



SIGGRAPH ASIA 2013

Shaping the Future of
Visual Computing

Multi-GPU Programming for Visual Computing

Wil Braithwaite - NVIDIA Applied Engineering

Talk Outline

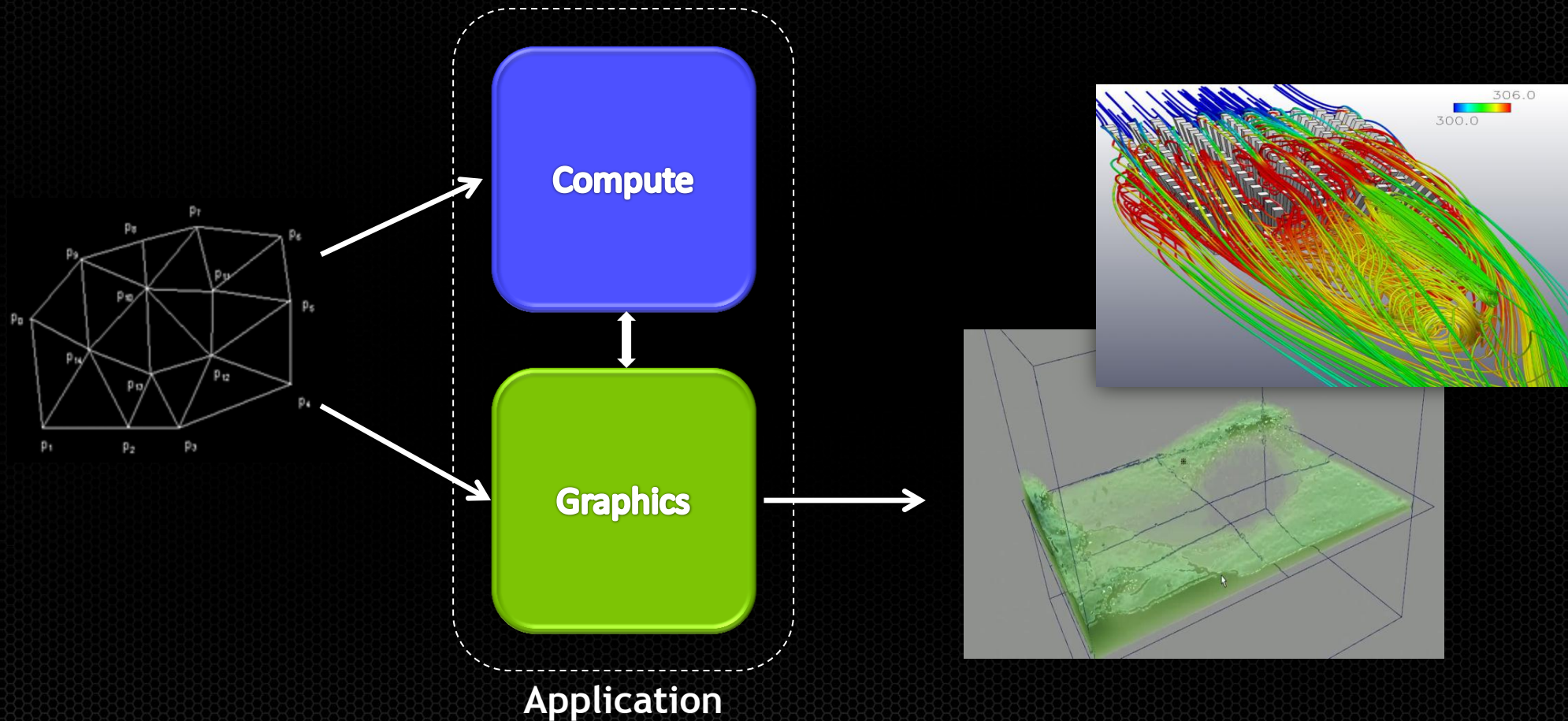
■ Part 1 (Wednesday 3pm)

- Compute and Graphics API Interoperability.
 - The basics.
 - Interoperability Methodologies & driver optimizations.
- Interoperability at a system level.
 - fine-grained control for managing scaling
- Demo.

■ Part 2 (Thursday 12pm)

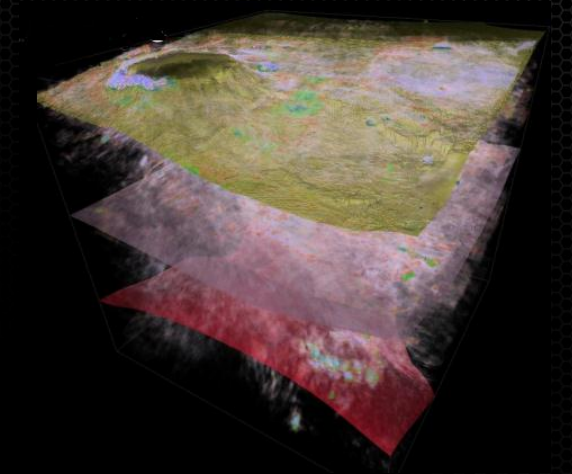
- Peer-to-peer & UVA, NUMA considerations
- Scaling beyond one-to-one, single-compute : single-render
 - Enumerating Graphics & Compute Resources
 - many-to-many, multiple-compute : multiple-render

Compute & Visualize the same data

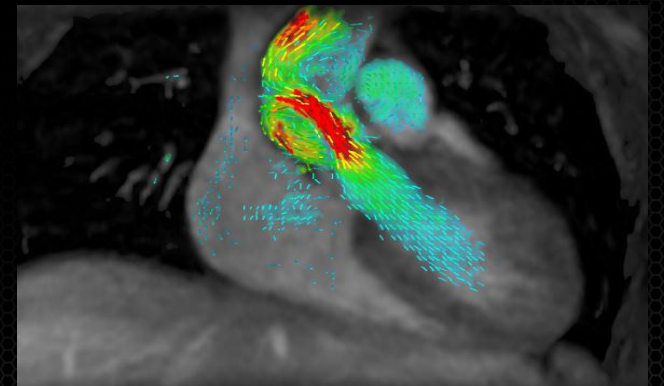


Multi-GPU Compute+Graphics Use Cases

- Image processing
 - Multiple compute GPUs and a low-end display GPU for blitting
- Mixing rasterization with compute
 - Polygonal rendering done in OpenGL and input to compute for further processing
- Visualization for HPC Simulation
 - Numerical simulation distributed across multiple compute GPUs, possibly on remote supercomputer



NVIDIA Index - Seismic Interpretation



Morpheus Medical - Real-time CFD Viz

NVIDIA Maximus Initiative

- Mixing Tesla and Quadro GPUs
 - Tight integration with OEMs and System Integrators
 - Optimized driver paths for GPU-GPU communication



Traditional
Workstation

Simulate (Cluster)

