



Introduction to Streams

Univ. of Pennsylvania CIS 565, Nov 2020

Slides by Steve Rennich

Presented by Tim Kaldewey tkaldewey@nvidia.com

NVIDIA - Developer Technology



Streams and Concurrency



- Concurrency
 - The ability to perform multiple operations simultaneously
 - Compute kernels on the GPU
 - Data transfer to device (H2D)
 - Data transfer to host (D2H)
 - Operations on the CPU
 - Enables improved performance
- Streams
 - How concurrency is achieved

Simple Accelerator Model

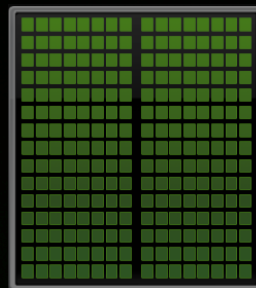


- Tasks are offloaded from CPU to GPU

CPU



GPU



Simple Accelerator Model

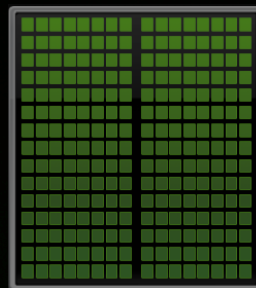


- Tasks are offloaded from CPU to GPU

CPU



GPU

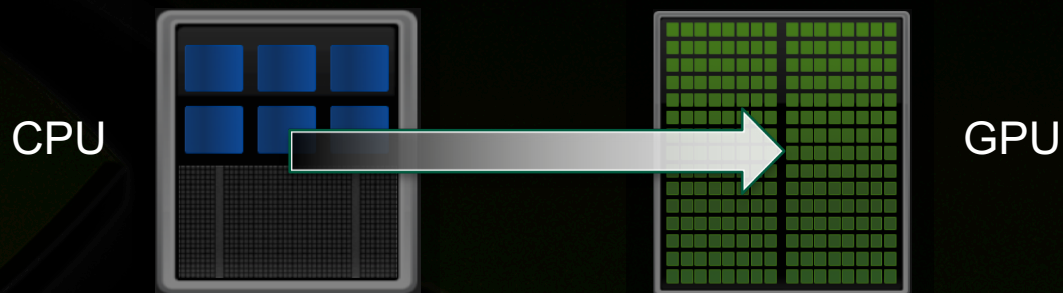


1. Setup data on CPU
2. Copy data from CPU to GPU (H2D)
3. Launch kernel on GPU
4. Copy result from GPU to CPU (D2H)
5. Repeat ...

Simple Accelerator Model



- Tasks are offloaded from CPU to GPU



1. Setup data on CPU
2. Copy data from CPU to GPU (H2D)
3. Launch kernel on GPU
4. Copy result from GPU to CPU (D2H)
5. Repeat ...

Simple Accelerator Model

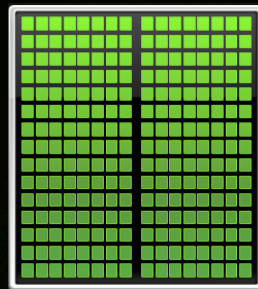


- Tasks are offloaded from CPU to GPU

CPU



GPU

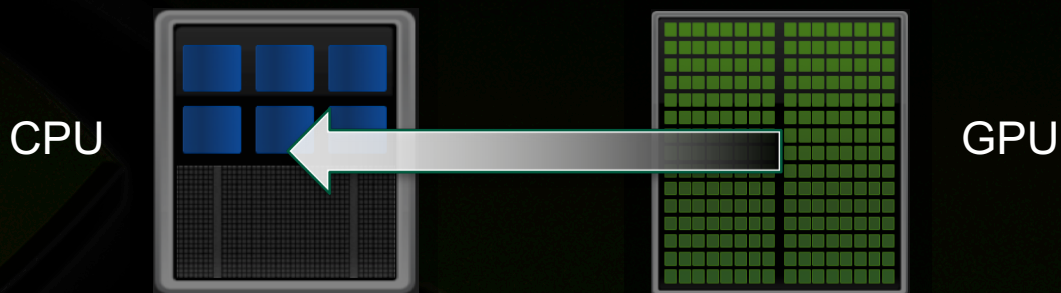


1. Setup data on CPU
2. Copy data from CPU to GPU (H2D)
3. Launch kernel on GPU
4. Copy result from GPU to CPU (D2H)
5. Repeat ...

Simple Accelerator Model



- Tasks are offloaded from CPU to GPU

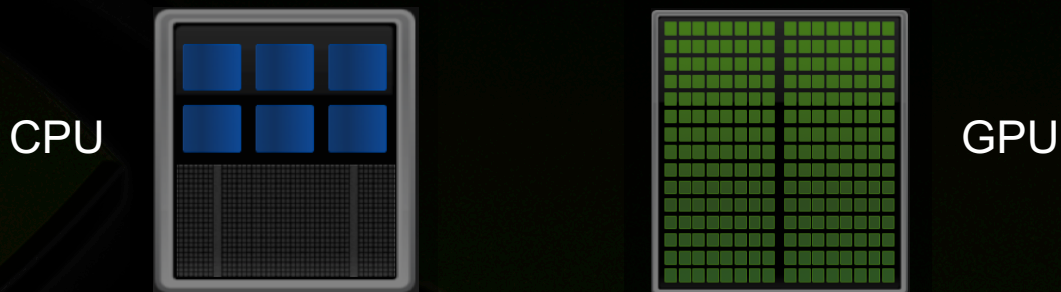


1. Setup data on CPU
2. Copy data from CPU to GPU (H2D)
3. Launch kernel on GPU
4. Copy result from GPU to CPU (D2H)
5. Repeat ...

Simple Accelerator Model

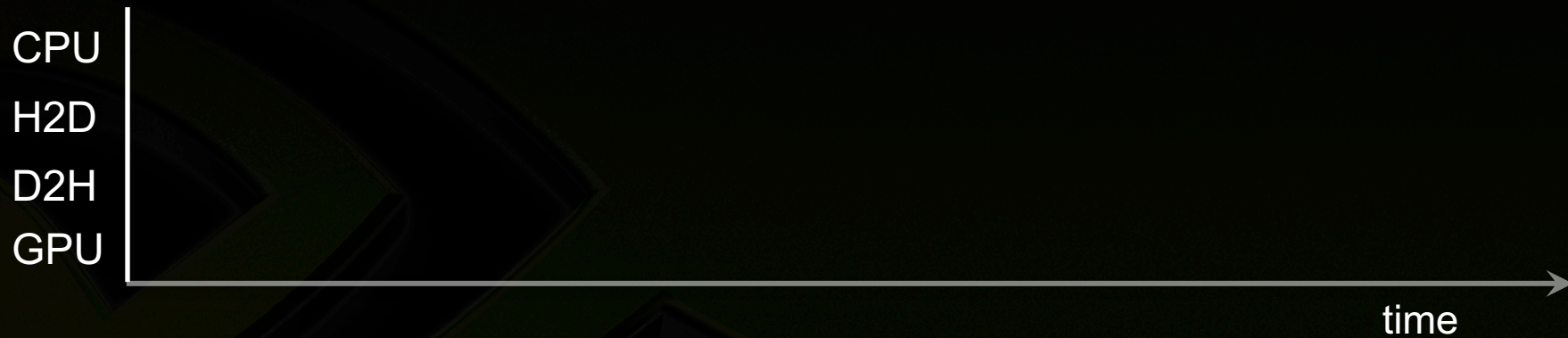


- Tasks are offloaded from CPU to GPU

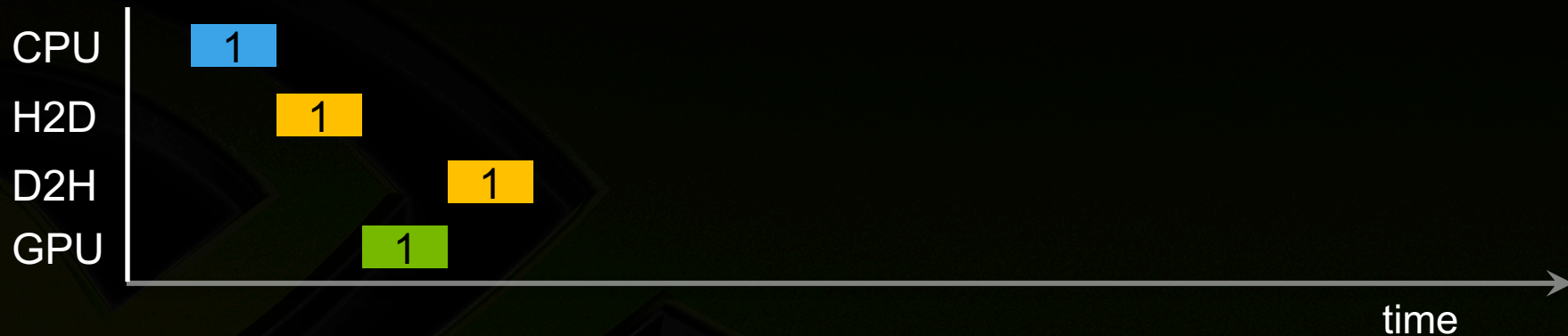


1. Setup data on CPU
2. Copy data from CPU to GPU (H2D)
3. Launch kernel on GPU
4. Copy result from GPU to CPU (D2H)
5. Repeat ...

