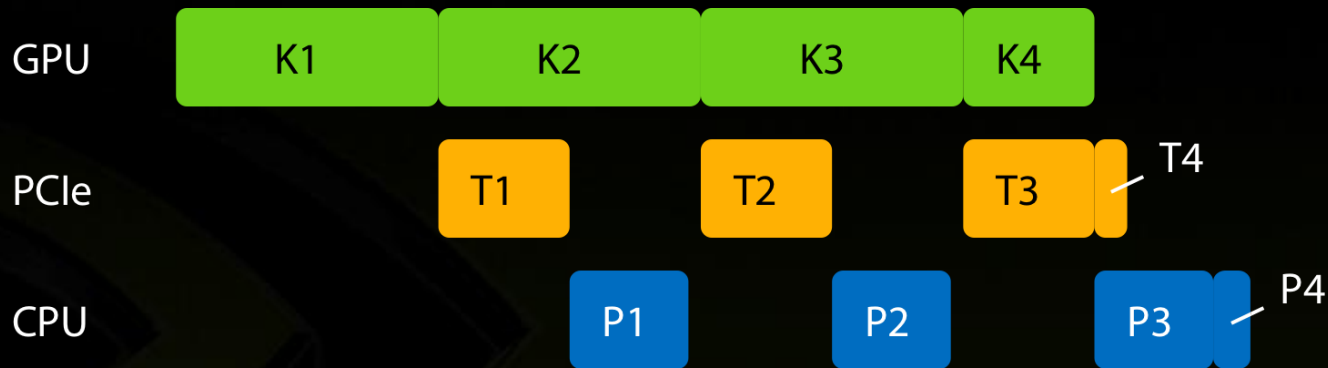
Phase C: Overlap Kernel and Transfer



Phase C: Overlap Kernel and Transfer

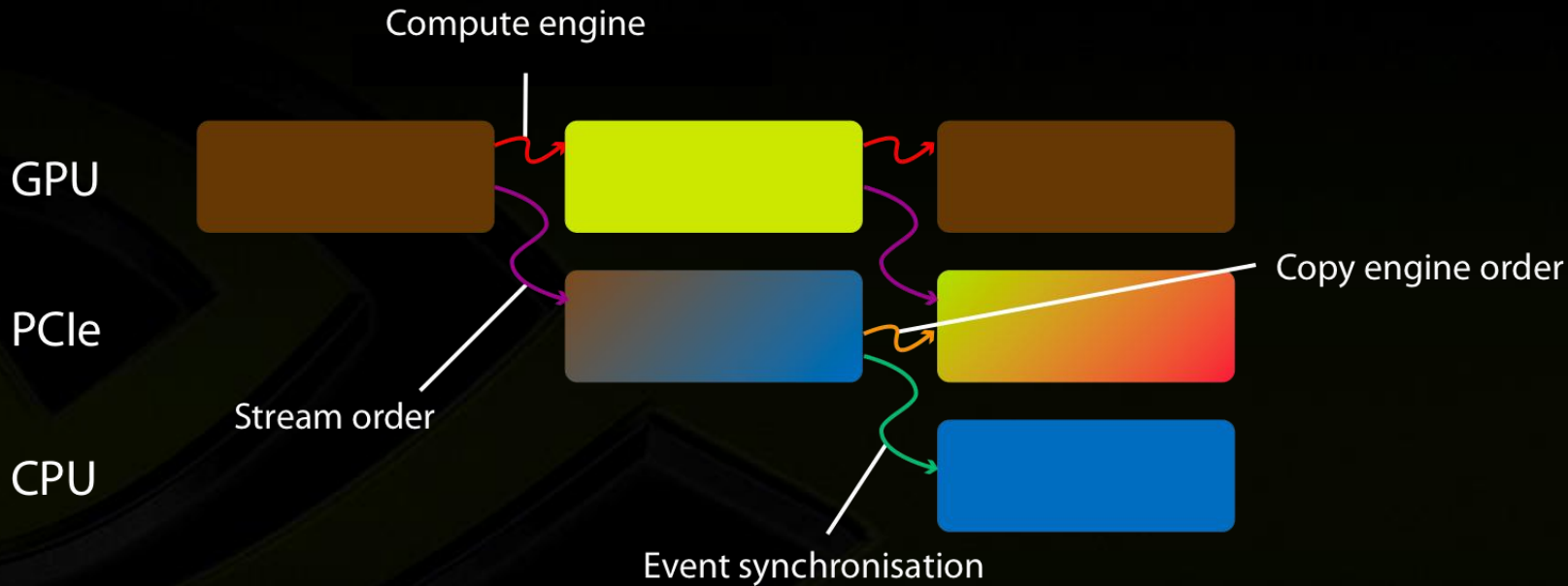


Phase D: Overlap Kernel, Transfer and CPU