Visualize



NVIDIA® MAXIMUS™
Workstation

Simulate₀

Simulate₁

Simulate₂

Simulate₃

+

+

+

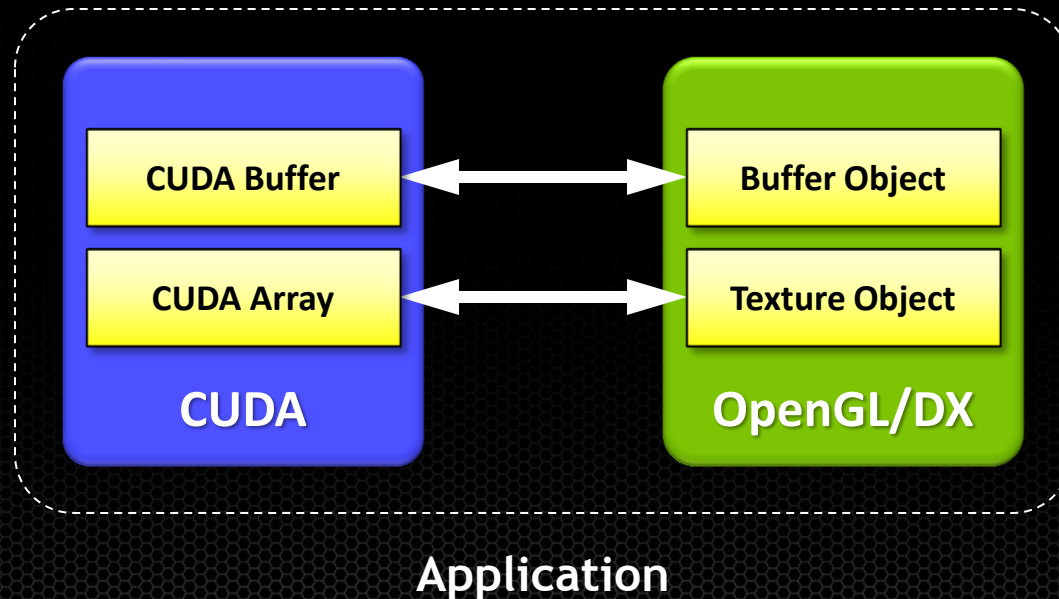
Visualize₀

Visualize₁

Visualize₂

Compute/Graphics interoperability

- Setup the objects in the graphics context.
- Register objects with the compute context.
- Map / Unmap the objects from the compute context.



■ Setup and Registration of Texture Objects:

[illegible]

Code Sample - Simple image interop

■ Mapping between contexts:

```
0
1  cudaArray* texArray;
2
3  while (!done)
4  {
5      cudaGraphicsMapResources(1, &texRes, 0);
6
7      cudaGraphicsSubResourceGetMappedArray(&texArray, texRes, 0, 0);
8      runCUDA(texArray, 0);
9
10     cudaGraphicsUnmapResources(1, &texRes, 0);
11
12     runGL(texId);
13 }
14
```


Code Sample - Simple buffer interop

■ Setup and Registration of Buffer Objects:

```
0  GLuint vboId;
1  cudaGraphicsResource_t vboRes;
2
3  // OpenGL buffer creation...
4  glGenBuffers(1, &vboId);
5  glBindBuffer(GL_ARRAY_BUFFER, vboId);
6  glBufferData(GL_ARRAY_BUFFER, vboSize, 0, GL_STREAM_DRAW);
7  glBindBuffer(GL_ARRAY_BUFFER, 0);
8
9  // Registration with CUDA.
10 cudaGraphicsGLRegisterBuffer(&vboRes, vboId, cudaGraphicsRegisterFlagsNone);
11
12
13
14
```

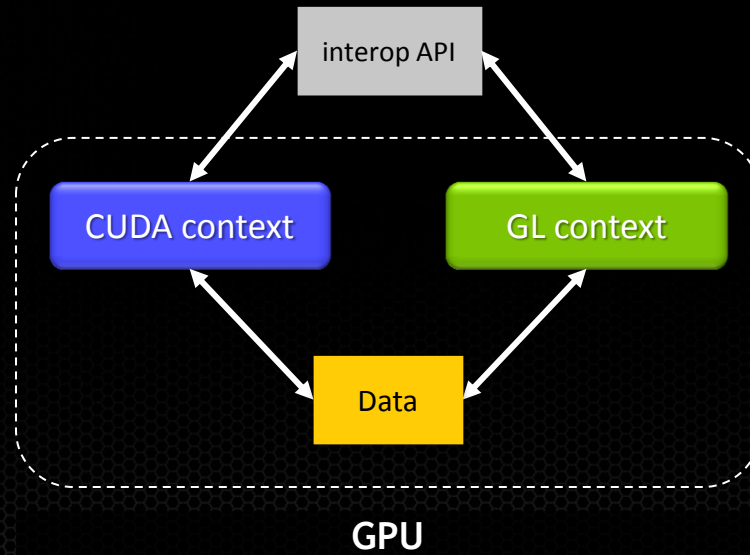

Code Sample - Simple buffer interop

■ Mapping between contexts:

```
0
1 float* vboPtr;
2
3 while (!done)
4 {
5     cudaGraphicsMapResources(1, &vboRes, 0);
6
7     cudaGraphicsResourceGetMappedPointer((void**)&vboPtr, &size, vboRes);
8     runCUDA(vboPtr, 0);
9
10    cudaGraphicsUnmapResources(1, &vboRes, 0);
11
12    runGL(vboId);
13 }
14
```


Resource Behavior: Single-GPU

- The resource is shared. 😊
- Context switch is fast and independent on data size.



Code Sample - Simple buffer interop

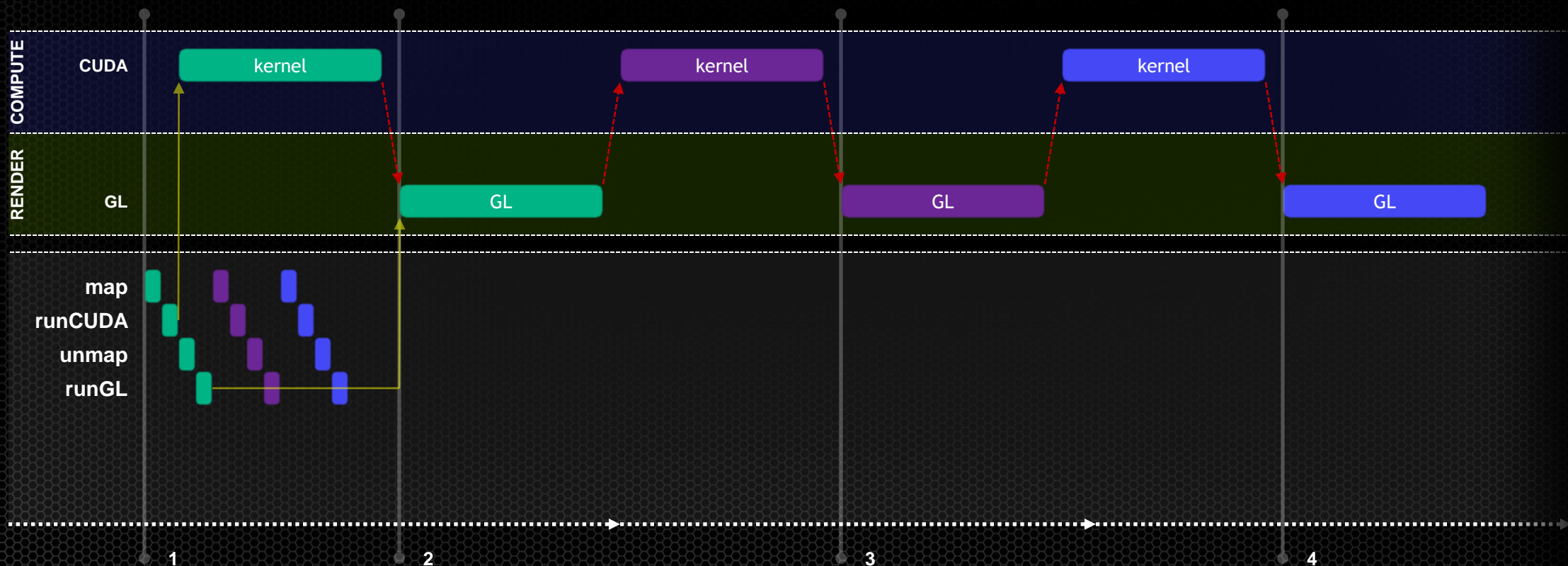
■ Mapping between contexts:

```
0  
1 float* vboPtr;  
2  
3 while (!done)  
4 {  
5     cudaGraphicsMapResources(1, &vboRes, 0);  
6  
7     cudaGraphicsResourceGetMappedPointer((void**)&vboPtr, &size, vboRes);  
8     runCUDA(vboPtr, 0);  
9  
10    cudaGraphicsUnmapResources(1, &vboRes, 0);  
11  
12    runGL(vboId);  
13 }  
14
```

Context-switching happens
when these commands are
processed.