Simple Processing Flow Timeline



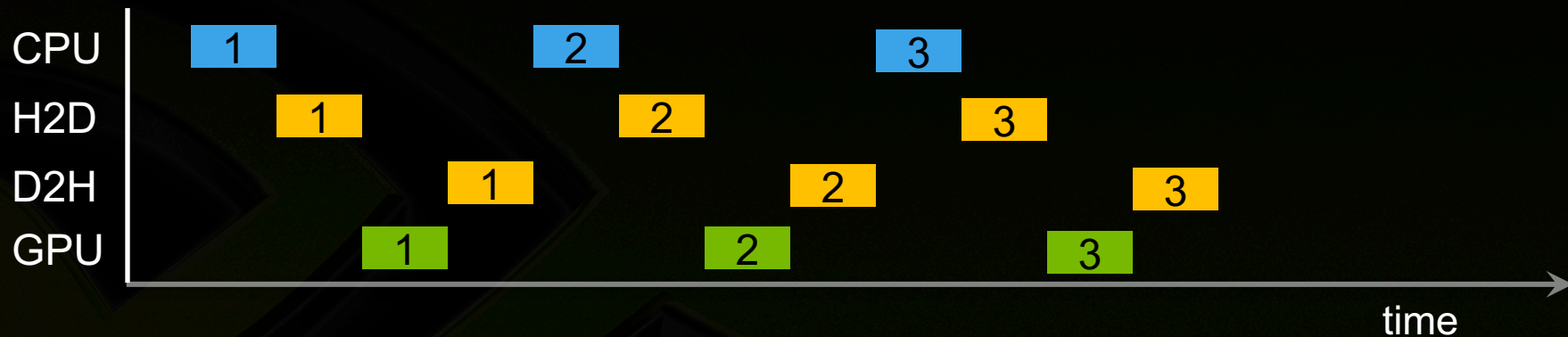
Simple Processing Flow Timeline



Simple Processing Flow Timeline



Simple Processing Flow Timeline



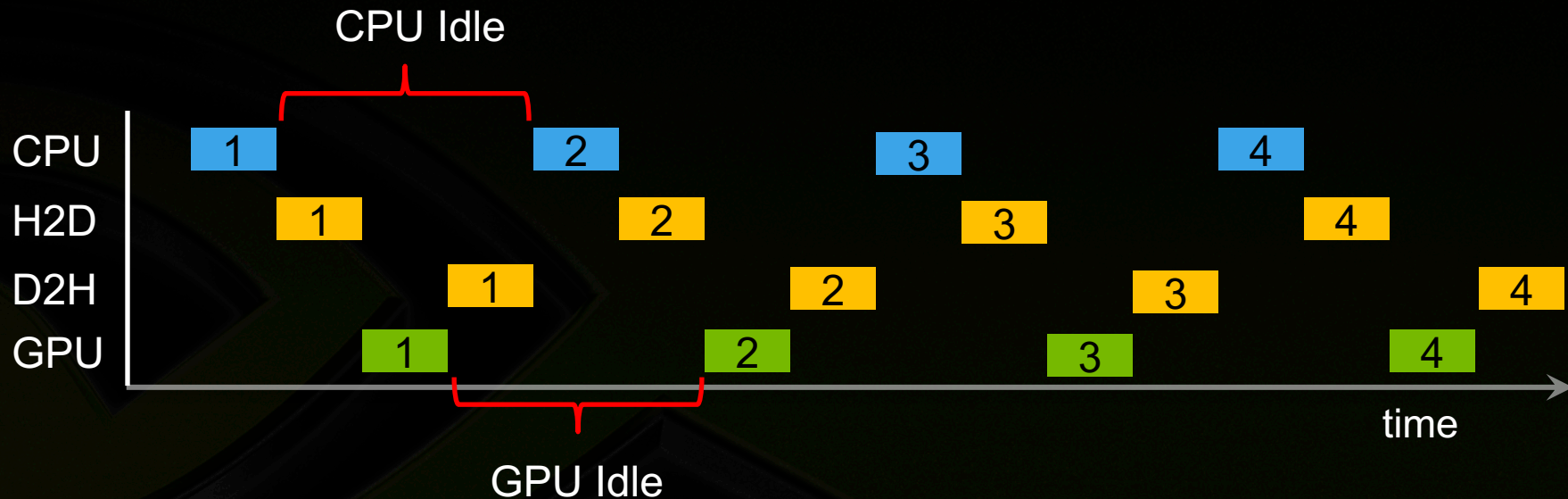
Simple Processing Flow Timeline



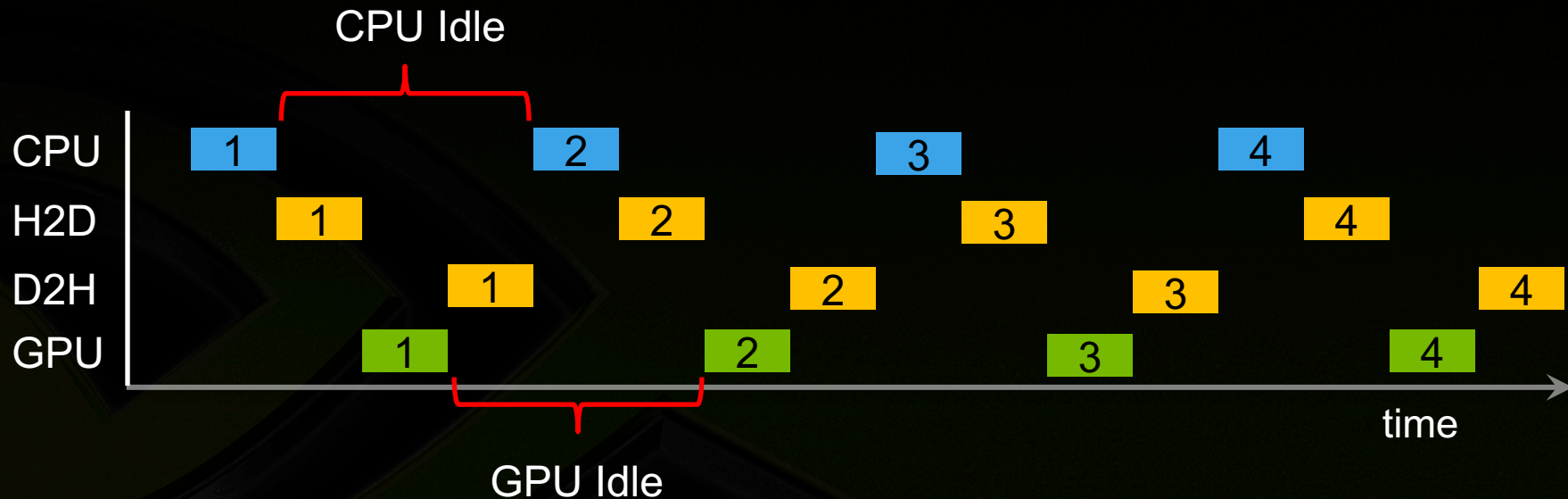
Simple Processing Flow Timeline



Simple Processing Flow Timeline



Simple Processing Flow Timeline



- Computing resources are poorly utilized
- We can use concurrency to improve utilization