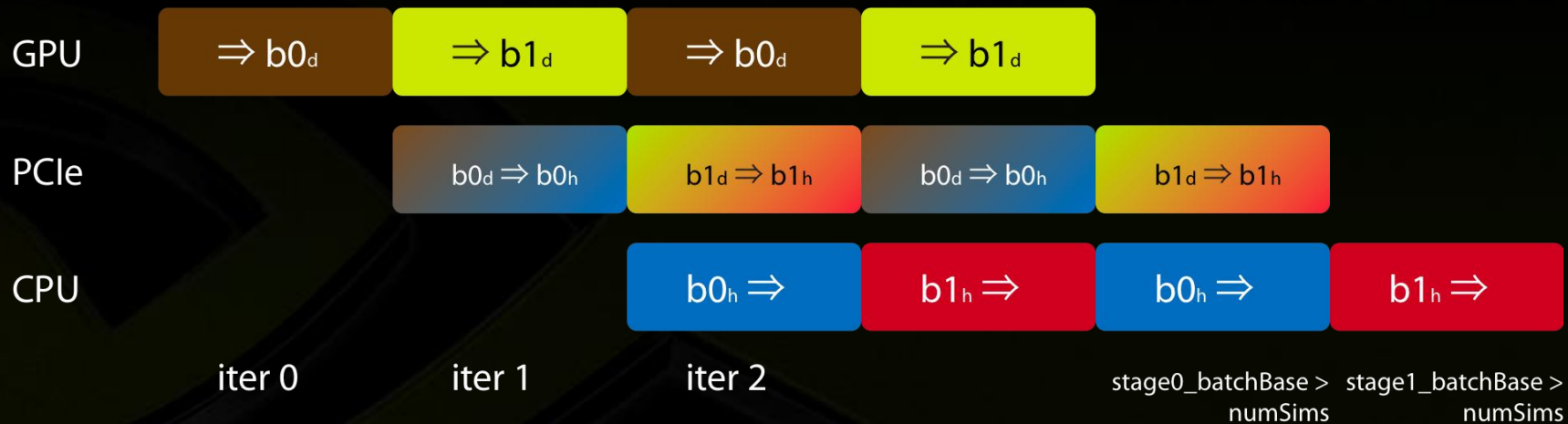


- Overlap GPU compute, PCIe transfer and CPU compute
- Use streams to specify K/T dependencies
- Use events to specify T/P dependencies
 - $T1 \rightarrow P1$, $T2 \rightarrow P2$ etc.