Timeline: Single-GPU

■ Driver-interop



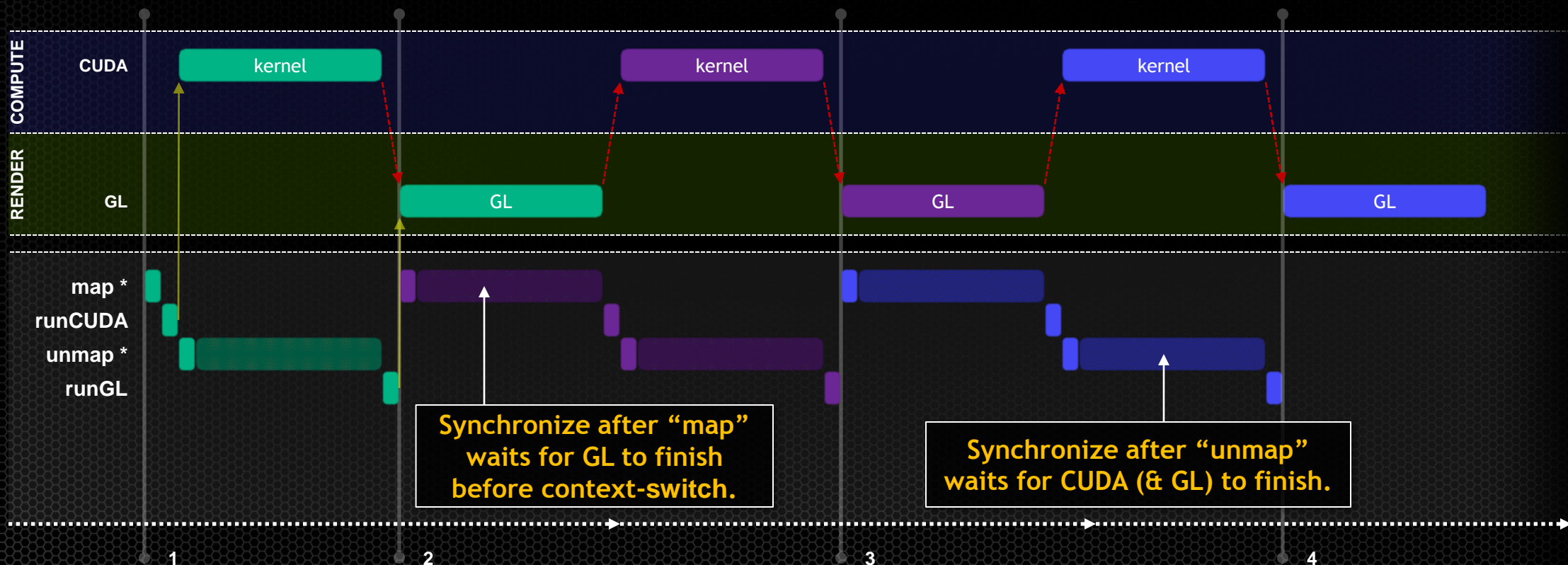
Code Sample - Simple buffer interop

- Adding synchronization for analysis:

```
0
1 float* vboPtr;
2
3 while (!done)
4 {
5     cudaGraphicsMapResources(1, &vboRes, 0);
6     cudaStreamSynchronize(0);
7
8     cudaGraphicsResourceGetMappedPointer((void**)&vboPtr, &size, vboRes);
9     runCUDA(vboPtr, 0);
10
11     cudaGraphicsUnmapResources(1, &vboRes, 0);
12     cudaStreamSynchronize(0);
13
14     runGL(vboId);
15 }
```

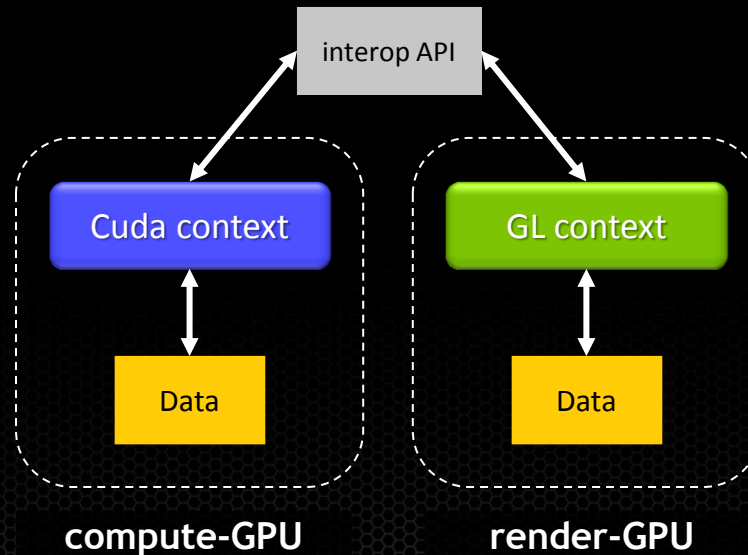

Timeline: Single-GPU

- Driver-interop, synchronous*
 - (we synchronize after map and unmap calls)



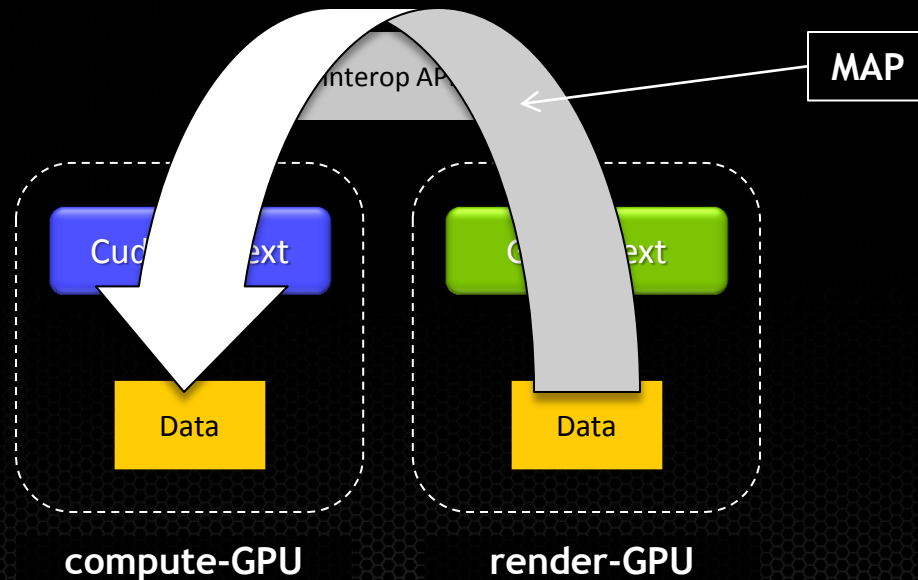
Resource Behavior: Multi-GPU

- Each GPU has a copy of the resource.
- Context-switch is dependent on data size, because driver must copy data. 😞