Concurrent Processing Flow Timeline



- Efficient utilization of all computing resources

Concurrent Processing Flow Timeline



- Efficient utilization of all computing resources

Concurrent Processing Flow Timeline



- Efficient utilization of all computing resources

Concurrent Processing Flow Timeline



- Efficient utilization of all computing resources

Concurrent Processing Flow Timeline



- Efficient utilization of all computing resources

Concurrent Processing Flow Timeline



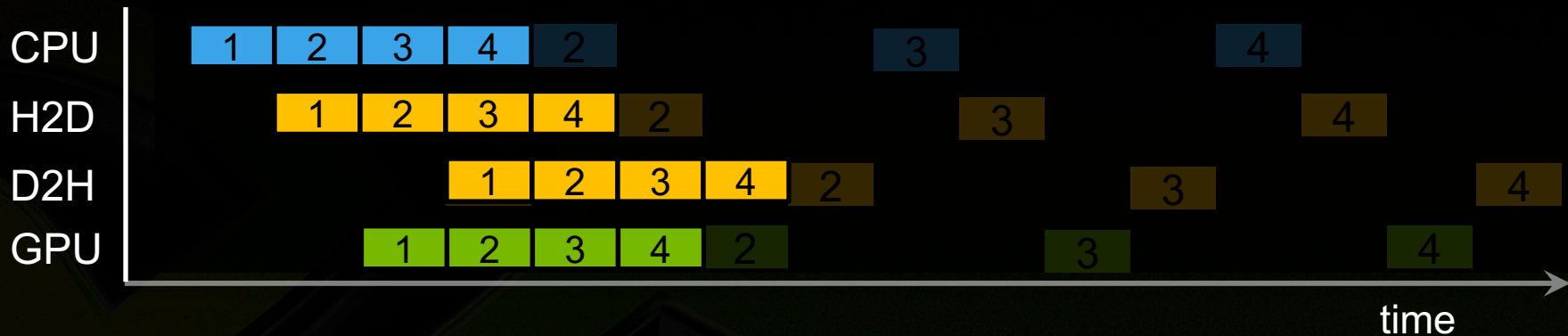
- Efficient utilization of all computing resources

Concurrent Processing Flow Timeline



- Efficient utilization of all computing resources

Concurrent Processing Flow Timeline



- Efficient utilization of all computing resources

Concurrent Processing Flow Timeline



- Efficient utilization of all computing resources

Streams and Concurrency



- Concurrency
 - The ability to perform multiple operations simultaneously
 - Compute kernels on the GPU
 - Data transfer to device (H2D)
 - Data transfer to host (D2H)
 - Operations on the CPU
 - Enables improved performance
- Streams
 - How concurrency is achieved

Enabling Concurrency with Streams

CUDA Streaming Model



- **Stream**: A sequence of operations that execute in issue-order (FIFO)
- All transfers and kernels are placed into a **stream**
 - Operations within a stream will not overlap
 - Operations in different streams may run concurrently
 - Operations from different streams may be interleaved
- Stream is the 4th launch parameter
 - `kernel <<< blocks, threads, smem , stream >>> ()`

Default Stream



- Stream used when no stream is specified
 - a.k.a. Stream '0'
 - a.k.a. 'Null Stream'
- Completely synchronous w.r.t. host and device
 - As if a `cudaDeviceSynchronize()` was inserted before and after every operation
- Exceptions – asynchronous w.r.t. host
 - Kernel launches
 - `cudaMemcpy*Async`
 - `cudaMemset*Async`
 - `cudaMemcpy` within the same device

Requirements for concurrency



- Operations must be in different, non-default, streams
- `cudaMemcpyAsync` with host from 'pinned' memory
 - page-locked memory
 - Allocated using `cudaMallocHost()` or `cudaHostAlloc()`
 - Or pin existing memory using `cudaHostRegister()`
- Sufficient resources must be available
 - `cudaMemcpyAsyncs` in different directions (#copy engines)
 - Device resources (SMEM, registers, etc.)

| | | |
|----------------------|----------------|----------------------|
| cudaMemcpyAsync(H2D) | Kernel <<< >>> | cudaMemcpyAsync(D2H) |
|----------------------|----------------|----------------------|

The diagram shows a 5-stream pipeline. Stream 0 is a yellow box labeled `cudaMemcpyAsync(H2D)`. Streams 1 through 4 are represented by dashed lines. A 4x4 kernel block, outlined by a dashed white box, is applied to streams 1 through 4. The kernel is composed of four 2x2 blocks: K1 (green), K2 (green), K3 (green), and K4 (green), each followed by a double-histogram operation (DH1, DH2, DH3, DH4) in blue boxes. The operations are arranged in a staggered fashion: K1 is on stream 1, K2 on stream 2, K3 on stream 3, and K4 on stream 4. Each K block is followed by a DH block on the same stream.

Diagram illustrating the execution of a loop on a GPU. The loop is divided into four iterations. The first iteration (HD1, K1, DH1) is completed. The second iteration (HD2, K2, DH2) is in progress. The third iteration (HD3, K3, DH3) is in progress. The fourth iteration (K4 on CPU) is in progress.

The diagram illustrates the execution of a 6x6 matrix multiplication on a 6x6 mesh of processors. The mesh is divided into two regions: a 5x5 region of processors (K1.1 to K5.5) and a 1x6 region of processors (K7 on CPU). The execution proceeds in 6 steps, each involving a row of processors (HD1 to HD6) and a column of processors (K1.1 to K6.3). The processors are color-coded: yellow for HD, green for K, and blue for DH. The diagram shows the sequence of operations from HD1 to HD6, with the final row HD6 being the only one that includes a processor from the K7 on CPU region.

Example – Tiled DGEMM



- CPU (dual 6 core SandyBridge E5-2667 @2.9 Ghz, MKL)

- 222 Gflop/s

- GPU (K20X)

- Serial: 519 Gflop/s (2.3x)

- 2-way: 663 Gflop/s (3x)

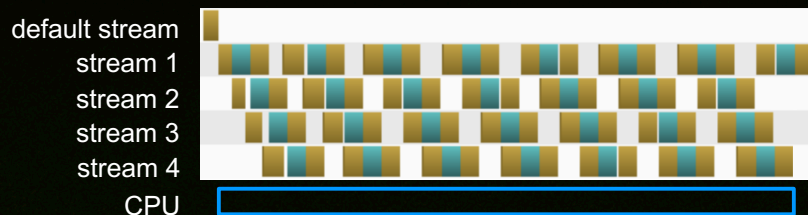
- 3-way: 990 Gflop/s (4x)

- GPU + CPU

- 4-way con.: 1180 Gflop/s (5.3x)

DGEMM: $m=n=16384$, $k=1408$

Nvidia Visual Profiler (nvvp)



- Obtain maximum performance by leveraging concurrency
- Removes impact of PCIe bandwidth
- Removes device memory size limitations

Managing Streams



- `cudaStream_t stream;`
 - Declares a stream handle
- `cudaStreamCreate(&stream);`
 - Allocates a stream
- `cudaStreamDestroy(stream);`
 - Deallocates a stream
 - Synchronizes host until work in stream has completed

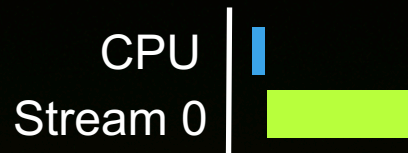
Kernel Concurrency



- Assume foo only utilizes 50% of the GPU
- Default stream

```
foo<<<blocks, threads>>>();
```

```
foo<<<blocks, threads>>>();
```



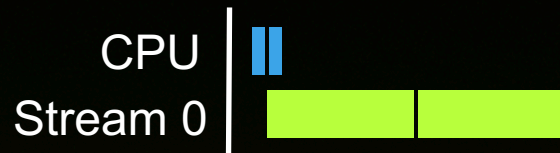
Kernel Concurrency



- Assume foo only utilizes 50% of the GPU
- Default stream

```
foo<<<blocks, threads>>>();
```

```
foo<<<blocks, threads>>>();
```

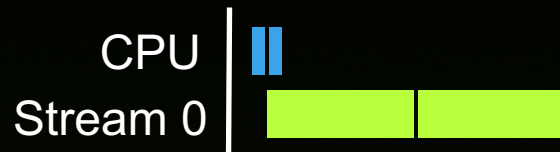


Kernel Concurrency



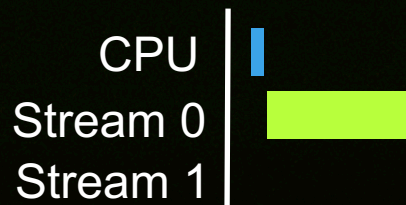
- Assume foo only utilizes 50% of the GPU
- Default stream

```
foo<<<blocks, threads>>>();  
foo<<<blocks, threads>>>();
```