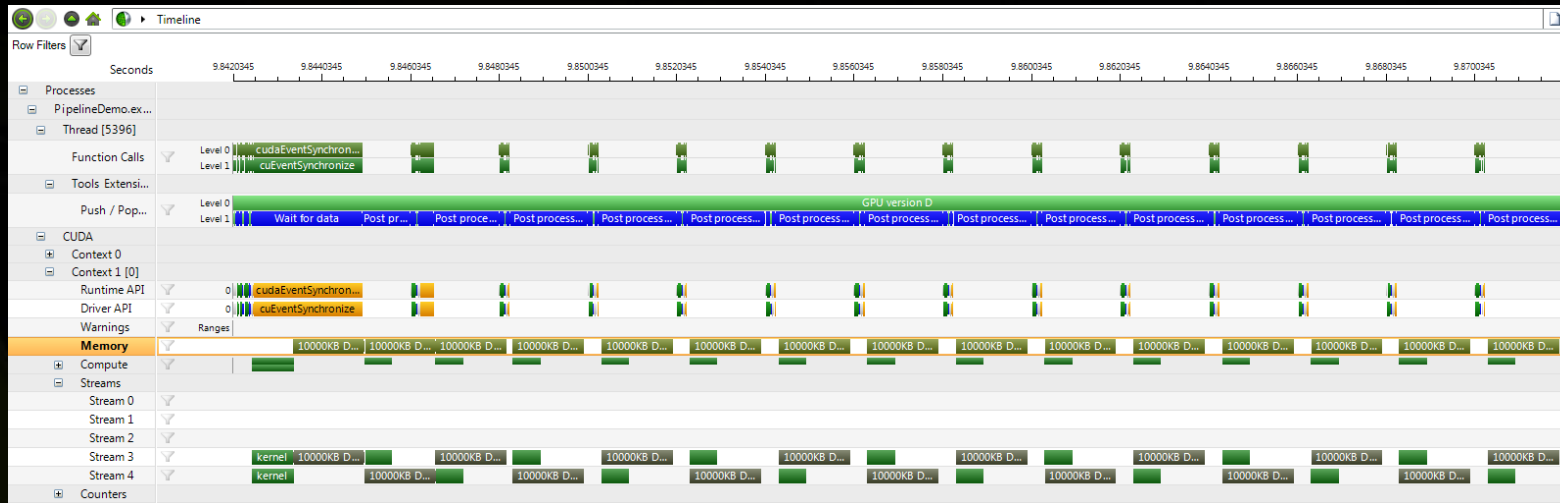
Phase D: Work Dependencies



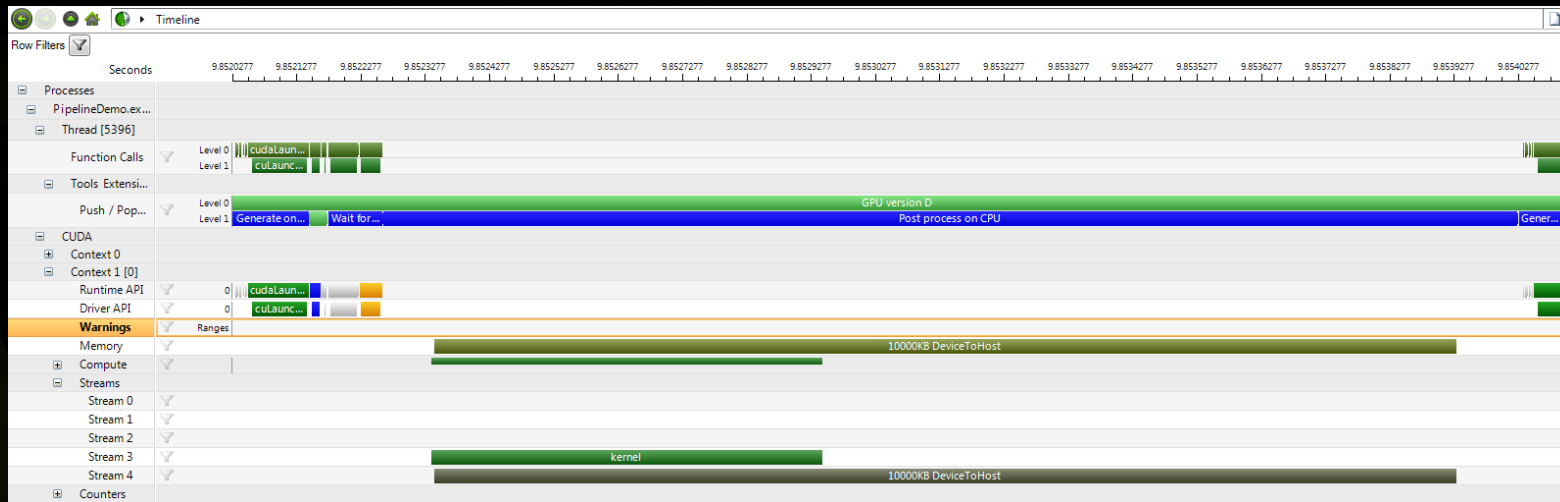
Phase D: Use Double Buffering



Phase D: Overlap Kernel, Transfer and CPU



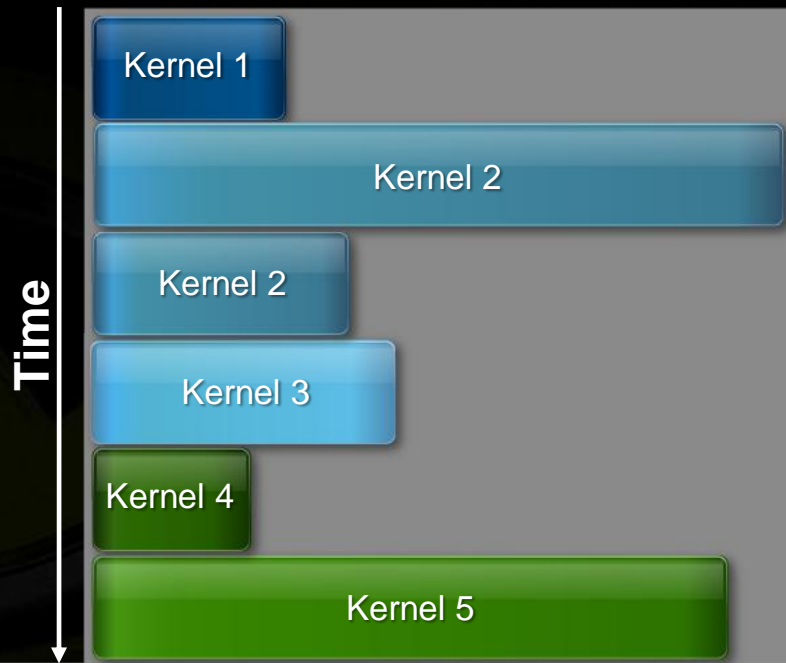
Phase D: Overlap Kernel, Transfer and CPU



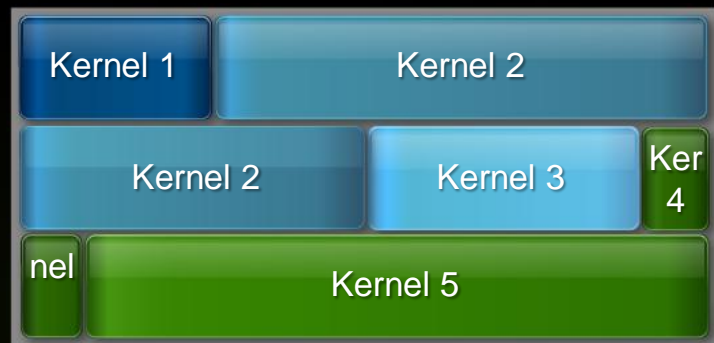


CONCURRENT KERNELS

Concurrent Kernel Execution



Sequential Kernel Execution



Parallel Kernel Execution



ANALYSIS-DRIVEN OPTIMIZATION

Performance Optimization Process



- **Use appropriate performance metric for each kernel**