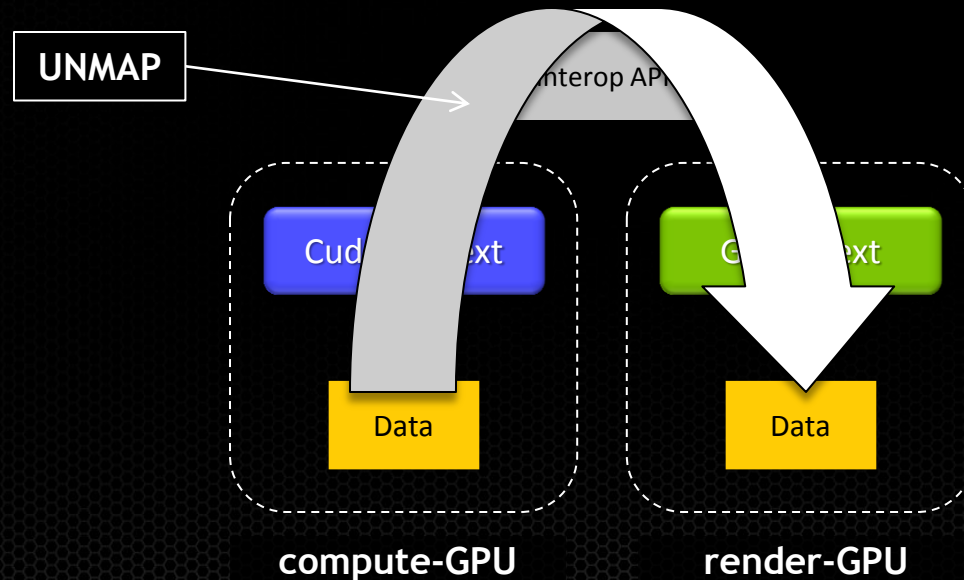
Resource Behavior: Multi-GPU

- Each GPU has a copy of the resource.
- Context-switch is dependent on data size, because driver must copy data. 😞



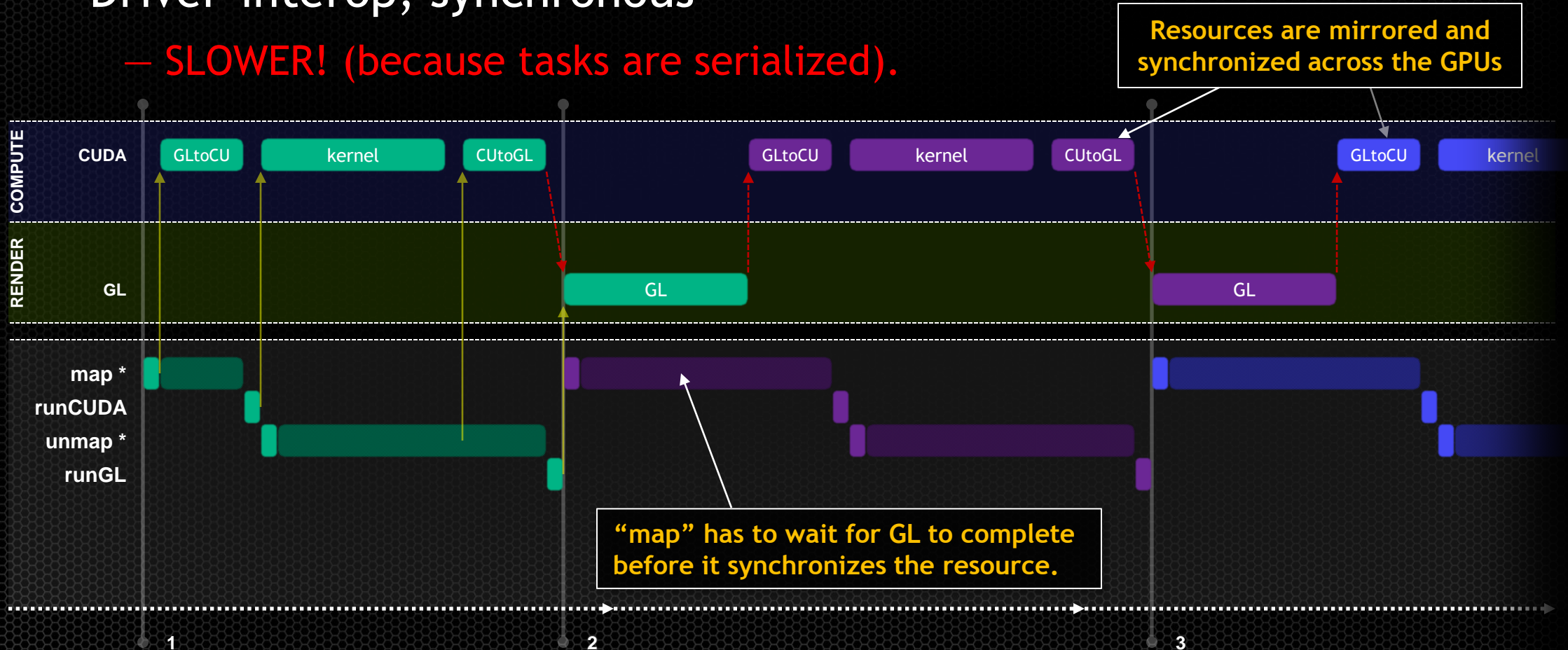
Resource Behavior: Multi-GPU

- Each GPU has a copy of the resource.
- Context-switch is dependent on data size, because driver must copy data. 😞



Timeline: Multi-GPU

- Driver-interop, synchronous*
 - SLOWER! (because tasks are serialized).



Interoperability Methodologies

■ READ-ONLY

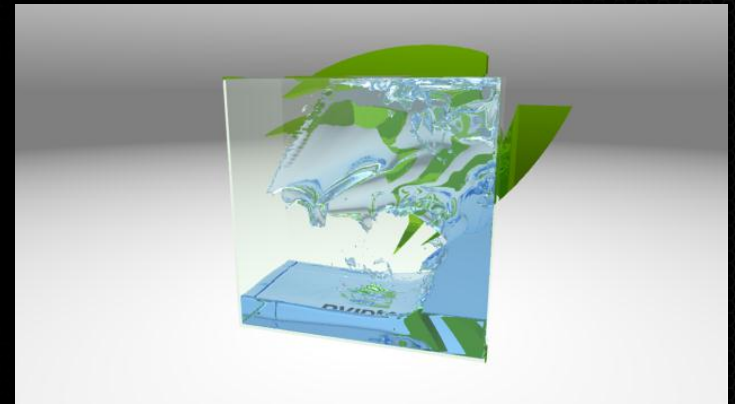
- GL produces... and CUDA consumes.
 - e.g. Post-process the GL render in CUDA.

■ WRITE-DISCARD

- CUDA produces... and GL consumes.
 - e.g. CUDA simulates fluid, and GL renders result.

■ READ & WRITE


- Useful if you want to use the rasterization pipeline.
 - e.g. Feedback loop:
 - `runGL(texture) → framebuffer`
 - `runCUDA(framebuffer) → texture`



Code Sample - WRITE-DISCARD

- CUDA produces... and OpenGL consumes:

```
0 float* vboPtr;
1
2 cudaGraphicsResourceSetMapFlags(vboRes, cudaGraphicsMapFlagsWriteDiscard);
3
4 while (!done)
5 {
6     cudaGraphicsMapResources(1, &vboRes, 0);
7     cudaStreamSynchronize(0);
8
9     cudaGraphicsResourceGetMappedPointer((void**)&vboPtr, &size, vboRes);
10    runCUDA(vboPtr, 0);
11
12    cudaGraphicsUnmapResources(1, &vboRes, 0);
13    cudaStreamSynchronize(0);
14
15    runGL(vboId);
16 }
```