- Default & user streams

```
cudaStream_t stream1;  
cudaStreamCreate(&stream1);  
foo<<<blocks, threads>>>();  
foo<<<blocks, threads, 0, stream1>>>();  
cudaStreamDestroy(stream1);
```

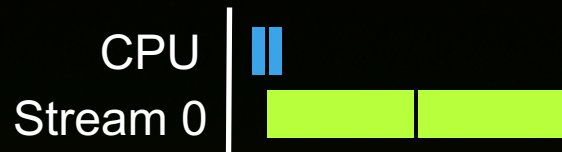


Kernel Concurrency



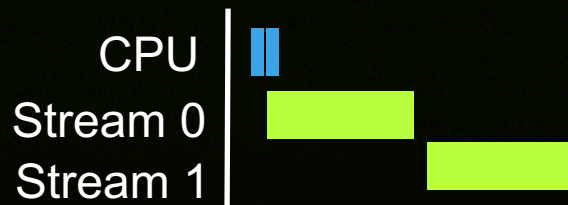
- Assume foo only utilizes 50% of the GPU
- Default stream

```
foo<<<blocks, threads>>>();  
foo<<<blocks, threads>>>();
```



- Default & user streams

```
cudaStream_t stream1;  
cudaStreamCreate(&stream1);  
foo<<<blocks, threads>>>();  
foo<<<blocks, threads, 0, stream1>>>();  
cudaStreamDestroy(stream1);
```

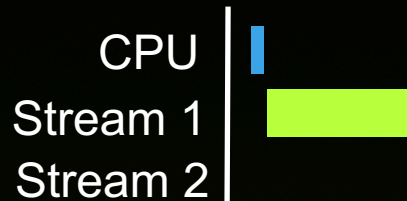


Kernel Concurrency



- User streams

```
cudaStream_t stream1, stream2;  
cudaStreamCreate(&stream1);  
cudaStreamCreate(&stream2);  
foo<<<blocks, threads, 0, stream1>>>();  
foo<<<blocks, threads, 0, stream2>>>();  
cudaStreamDestroy(stream1);  
cudaStreamDestroy(stream2);
```

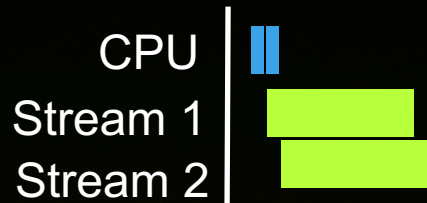


Kernel Concurrency



- User streams

```
cudaStream_t stream1, stream2;  
cudaStreamCreate(&stream1);  
cudaStreamCreate(&stream2);  
foo<<<blocks, threads, 0, stream1>>>();  
foo<<<blocks, threads, 0, stream2>>>();  
cudaStreamDestroy(stream1);  
cudaStreamDestroy(stream2);
```



Concurrent Memory Copies

- `cudaMemcpy(...)`
 - Places transfer into default stream
 - **Synchronous**: Must complete prior to returning
- `cudaMemcpyAsync(..., stream)`
 - Places transfer into stream and returns immediately
- To achieve concurrency
 - Transfers must be in a non-default stream
 - Only 1 transfer per direction at a time
 - Memory on the host must be **pinned**

Pinned Memory

- Pageable Memory (malloc, new, etc)
 - Can be paged in and out
 - Achieves a low % of peak PCIe bandwidth
- Pinned Memory
 - Cannot be paged in and out
 - Achieves a high % of peak PCIe bandwidth
- **cudaMallocHost(), cudaHostAlloc(), cudaFreeHost()**
 - Allocate/Free pinned memory on the host
 - Replaces malloc/free
- **cudaHostRegister(), cudaHostUnregister()**
 - Pins/Unpins existing memory

Paged Memory Example

```
int *h_ptr, *d_ptr;
```

```
h_ptr=malloc(bytes);
```

```
cudaMalloc(&d_ptr, bytes);
```

```
cudaMemcpyAsync(d_ptr, h_ptr, bytes, cudaMemcpyHostToDevice, stream);
```

```
free(h_ptr);
```

```
cudaFree(d_ptr);
```

(synchronous)

Pinned Memory Example

```
int *h_ptr, *d_ptr;  
  
cudaMallocHost(&h_ptr, bytes);  
cudaMalloc(&d_ptr, bytes);  
  
cudaMemcpyAsync(d_ptr, h_ptr, bytes, cudaMemcpyHostToDevice, stream);  
  
cudaFreeHost(h_ptr);  
cudaFree(d_ptr);
```

(asynchronous)

Pinned Memory Example 2

```
int *h_ptr, *d_ptr;  
  
h_ptr=malloc(bytes);  
cudaHostRegister(h_ptr,bytes,0);  
cudaMalloc(&d_ptr,bytes);  
  
cudaMemcpyAsync(d_ptr,h_ptr, bytes, cudaMemcpyHostToDevice, stream);  
  
cudaHostUnregister(h_ptr);  
free(h_ptr);  
cudaFree(d_ptr);
```

(asynchronous)

Concurrency Examples



Synchronous

```
cudaMemcpy(...);  
foo<<<...>>>();
```



Asynchronous Same Stream

```
cudaMemcpyAsync(...,stream1);  
foo<<<...,stream1>>>();
```

Asynchronous Different Streams

```
cudaMemcpyAsync(...,stream1);  
foo<<<...,stream2>>>();
```

Concurrency Examples



Synchronous

```
cudaMemcpy(...);  
foo<<<...>>>();
```



Asynchronous Same Stream

```
cudaMemcpyAsync(...,stream1);  
foo<<<...,stream1>>>();
```



Asynchronous Different Streams

```
cudaMemcpyAsync(...,stream1);  
foo<<<...,stream2>>>();
```


Concurrency Examples



Synchronous

```
cudaMemcpy(...);  
foo<<<...>>>();
```



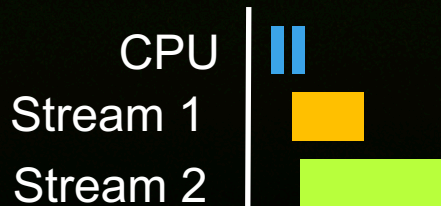
Asynchronous Same Stream

```
cudaMemcpyAsync(...,stream1);  
foo<<<...,stream1>>>();
```



Asynchronous Different Streams

```
cudaMemcpyAsync(...,stream1);  
foo<<<...,stream2>>>();
```



Implicit Synchronization



- These functions *implicitly* synchronize the device
 - `cudaMalloc`, `cudaFree`
 - `cudaEventCreate`, `cudaEventDestroy`,
 - `cudaStreamCreate`, `cudaStreamDestroy`
 - `cudaHostRegister`, `cudaHostUnregister`
 - `cudaFuncSetCacheConfig`
- Avoid by reusing memory and data structures as much as possible

Explicit Synchronization



- CUDA provides mechanisms for expressing synchronization between the host, device, and streams.
- Synchronize everything
 - `cudaDeviceSynchronize()`
 - Blocks host until all issued CUDA calls are complete
- Synchronize host w.r.t. a specific stream
 - `cudaStreamSynchronize (streamid)`
 - Blocks host until all issued CUDA calls in streamid are complete

Explicit Synchronization using Events



- Mechanism for arbitrary synchronization
 - Create 'events' at specific points within streams
 - boolean state: occurred / not occurred, default = occurred
- Synchronize using events
 - `cudaEventRecord` (event, streamid)
 - `cudaEventQuery` (event)
 - `cudaEventSynchronize` (event)
 - `cudaStreamWaitEvent` (stream, event)

Explicit Synchronization Example



```
cudaEvent_t event;  
cudaEventCreate (&event); // create event  
  
cudaMemcpyAsync (d_in, in, size, H2D, stream1); // 1) H2D copy of new input  
cudaEventRecord (event, stream1); // record event  
  
cudaMemcpyAsync (out, d_out, size, D2H, stream2); // 2) D2H copy of previous  
// result  
  
cudaStreamWaitEvent (stream2, event); // wait for event in stream1  
kernel <<< , , , stream2 >>> (d_in, d_out); // 3) must wait for 1 and 2  
  
asynchronousCPUMethod ( ... ) // Async CPU method
```