 - For example, Gflops/s don't make sense for a bandwidth-bound kernel
- **Determine what limits kernel performance**
 - Memory throughput
 - Instruction throughput
 - Latency
 - Combination of the above
- **Address the limiters in the order of importance**
 - Determine how close to the HW limits the resource is being used
 - Analyze for possible inefficiencies
 - Apply optimizations
 - Often these will just fall out from how HW operates

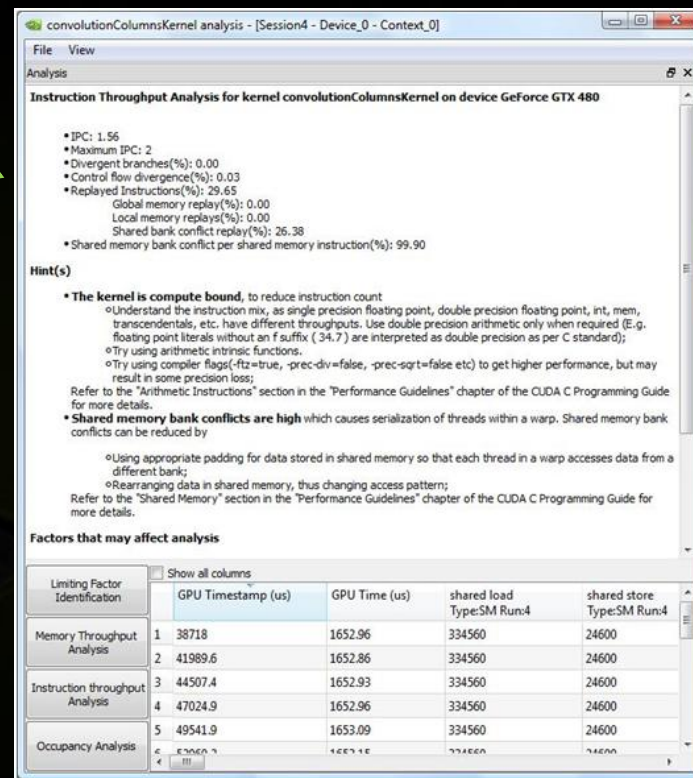
New in CUDA Toolkit 4.0: Automated Performance Analysis using Visual Profiler