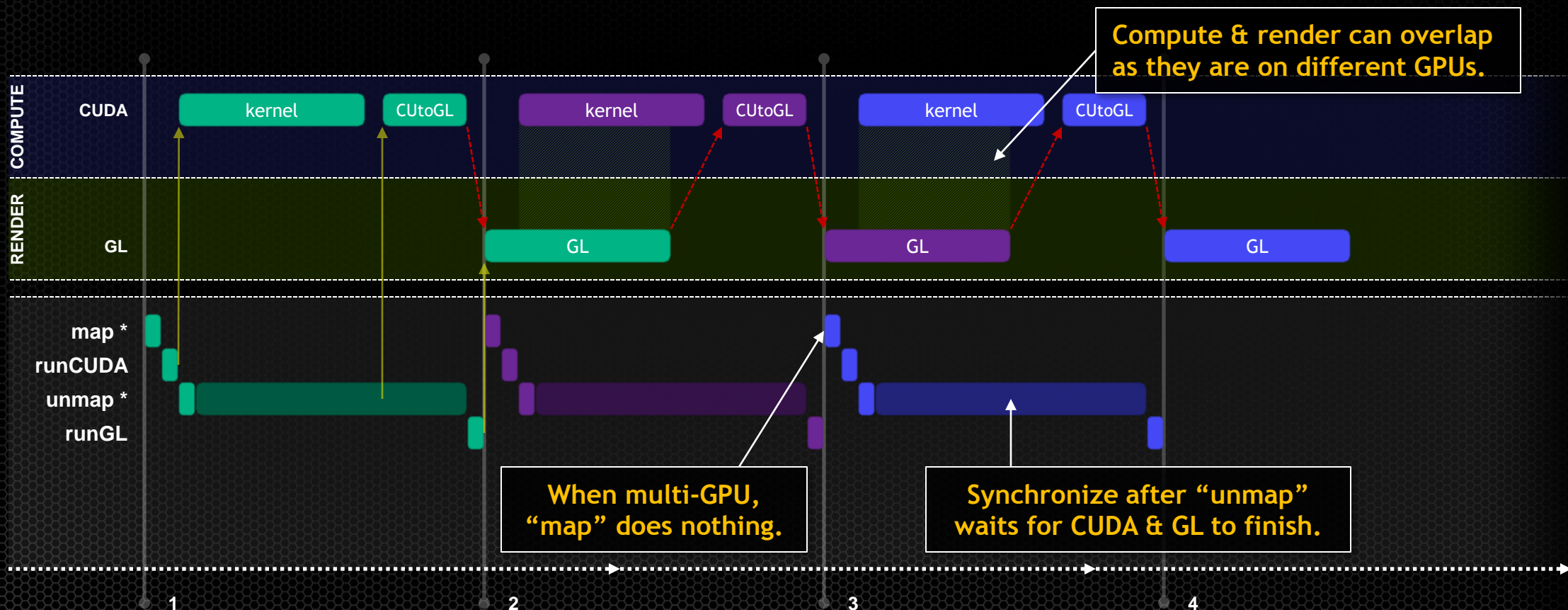


Hint that we do not care about
the previous contents of buffer.

-
- COMPUTE
- CUDA
- kernel
- RENDER
- GL
- GL
- map *
- runCUDA
- unmap *
- runGL
- Context-switch forces serialization.
- Synchronize after "map" waits for GL to finish before context-switch.
- Synchronize after "unmap" waits for CUDA (& GL) to finish.
- 1 2 3 4

Timeline: Multi-GPU

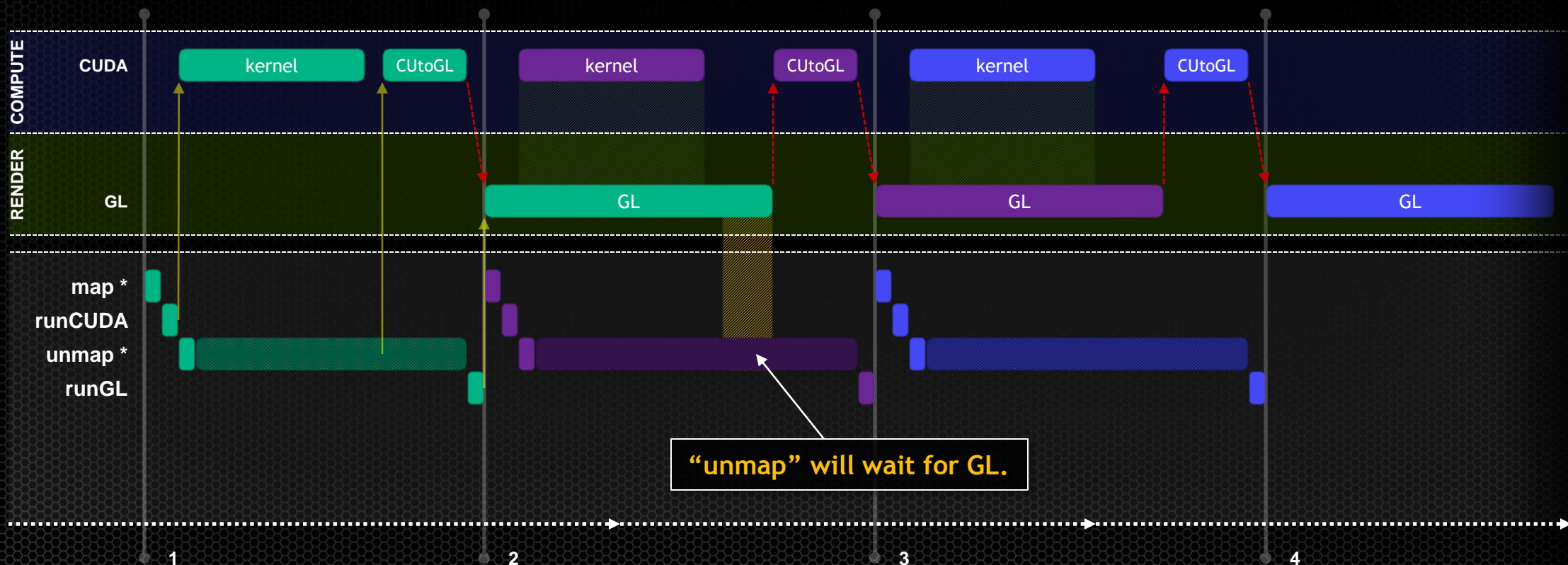
- Driver-interop, synchronous*, **WRITE-DISCARD**