Explicit Synchronization Example

```

cudaEvent_t event;
cudaEventCreate (&event);                                // create event

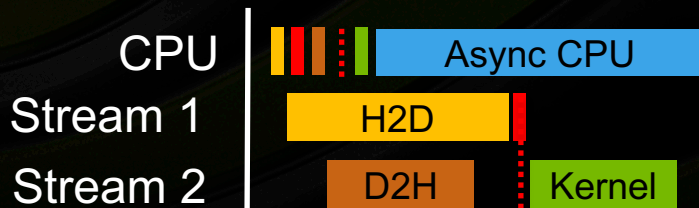
cudaMemcpyAsync (d_in, in, size, H2D, stream1);           // 1) H2D copy of new input
cudaEventRecord (event, stream1);                         // record event

cudaMemcpyAsync (out, d_out, size, D2H, stream2);          // 2) D2H copy of previous
                                                         // result

cudaStreamWaitEvent (stream2, event);                     // wait for event in stream1
kernel <<< , , , stream2 >>> (d_in, d_out);               // 3) must wait for 1 and 2

asynchronousCPUMethod ( ... )                             // Async CPU method

```



CUDA_LAUNCH_BLOCKING



- Environment variable which forces synchronization
 - `export CUDA_LAUNCH_BLOCKING=1`
 - All CUDA operations are synchronous w.r.t the host
- Useful for debugging race conditions
 - If it runs successfully with `CUDA_LAUNCH_BLOCKING` set but doesn't without you have a race condition.

Advanced: Relaxed Default Streams

- Libraries sometimes use the default stream internally
- Explicitly disable blocking from default stream (CUDA 5.0+)
 - `cudaStreamCreateWithFlags(&stream, cudaStreamNonBlocking);`
- CUDA 7.0 can create separate default streams for each host thread
 - NVCC option `--default-stream per-thread`
- Explicit handles for default streams
 - `cudaStreamLegacy`, `cudaStreamPerThread`
 - Use case:
 - `cudaDeviceSynchronize()` blocks all streams on the device!
 - Use `cudaStreamSynchronize(cudaStreamPerThread)` instead

Advanced: Priority Streams (K20+)



- You can give streams priority
 - CUDA 5.5+
 - High priority streams will preempt lower priority streams.
 - Currently executing blocks will complete but new blocks will only be scheduled after higher priority work has been scheduled.
- Query available priorities:
 - `cudaDeviceGetStreamPriorityRange(&low, &high)`
 - Kepler: low = 0, high = -1
 - Lower number is higher priority
- Create using special API:
 - `cudaStreamCreateWithPriority(&stream, flags, priority)`

Advanced: Stream Callbacks



- Cuda 5.0+ allows you to add stream callbacks (K20 or newer)
 - Useful for launching work on the host when something has completed

```
void CUDART_CB MyCallback(void *data){  
    ...  
}  
...  
MyKernel<<<100, 512, 0, stream>>>();  
    cudaStreamAddCallback(stream, MyCallback, (void*)i, 0);
```

- Callbacks are processed by a driver thread
 - The same thread processes all callbacks
 - You can use this thread to signal other threads
 - Cannot make any CUDA calls from callback

Advanced: Multiple GPUs



- `cudaDeviceSynchronize` syncs with current device only
- Streams are associated with a particular device
 - Current device when stream was created
 - Error if stream is referenced when its device is not current
 - Each device has its own default stream
- Events are associated with a particular device
 - `cudaEventRecord` will fail if event and stream associate with different devices
 - `cudaEventElapsedTime` must take two events associated with the same device
- **Synchronization works between devices with any event**
 - `cudaEventQuery`, `cudaEventSynchronize`, `cudaStreamWaitEvent`

Advanced: Multi-Process Service (MPS)



- Background:
 - Each process has a unique context.
 - Only a single context can be active on a device at a time.
 - Multiple processes (e.g. MPI) on a single GPU could not operate concurrently
- MPS: Software layer that sits between the driver and your application.
 - Routes all CUDA calls through a single context
 - Multiple processes can execute concurrently

Advanced: Multi-Process Service (cont)



- Advantages:
 - Oversubscribe MPI processes and concurrency occurs automatically
 - E.g. 1 MPI process per core sharing a single GPU
 - Simple and natural path to acceleration (especially if your application is MPI ready)
- Disadvantage:
 - MPS adds extra launch latency
 - Not supported on older hardware (Kepler and newer)
 - Linux Only

Common Streaming Issues

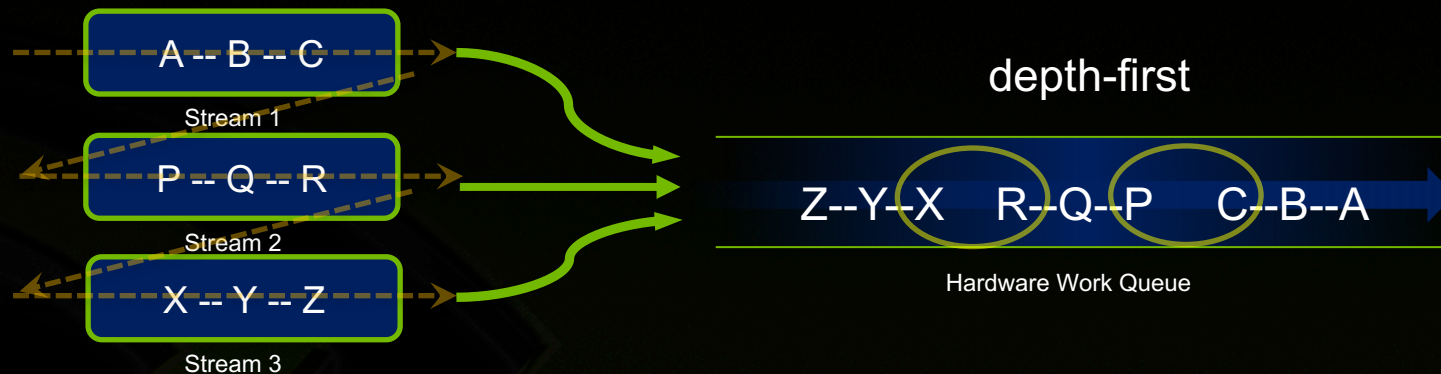
- Using the default stream
- Not using asynchronous version of memcpy
- Not using pinned host memory for memcpy
- Implicit synchronization
- Concurrency can be disabled for debugging
 - `CUDA_LAUNCH_BLOCKING=1`

Potential Hazards



- Concurrency is broken if there are more than 62 outstanding operations
 - Kernels, memory copies, recorded events
 - That is, in 'issue order' concurrent operations must not be separated by more than 62 other issues
 - Further operations are serialized
 - Avoid by changing issue order
- Kernels using more than 8 textures cannot run concurrently
- Switching L1/Shared configuration may break concurrency
- Hardware older than SM3.5 suffers from false serialization

Fermi – Concurrent Kernels

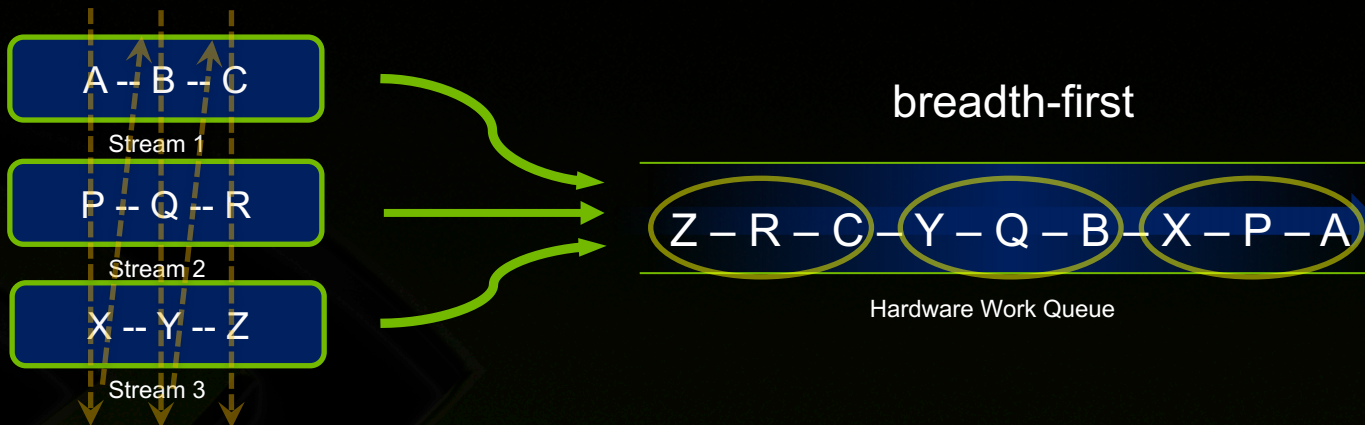


Fermi allows 16-way concurrency

- But CUDA kernels multiplex into a single queue
- Issue order matters for concurrency see “CUDA Concurrency & Streams”

