

Applied Parallel Programming

Lecture 21: Data Transfer and CUDA Streams (Task Parallelism)

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

1

Objective

To learn more advanced features of the CUDA APIs
for data transfer and kernel launch

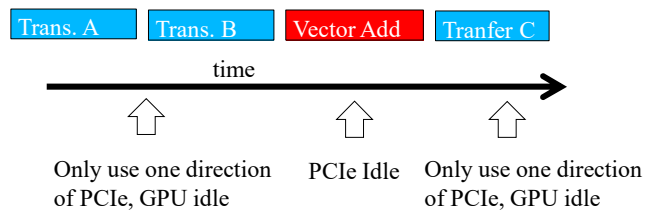
- Task parallelism for overlapping data transfer with kernel computation
- CUDA streams

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

2

Serialized Data Transfer and GPU computation

- So far, the way we use cudaMemcpy serializes data transfer and GPU computation



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

3

Device Overlap

- Most CUDA devices support *device overlap*
 - *Simultaneously execute a kernel while performing a copy between device and host memory*

```
int dev_count;
cudaDeviceProp prop;

cudaGetDeviceCount(&dev_count);
for (int i = 0; i < dev_count; i++) {
    cudaGetDeviceProperties(&prop, i);

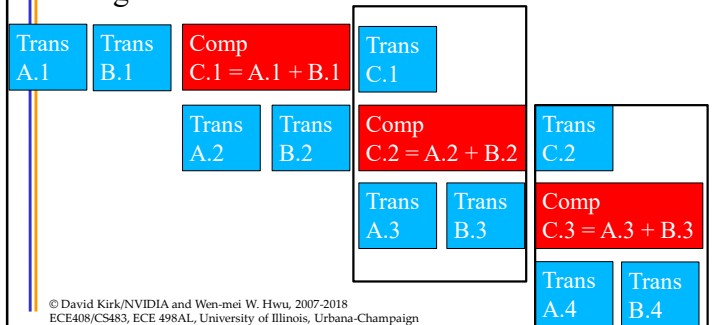
    if (prop.deviceOverlap) ...
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

4

Overlapped (Pipelined) Timing

- Divide large vectors into segments
- Overlap transfer and compute of adjacent segments



5

Using CUDA Streams and Asynchronous Memcpy

- CUDA supports parallel execution of kernels and cudaMemcpy with **streams**
- Each stream **is a queue of operations** (kernel launches and cudaMemcpy's)
- Operations (tasks) in different streams
 - can execute in parallel
 - a version of **task parallelism**

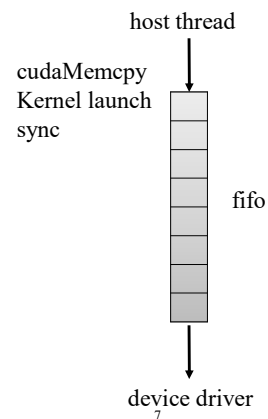
© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

6

Streams

Queue operation:

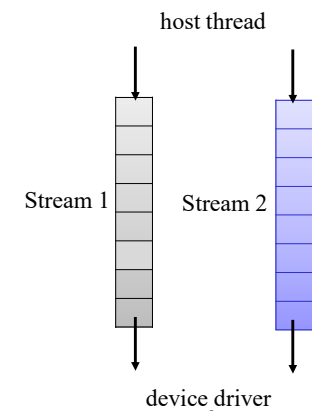
1. device requests (host code) placed into queue
2. queue processed asynchronously by driver and device
3. each queue processed in order (no overlap), so memory copies end before kernel launch, and so forth



7

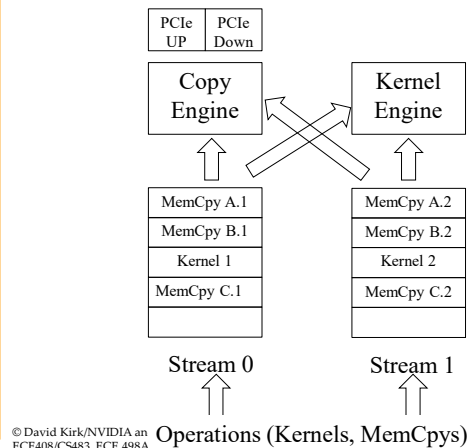
Multiple Streams Enable Parallelism

- To allow concurrent copying and kernel execution, you need to use multiple streams.



8

Conceptual View of Streams



9

A Simple Multi-Stream Host Code

```
cudaStream_t stream0, stream1;
cudaStreamCreate( &stream0);
cudaStreamCreate( &stream1);
float *d_A0, *d_B0, *d_C0; // device memory for stream 0
float *d_A1, *d_B1, *d_C1; // device memory for stream 1