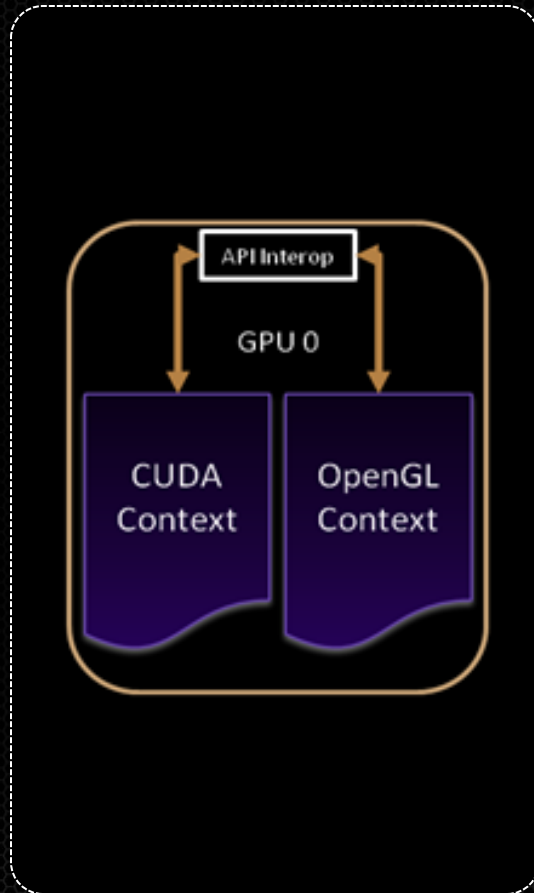
Timeline: Multi-GPU

- Driver-interop, synchronous*, **WRITE-DISCARD**

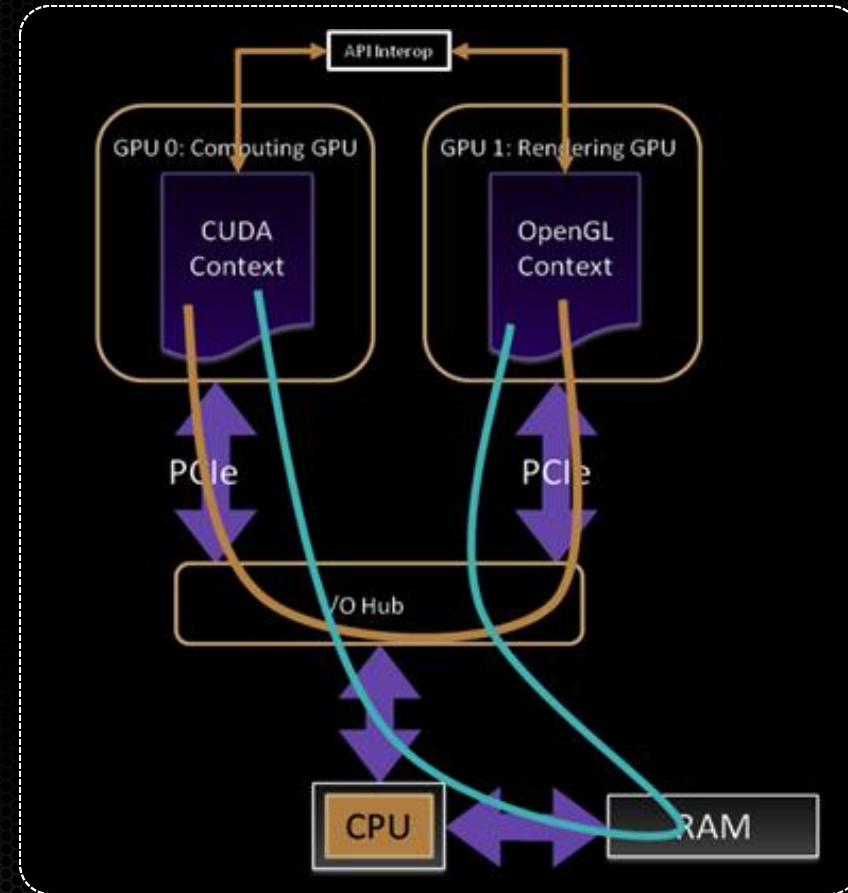
— if render is long...



Driver-Interop: System View



Single-GPU

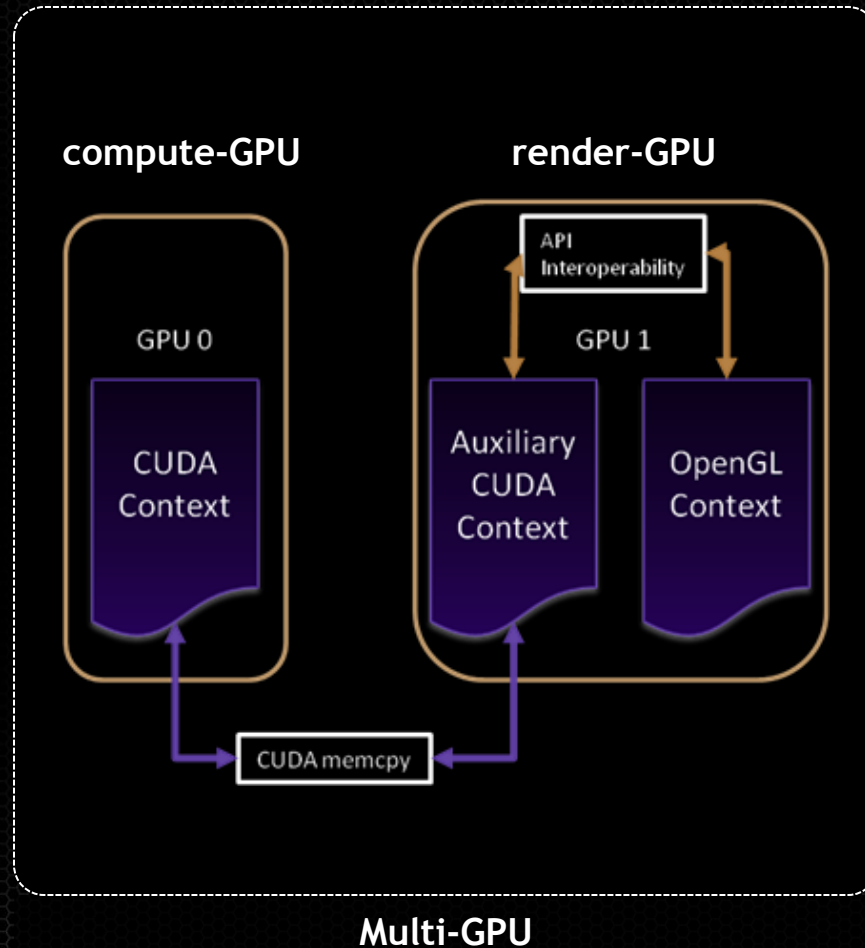


Multi-GPU

Manual-Interop

- Driver-Interop hides all the complexity, but sometimes...
 - We don't want to transfer all data
 - Kernel may only compute sub-regions.
 - We may have a complex system with multiple compute-GPUs and/or render-GPUs.
 - We have application specific pipelining and multi-buffering
 - May have some CPU code in your algorithm between compute and graphics.

Manual-Interrop: System View



Code Sample - Manual-Interop

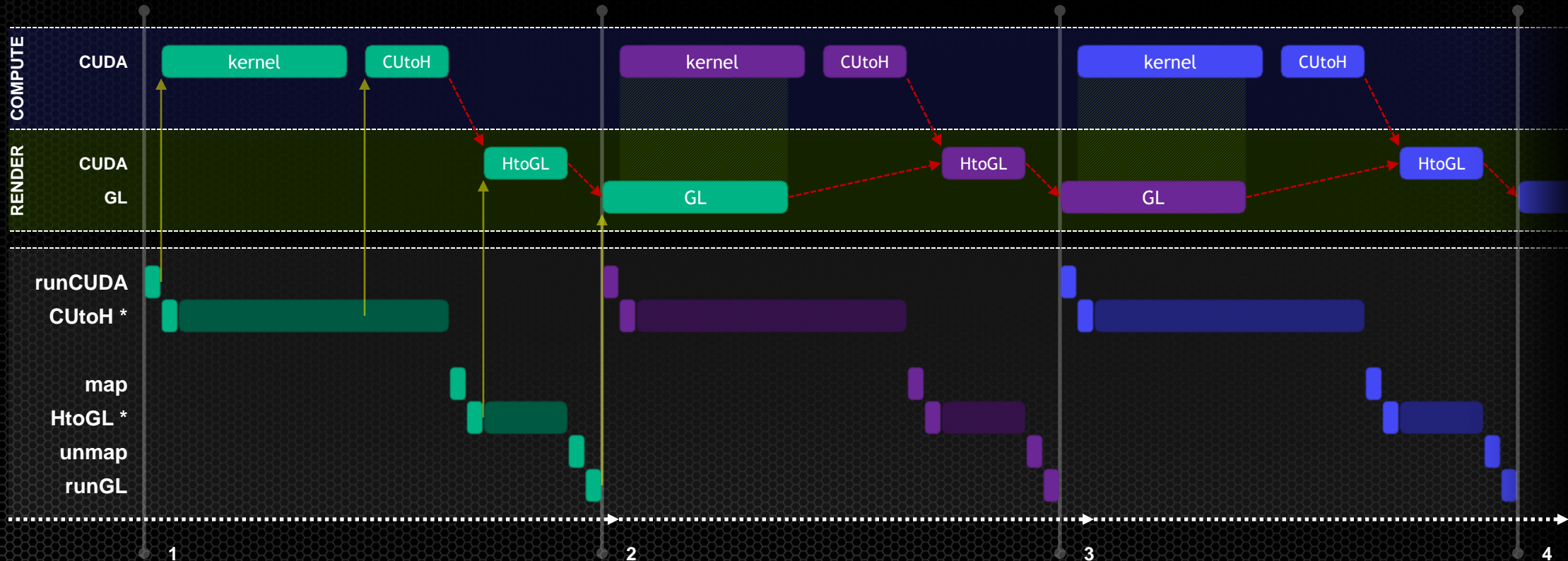
```
0  cudaMalloc((void**)&d_data, vboSize);
1  cudaHostAlloc((void**)&h_data, vboSize, cudaHostAllocPortable);
2
3  while (!done) {
4      // Compute data in temp buffer, and copy to host...
5      runCUDA(d_data, 0);
6      cudaMemcpyAsync(h_data, d_data, vboSize, cudaMemcpyDeviceToHost, 0);
7      cudaStreamSynchronize(0);
8
9      // Map the render-GPU's resource and upload the host buffer...
10     cudaSetDevice(renderGPU);
11     cudaGraphicsMapResources(1, &vboRes, 0);
12     cudaGraphicsResourceGetMappedPointer((void**)&vboPtr, &size, vboRes);
13     cudaMemcpyAsync(vboPtr, h_data, size, cudaMemcpyHostToDevice, 0);
14     cudaGraphicsUnmapResources(1, &vboRes, 0);
15     cudaStreamSynchronize(0);
16     cudaSetDevice(computeGPU);
17
18     runGL(vboId);
19 }
```