<https://developer.nvidia.com/gpu-computing-webinars>

Fermi – Concurrent Kernels

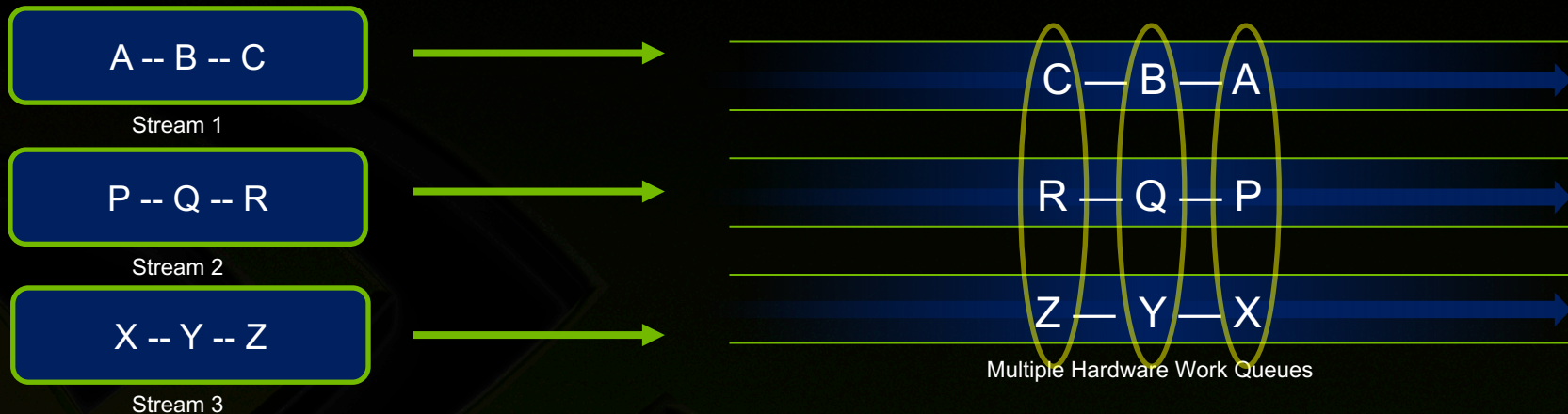


Fermi allows 16-way concurrency

- But CUDA kernels multiplex into a single queue
- Issue order matters for concurrency see “CUDA Concurrency & Streams”

<https://developer.nvidia.com/gpu-computing-webinars>

K20 Improved Concurrency

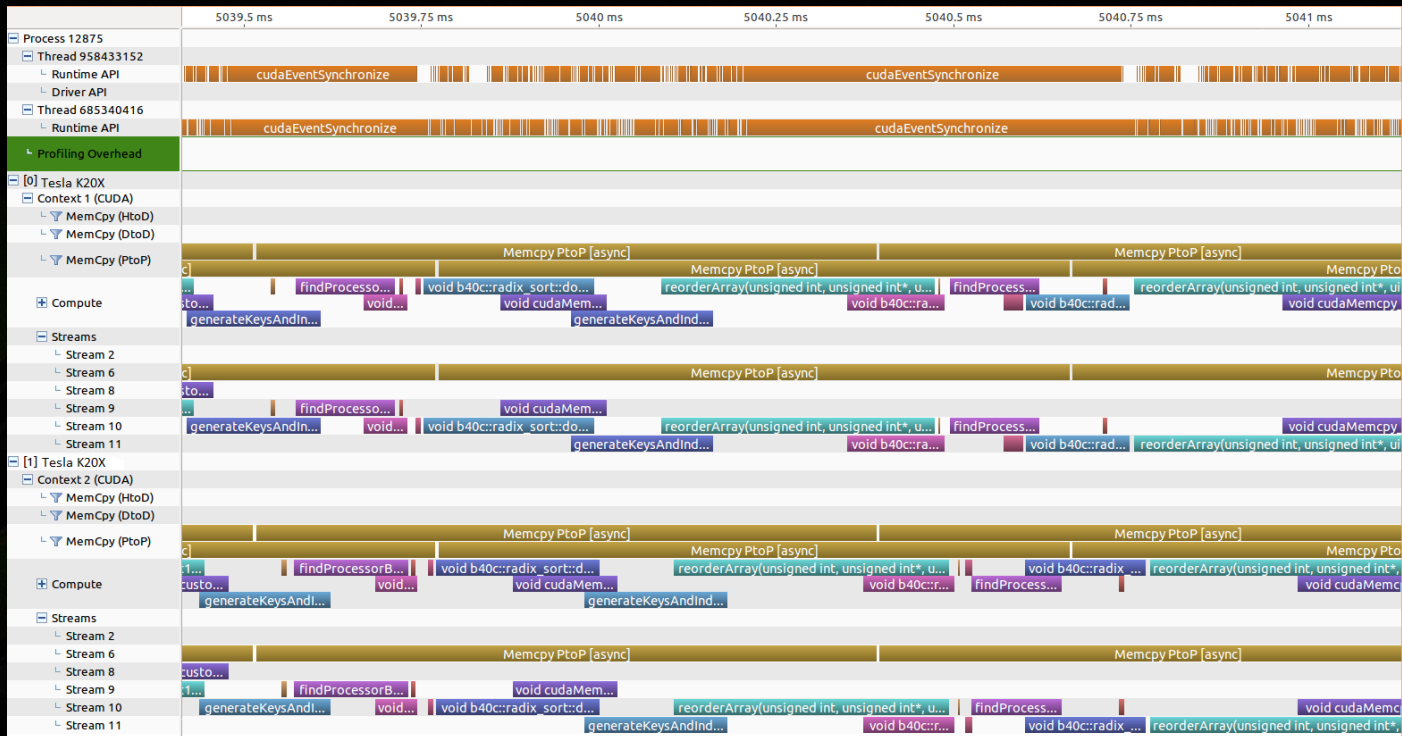


Kepler allows 32-way concurrency

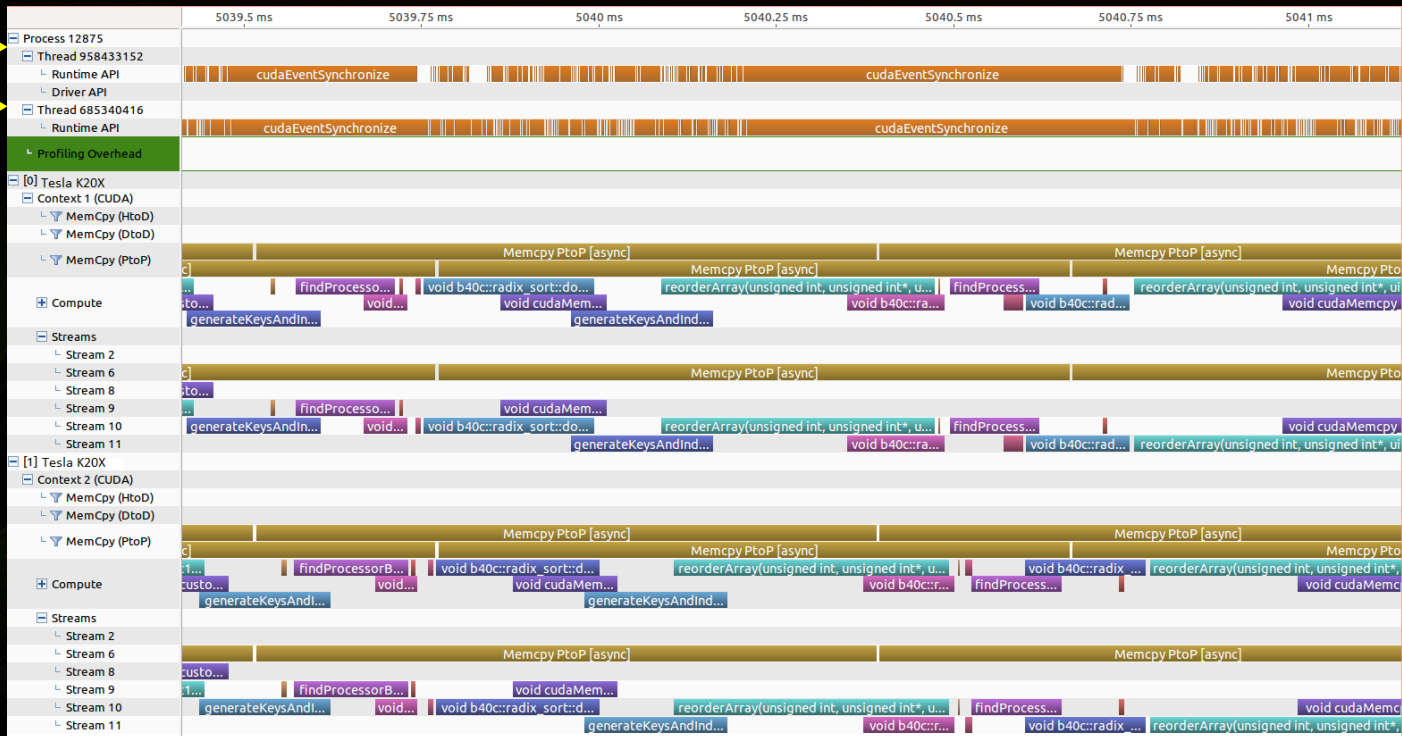
- One kernel queue per stream
- No inter-stream dependencies