Summary analysis & hints

- Per-Session
- Per-Device
- Per-Context
- Per-Kernel

New UI for kernel analysis

- Identify the limiting factor
- Analyze instruction throughput
- Analyze memory throughput
- Analyze kernel occupancy



convolutionColumnsKernel analysis - [Session4 - Device_0 - Context_0]

File View

Analysis

Instruction Throughput Analysis for kernel convolutionColumnsKernel on device GeForce GTX 480

- IPC: 1.56
- Maximum IPC: 2
- Divergent branches(%): 0.00
- Control flow divergence(%): 0.03
- Replayed Instructions(%): 29.65
 - Global memory replay(%): 0.00
 - Local memory replays(%): 0.00
 - Shared bank conflict replay(%): 26.38
- Shared memory bank conflict per shared memory instruction(%): 99.90

Hint(s)

- **The kernel is compute bound.** to reduce instruction count
 - Understand the instruction mix, as single precision floating point, double precision floating point, int, mem, transcendentals, etc. have different throughputs. Use double precision arithmetic only when required (E.g. floating point literals without an f suffix (34.7) are interpreted as double precision as per C standard);
 - Try using arithmetic intrinsic functions.
 - Try using compiler flags (-ftz=true, -prec-div=false, -prec-sqrt=false etc) to get higher performance, but may result in some precision loss;
 Refer to the "Arithmetic Instructions" section in the "Performance Guidelines" chapter of the CUDA C Programming Guide for more details.
- **Shared memory bank conflicts are high** which causes serialization of threads within a warp. Shared memory bank conflicts can be reduced by
 - Using appropriate padding for data stored in shared memory so that each thread in a warp accesses data from a different bank;
 - Rearranging data in shared memory, thus changing access pattern;
 Refer to the "Shared Memory" section in the "Performance Guidelines" chapter of the CUDA C Programming Guide for more details.

Factors that may affect analysis

Show all columns

Limiting Factor Identification	GPU Timestamp (us)	GPU Time (us)	shared load Type:SM Run:4	shared store Type:SM Run:4
Memory Throughput Analysis	1 38718	1652.96	334560	24600
	2 41989.6	1652.86	334560	24600
Instruction throughput Analysis	3 44507.4	1652.93	334560	24600
	4 47024.9	1652.96	334560	24600
Occupancy Analysis	5 49541.9	1653.09	334560	24600

Notes on profiler



- **Most counters are reported per Streaming Multiprocessor (SM)**
 - Not entire GPU
 - Exceptions: L2 and DRAM counters
- **A single run can collect a few counters**
 - Multiple runs are needed when profiling more counters
 - Done automatically by the Visual Profiler
 - Have to be done manually using command-line profiler
- **Counter values may not be exactly the same for repeated runs**
 - Threadblocks and warps are scheduled at run-time
 - So, “two counters being equal” usually means “two counters within a small delta”
- **See the profiler documentation for more information**

ANALYSIS-DRIVEN OPTIMIZATION: IDENTIFYING PERF LIMITERS

Limited by Bandwidth or Arithmetic?