// cudaMalloc for d_A0, d_B0, d_C0, d_A1, d_B1, d_C1 go here

for (int i=0; i<n; i+=SegSize*2) {
    cudaMemcpyAsync(d_A0, h_A+i, SegSize*sizeof(float),..., stream0);
    cudaMemcpyAsync(d_B0, h_B+i, SegSize*sizeof(float),..., stream0);
    vecAdd<<<SegSize/256, 256, 0, stream0>>>(d_A0, d_B0, ...);
    cudaMemcpyAsync(h_C+i, d_C0, SegSize*sizeof(float),..., stream0);
    cudaMemcpyAsync(d_A1, h_A+i+SegSize, SegSize*sizeof(float),..., stream1);
    cudaMemcpyAsync(d_B1, h_B+i+SegSize, SegSize*sizeof(float),..., stream1);
    vecAdd<<<SegSize/256, 256, 0, stream1>>>(d_A1, d_B1, ...);
    cudaMemcpyAsync(h_C+i+SegSize, d_C1, SegSize*sizeof(float),..., stream1);
}
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

10

A Simple Multi-Stream Host Code (Cont.)

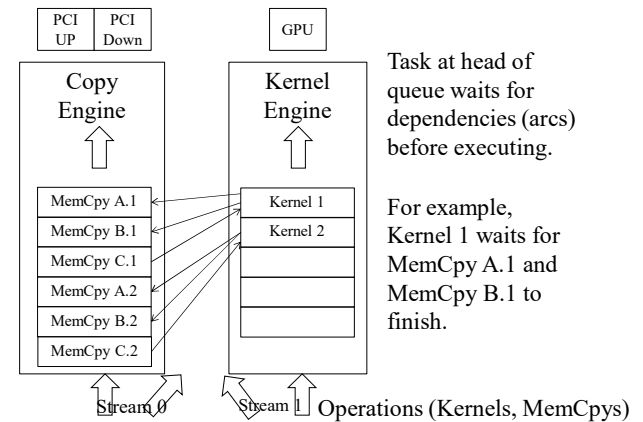
```
for (int i=0; i<n; i+=SegSize*2) {
    cudaMemcpyAsync(d_A0, h_A+i, SegSize*sizeof(float),..., stream0);
    cudaMemcpyAsync(d_B0, h_B+i, SegSize*sizeof(float),..., stream0);
    vecAdd<<<SegSize/256, 256, 0, stream0>>>(d_A0, d_B0, ...);
    cudaMemcpyAsync(h_C+i, d_C0, SegSize*sizeof(float),..., stream0);

    cudaMemcpyAsync(d_A1, h_A+i+SegSize, SegSize*sizeof(float),..., stream1);
    cudaMemcpyAsync(d_B1, h_B+i+SegSize, SegSize*sizeof(float),..., stream1);
    vecAdd<<<SegSize/256, 256, 0, stream1>>>(d_A1, d_B1, ...);
    cudaMemcpyAsync(h_C+i+SegSize, d_C1, SegSize*sizeof(float),..., stream1);
}
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

11

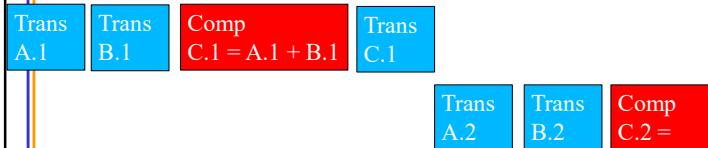
Older GPUs Support Streams in Software



12

Not quite the overlap we want

- C.1 blocks A.2 and B.2 in the copy engine queue (head-of-line blocking).



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

13

A Better Multi-Stream Host Code

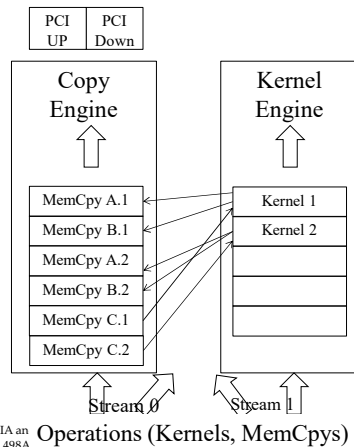
```
for (int i=0; i<n; i+=SegSize*2) {
    cudaMemcpyAsync(d_A0, h_A+i, SegSize*sizeof(float),..., stream0);
    cudaMemcpyAsync(d_B0, h_B+i, SegSize*sizeof(float),..., stream0);
    cudaMemcpyAsync(d_A1, h_A+i+SegSize,
                    SegSize*sizeof(float),..., stream1);
    cudaMemcpyAsync(d_B1, h_B+i+SegSize,
                    SegSize*sizeof(float),..., stream1);

    vecAdd<<<SegSize/256, 256, 0, stream0>>>(d_A0, d_B0, ...);
    vecAdd<<<SegSize/256, 256, 0, stream1>>>(d_A1, d_B1, ...);
    cudaMemcpyAsync(d_C0, h_C+i, SegSize*sizeof(float),..., stream0);
    cudaMemcpyAsync(d_C1, h_C+i+SegSize,
                    SegSize*sizeof(float),..., stream1);
}
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