Create a temporary buffer
in pinned host-memory.

Timeline: Multi-GPU

- Manual-interop, synchronous*, WRITE-DISCARD

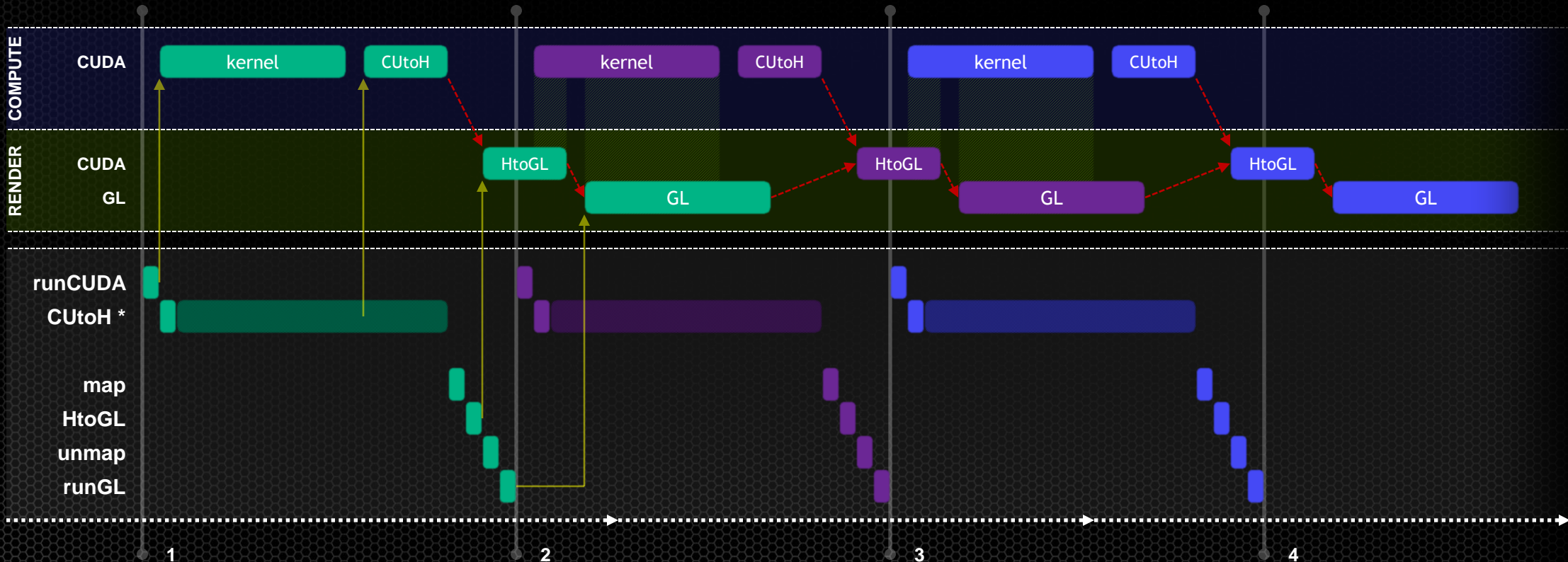


Code Sample - Manual Interop (Async)

```
0  cudaMalloc((void**)&d_data, vboSize);
1  cudaHostAlloc((void**)&h_data, vboSize, cudaHostAllocPortable);
2
3  while (!done) {
4      // Compute data in temp buffer, and copy to host...
5      runCUDA(d_data, 0);
6      cudaMemcpyAsync(h_data, d_data, vboSize, cudaMemcpyDeviceToHost, 0);
7      cudaStreamSynchronize(0);
8
9      // Map the render-GPU's resource and upload the host buffer...
10     cudaSetDevice(renderGPU);
11     cudaGraphicsMapResources(1, &vboRes, 0);
12     cudaGraphicsResourceGetMappedPointer((void**)&vboPtr, &size, vboRes);
13     cudaMemcpyAsync(vboPtr, h_data, size, cudaMemcpyHostToDevice, 0);
14     cudaGraphicsUnmapResources(1, &vboRes, 0);
15     cudaSetDevice(computeGPU);
16     // cudaStreamSynchronize(0);
17
18     runGL(vboId);
19 }
```


Timeline: Multi-GPU

- Manual-interop, **asynchronous**, WRITE-DISCARD



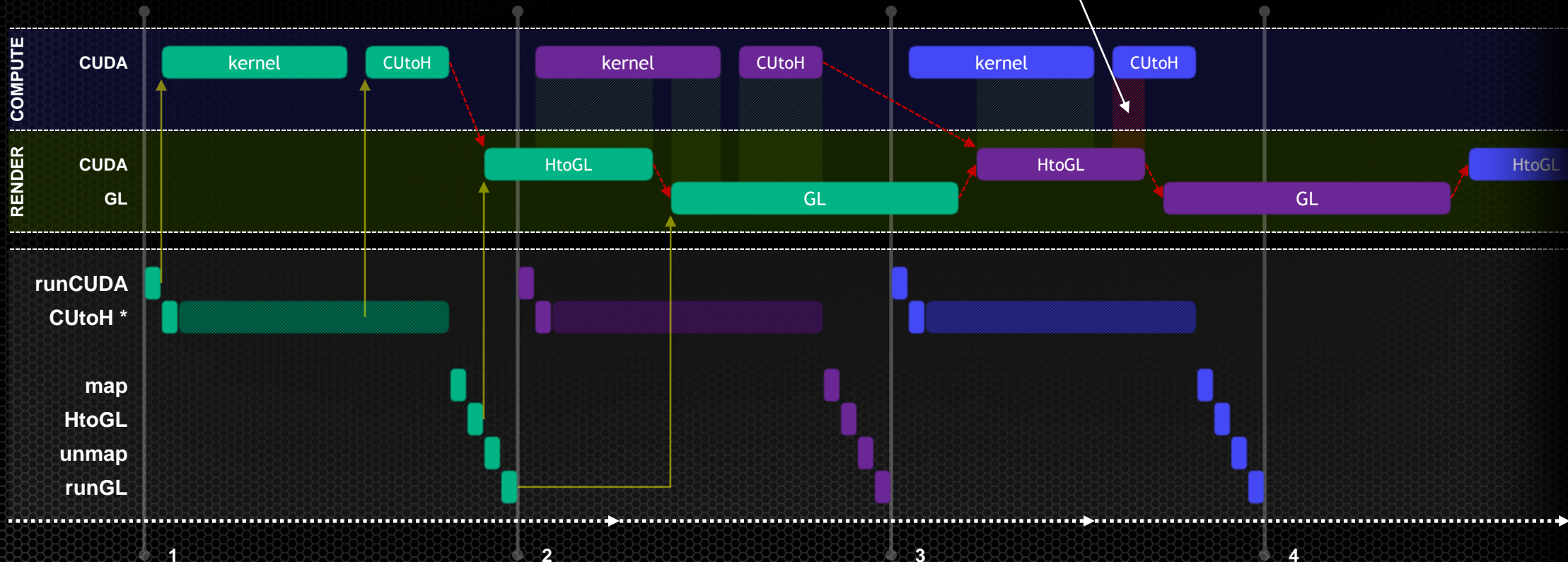
Timeline: Multi-GPU

- Manual-interop, **asynchronous**, WRITE-DISCARD

— if render is long...