- **Perfect instructions:bytes ratio for Fermi C2050:**
 - ~4.5 : 1 with ECC on
 - ~3.6 : 1 with ECC off
 - These assume fp32 instructions, throughput for other instructions varies
- **Algorithmic analysis:**
 - Rough estimate of arithmetic to bytes ratio
- **Code likely uses more instructions and bytes than algorithm analysis suggests:**
 - Instructions for loop control, pointer math, etc.
 - Address pattern may result in more memory fetches
 - Two ways to investigate:
 - Use the profiler (quick, but approximate)
 - Use source code modification (more accurate, more work intensive)

A Note on Counting Global Memory Accesses



- **Load/store instruction count can be lower than the number of actual memory transactions**
 - Address pattern, different word sizes
- **Counting requests from L1 to the rest of the memory system makes the most sense**
 - Caching-loads: count L1 misses
 - Non-caching loads and stores: count L2 read requests
 - Note that L2 counters are for the entire chip, L1 counters are per SM
- **Some shortcuts, assuming “coalesced” address patterns:**
 - **One 32-bit** access instruction -> **one 128-byte** transaction per warp
 - **One 64-bit** access instruction -> **two 128-byte** transactions per warp
 - **One 128-bit** access instruction -> **four 128-byte** transactions per warp

Analysis with Profiler



- **Profiler counters:**

- **instructions_issued, instructions_executed**
 - Both incremented by 1 per warp
 - “issued” includes replays, “executed” does not
- **gld_request, gst_request**
 - Incremented by 1 per warp for each load/store instruction
 - Instruction may be counted if it is “predicated out”
- **l1_global_load_miss, l1_global_load_hit, global_store_transaction**
 - Incremented by 1 per L1 line (line is 128B)
- **uncached_global_load_transaction**
 - Incremented by 1 per group of 1, 2, 3, or 4 transactions
 - Better to look at **L2_read_request** counter (incremented by 1 per 32B transaction; per GPU, not per SM)

- **Compare:**

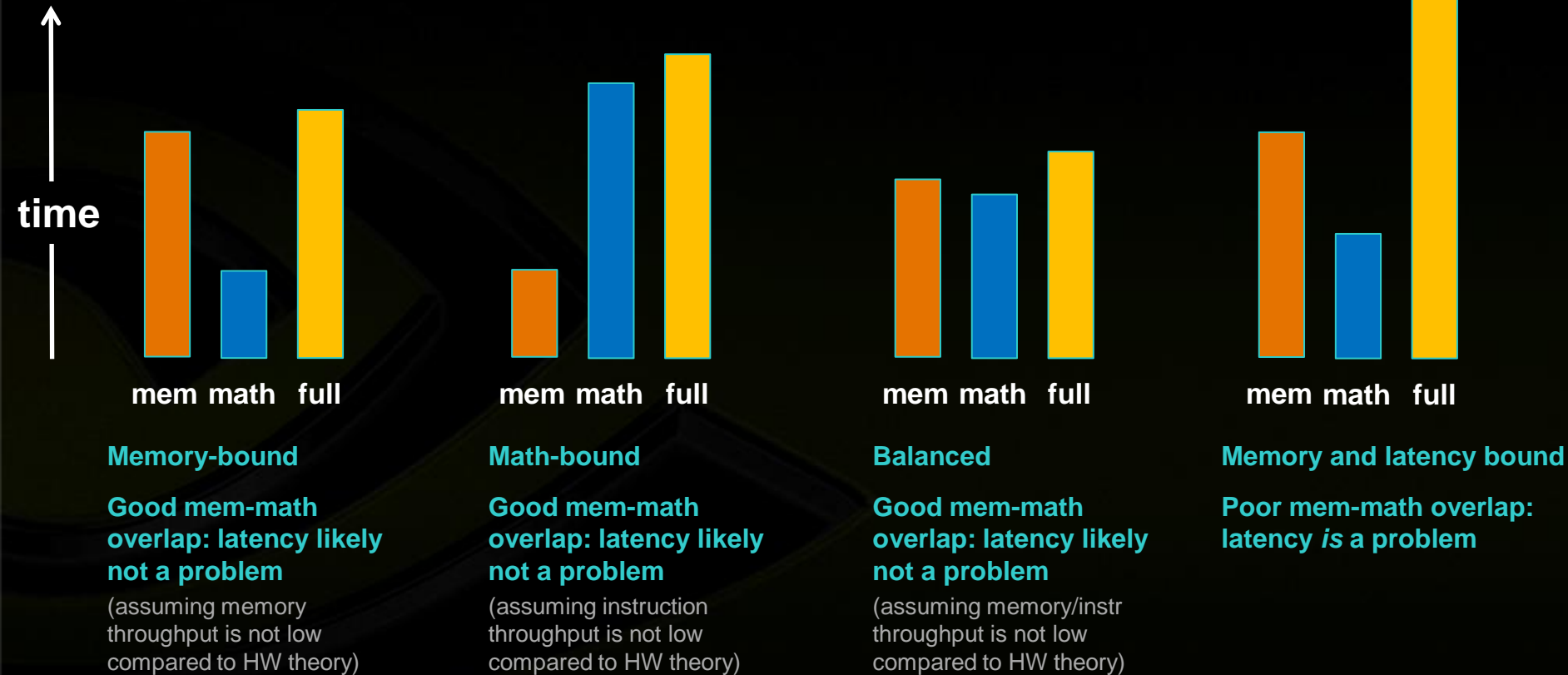
- $32 * \text{instructions_issued}$ /* 32 = warp size */
- $128B * (\text{global_store_transaction} + \text{l1_global_load_miss})$