14

A View Closer to Reality

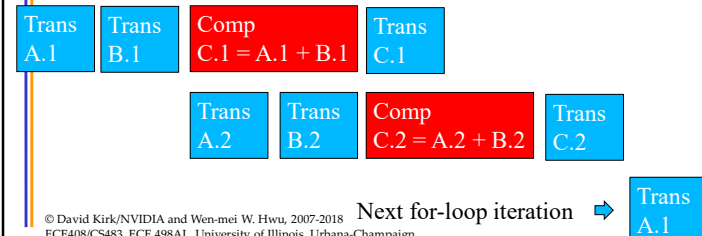


© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

15

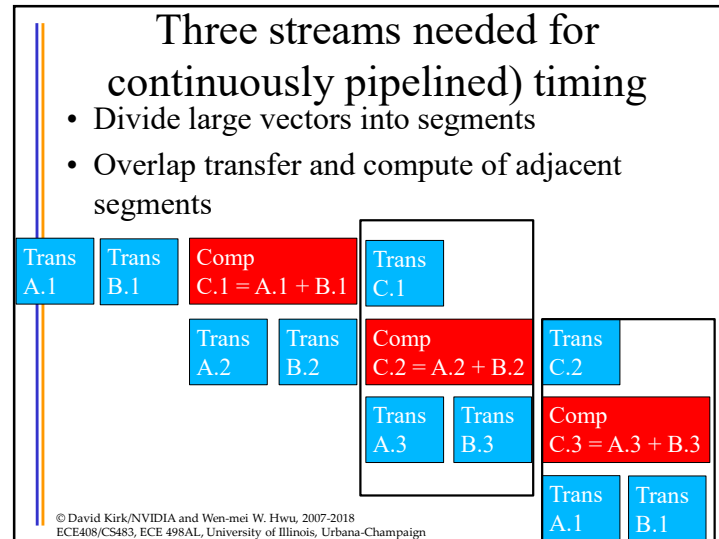
Better Overlap with Two Streams

- C.1 no longer blocks A.2 and B.2 in the copy engine queue
- However, C.2 still blocks A.1 and A.2 from the next iteration – PCIe used for only one direction