depth-first or
breadth-first

Nvidia Visual Profiler (nvvp)



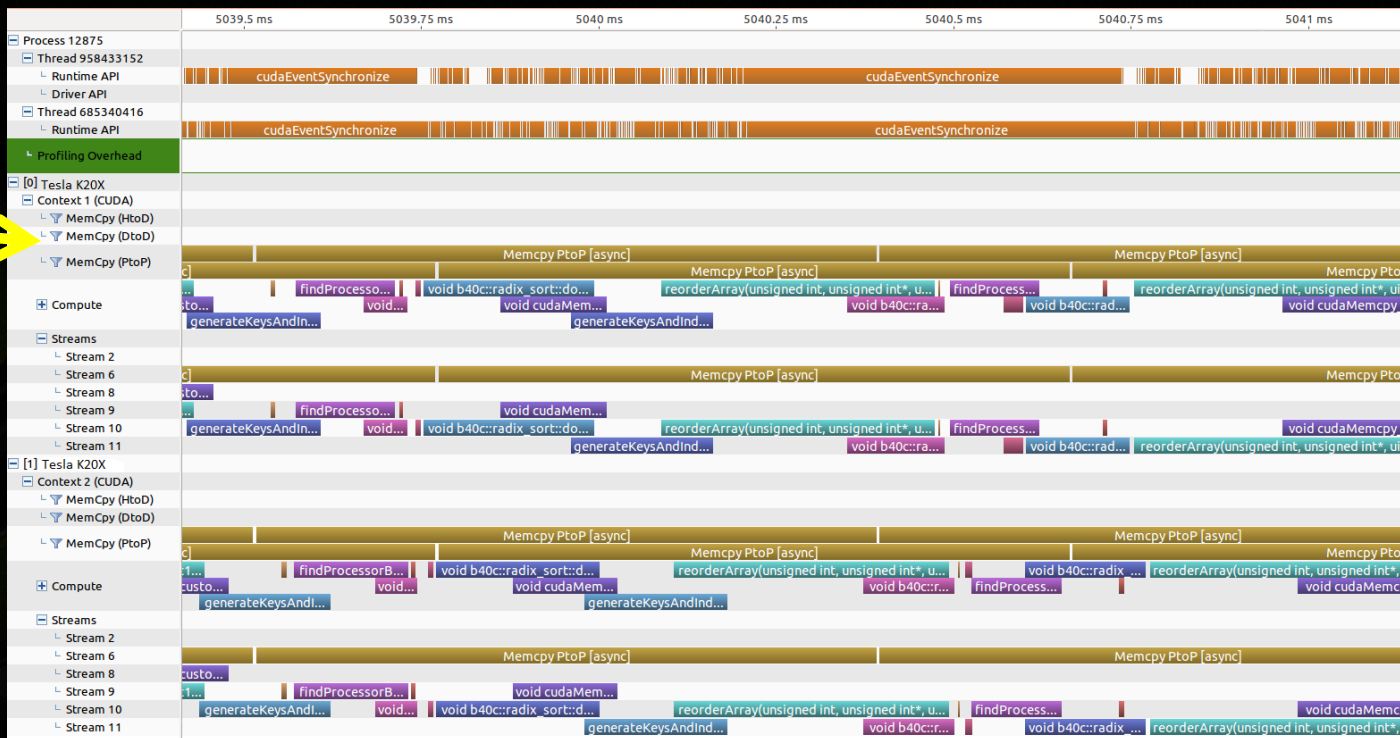
What is the host doing?



Nvidia Visual Profiler (nvvp)



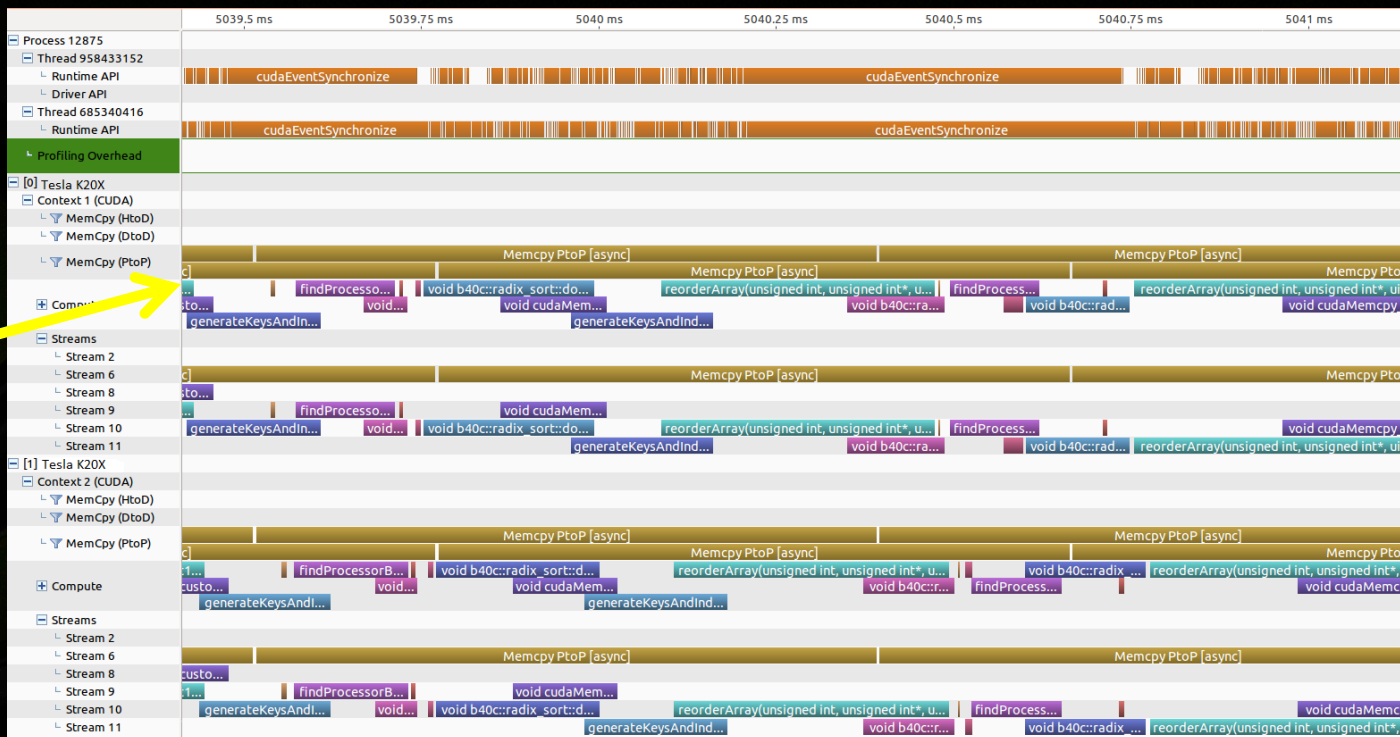
What copies are happening?