Analysis with Modified Source Code



- **Time memory-only and math-only versions of the kernel**
 - Easier for codes that don't have data-dependent control-flow or addressing
 - Gives you good estimates for:
 - Time spent accessing memory
 - Time spent in executing instructions
- **Compare the times taken by the modified kernels**
 - Helps decide whether the kernel is mem or math bound
- **Compare the sum of mem-only and math-only times to full-kernel time**
 - Shows how well memory operations are overlapped with arithmetic
 - Can reveal latency bottleneck

Some Example Scenarios



Source Modification



- **Memory-only:**
 - Remove as much arithmetic as possible
 - Without changing access pattern
 - Use the profiler to verify that load/store instruction count is the same
- **Store-only:**
 - Also remove the loads (to compare read time vs. write time)
- **Math-only:**
 - Remove global memory accesses
 - Need to trick the compiler:
 - Compiler throws away all code that it detects as not contributing to stores
 - Put stores inside conditionals that always evaluate to false
 - Condition should depend on the value about to be stored (prevents other optimizations)
 - Condition outcome should not be known to the compiler

Source Modification for Math-only



```
__global__ void fwd_3D( ..., int flag)
{
    ...
    value = temp + coeff * vsq;
    if( 1 == value * flag )
        g_output[out_idx] = value;
}
```

If you compare only the flag, the compiler may move the computation into the conditional as well



Source Modification and Occupancy



- **Removing pieces of code is likely to affect register count**
 - This could increase occupancy, skewing the results
 - See slide 23 to see how that could affect throughput
- **Make sure to keep the same occupancy**
 - Check the occupancy with profiler before modifications
 - After modifications, if necessary add shared memory to match the unmodified kernel's occupancy

```
kernel<<< grid, block, smem, ...>>>(...
```

Case Study: Limiter Analysis



- 3DFD of the wave equation, fp32

- Time (ms):

- Full-kernel: 35.39
- Mem-only: 33.27
- Math-only: 16.25

- Instructions issued:

- Full-kernel: 18,194,139
- Mem-only: 7,497,296
- Math-only: 16,839,792

- Memory access transactions:

- Full-kernel: 1,708,032
- Mem-only: 1,708,032
- Math-only: 0

- Analysis:

- Instr:byte ratio = ~2.66
 - $32 * 18,194,139 / 128 * 1,708,032$
- Good overlap between math and mem:
 - 2.12 ms of math-only time (13%) are not overlapped with mem
- App memory throughput: 62 GB/s
 - HW theory is 114 GB/s, so we're off

- Conclusion:

- Code is memory-bound
- Latency could be an issue too
- Optimizations should focus on memory throughput first
 - math contributes very little to total time (2.12 out of 35.39ms)

Summary: Limiter Analysis



- **Rough algorithmic analysis:**
 - How many bytes needed, how many instructions
- **Profiler analysis:**
 - Instruction count, memory request/transaction count
- **Analysis with source modification:**
 - Memory-only version of the kernel
 - Math-only version of the kernel
 - Examine how these times relate and overlap

Questions?