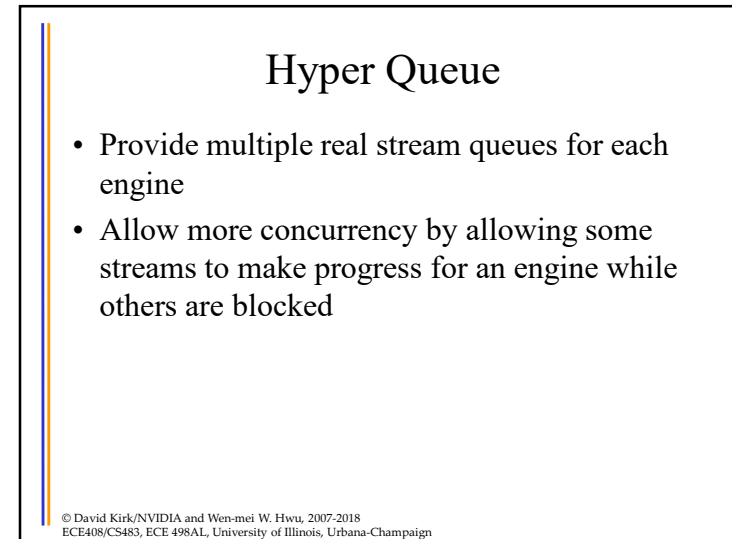


© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

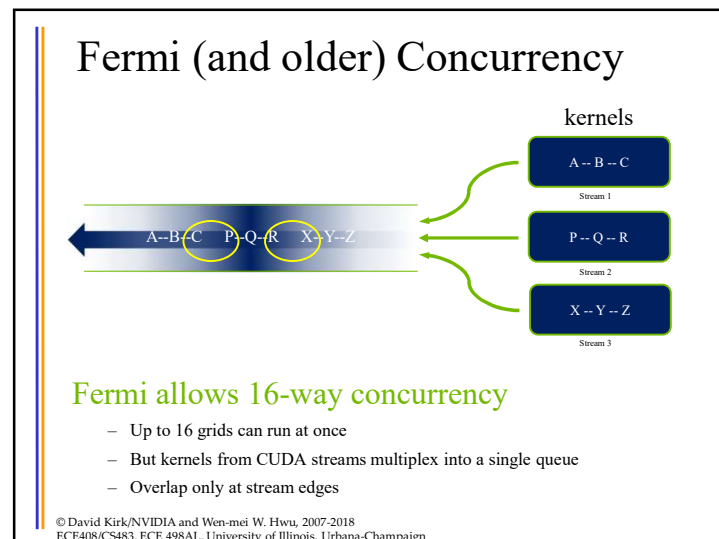
16



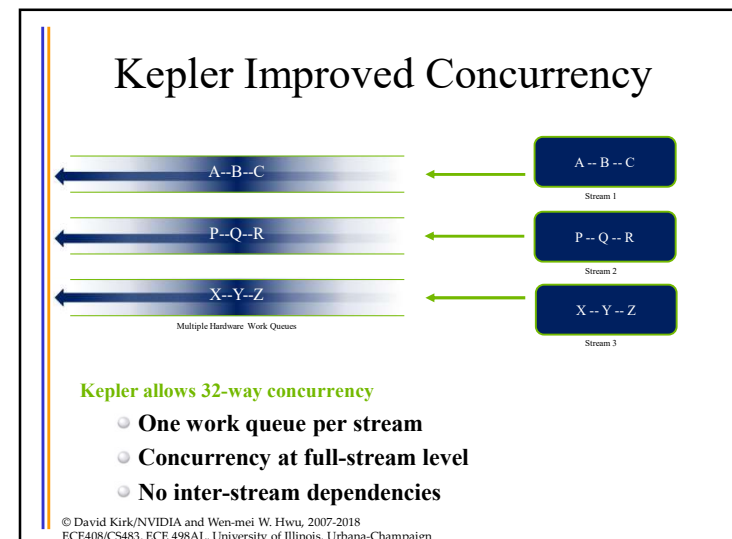
17



18



19



20

Smaller Segments Reduce Boundary Effects

How small should segments be?

- If we **overlap**
 - **transfer of** segment N's **inputs**,
 - **computation** of segment N – 1, **and**
 - **transfer** of segment N – 2's **results**,
- we **still have non-overlapping work** at the beginning and the end.

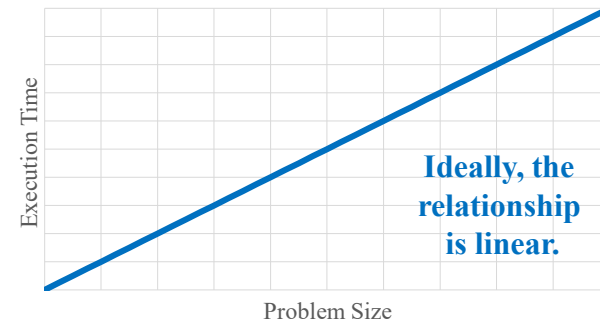
So segments should be really small?

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

21

Execution Time is Ideally Linear in Size

Think about execution time as a function of segment size.

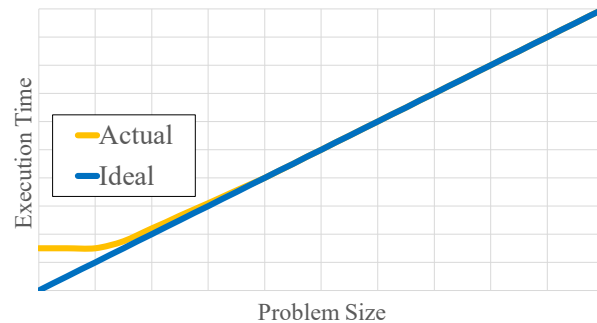


© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

22

Execution Time Never Reaches Zero

But real execution time has a minimum.

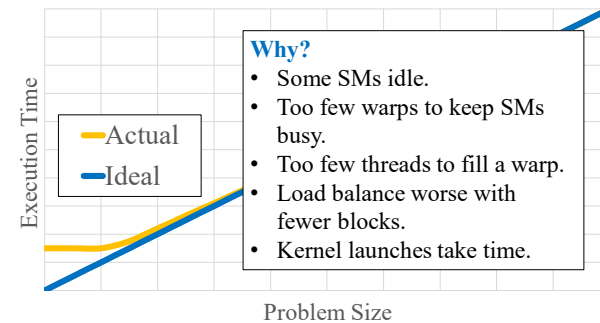


© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

23

Execution Time Never Reaches Zero

But real execution time has a minimum.



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

24

Use Moderate Segment Size and Device Query

Data transfers

- **have similar non-linearities** for small sizes
- due to startup costs on host and DMA.

So how small should segments be?

Moderately sized.

Best size likely to **depend on GPU**.

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

25

ANY QUESTIONS?

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018
ECE408/CS483, ECE 498AL, University of Illinois, Urbana-Champaign

26