Nvidia Visual Profiler (nvvp)



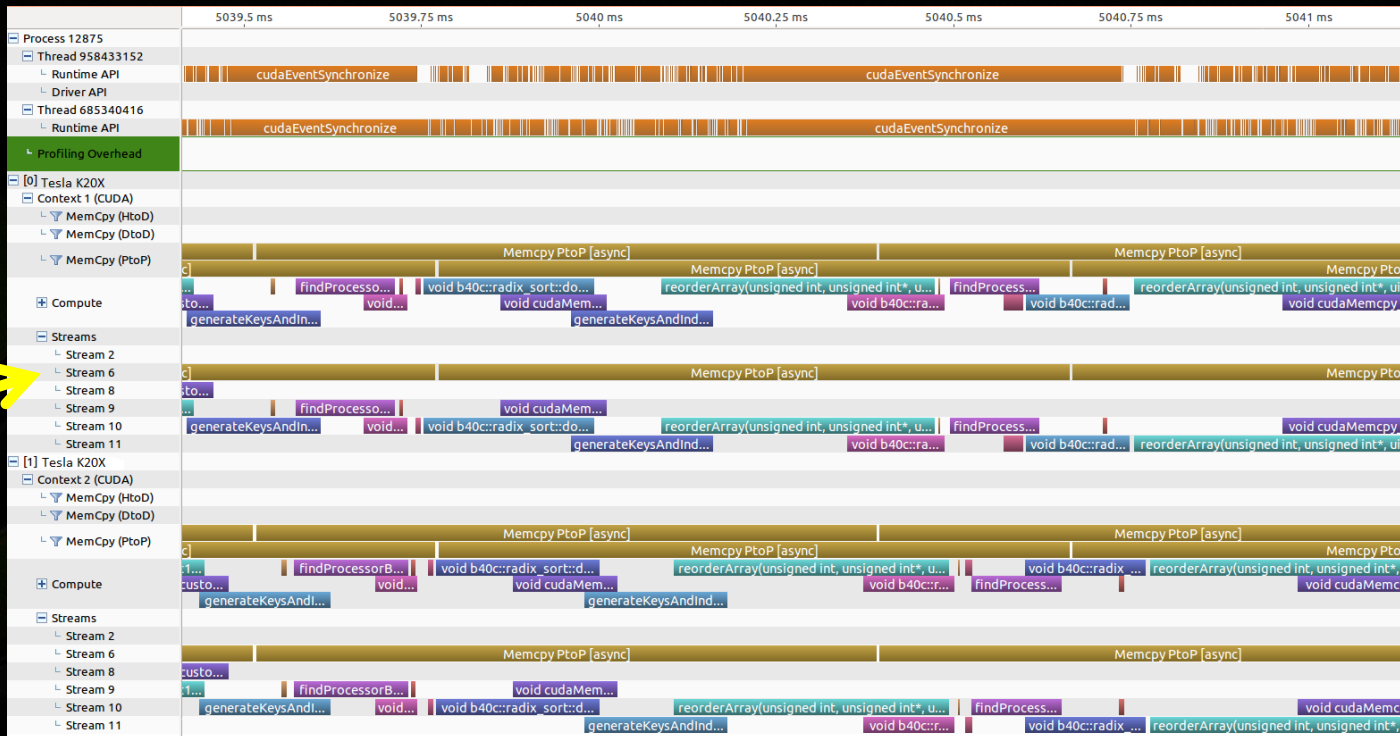
How much overlap?



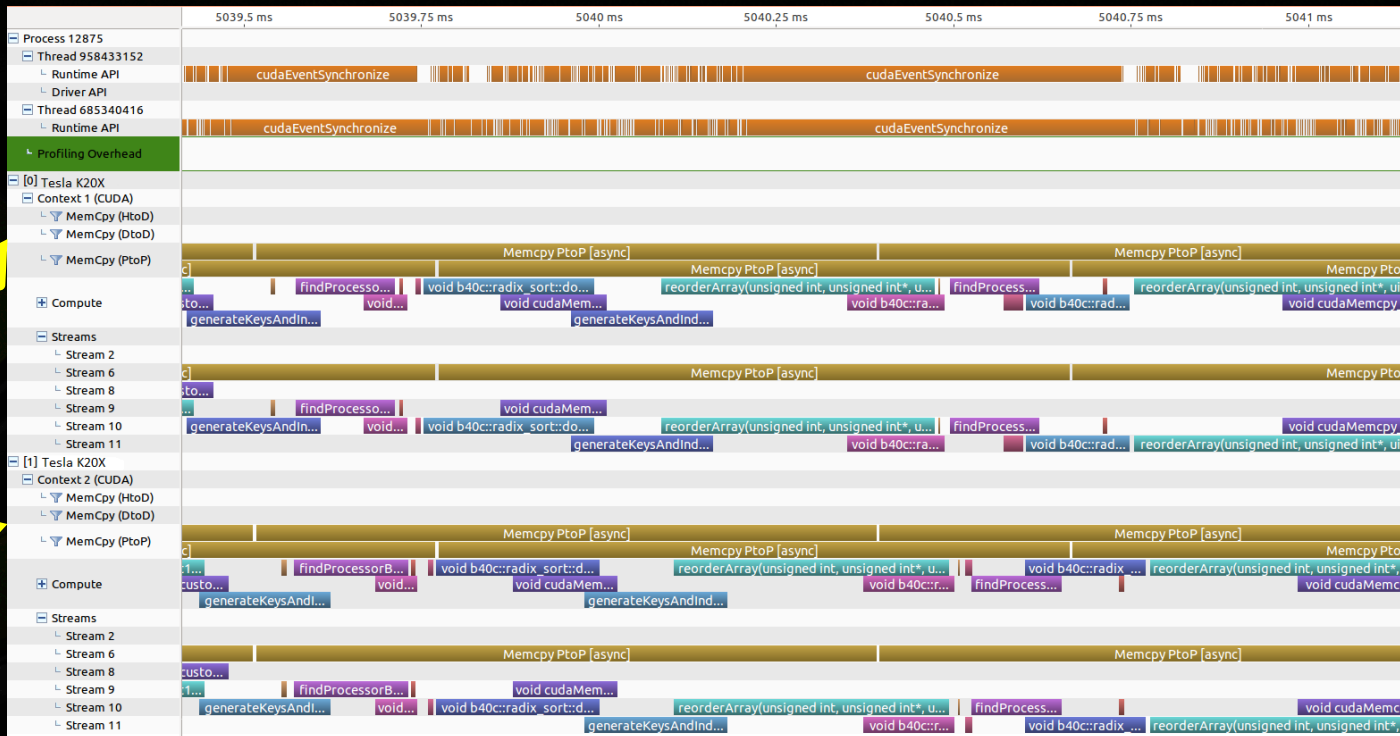
Nvidia Visual Profiler (nvvp)



Streams



Nvidia Visual Profiler (nvvp)



Multi GPU

Concurrency Guidelines



- Code to the programming model – streams
 - Future devices will continually improve HW rep. of programming model
 - Often possible to hide host-device communication overhead
- Pay attention to operations which can break concurrency
 - Use of default stream
 - Implicit synchronization
- Synchronize only when required
 - Excessive synchronization imparts overhead, limits scheduler
- **Use profilers! (nvvp, Nsight, ...)**



Questions?

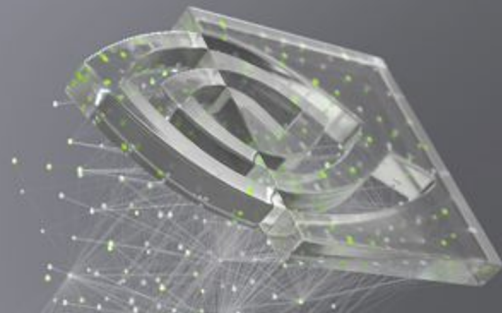
docs.nvidia.com
google[GTC on demand]



An MPS use case for deep learning



- <http://on-demand.gputechconf.com/gtc/2017/presentation/S7320-tim-kaldewey-optimizing-efficiency-of-deep-learning-workloads-through-gpu-virtualization.pdf>



YOUR FUTURE STARTS HERE.

NVIDIA is **hiring interns** and **new college grads**. Come join the industry leader in virtual reality, artificial intelligence, self-driving cars, and gaming.

Learn more at www.nvidia.com/university

General hiring areas bit.ly/nvidiaur-jd

Email your resume to Chau Luu cluu@nvidia.com



NVIDIA.