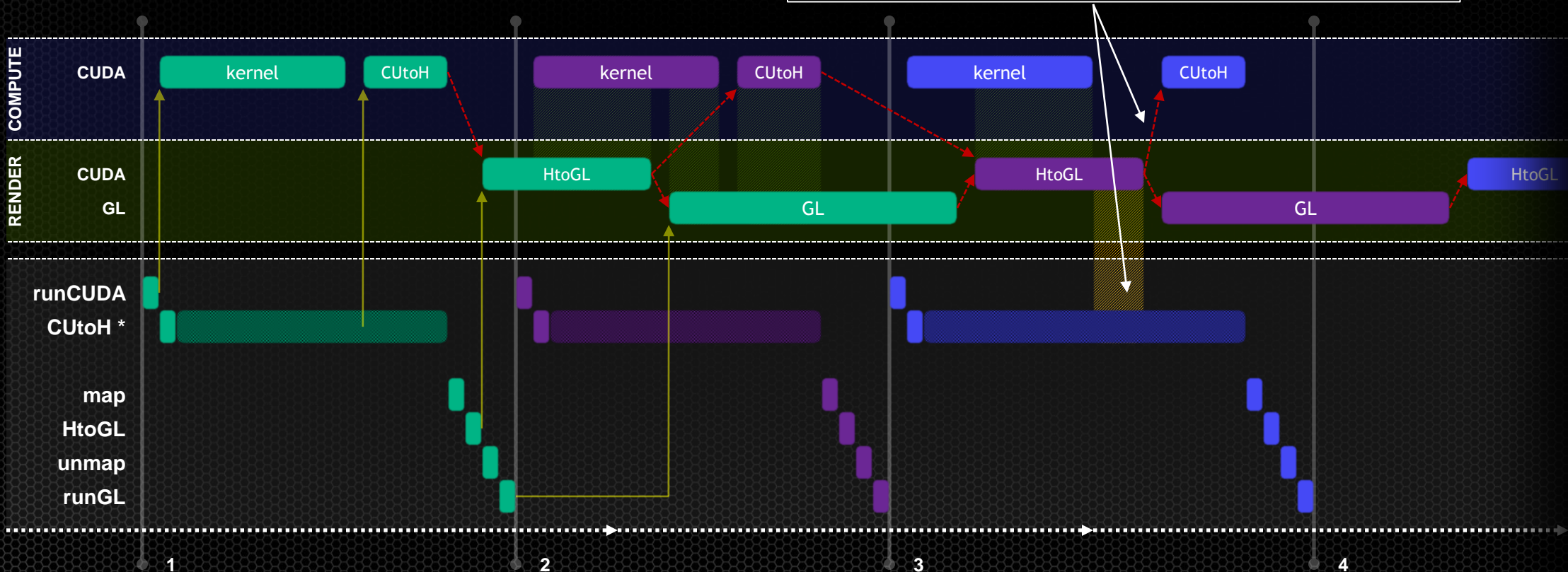
We are downloading while uploading!

Drifting out of sync!



- if render is long...

Synchronization must *also* wait for HtoGL to finish



Synchronization across GPUs

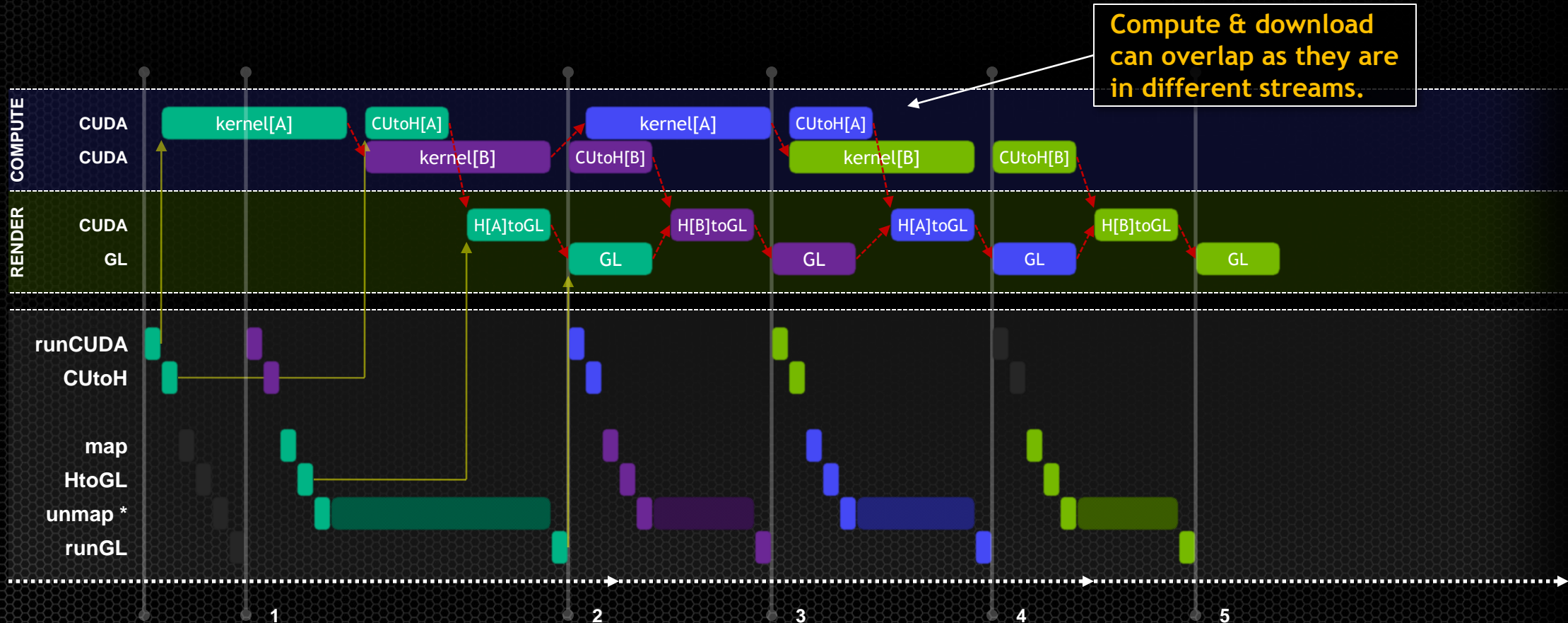
- Streams and events are per device
 - Determined by the GPU that is current when created.
- Stream GPU must be set current for
 - Launching kernel to a stream
 - Recording events to a stream
- Agnostic to the current GPU
 - Memcpy can be launched on any stream
 - Synchronize/Query of Events

Code Sample - Manual Interop (safe Async)

```
0  while (!done) {
1      // Compute the data in a temp buffer, and copy to a host buffer...
2      runCUDA(d_data, 0);
3      cudaStreamWaitEvent(0, uploadFinished, 0);
4      cudaMemcpyAsync(h_data, d_data, vboSize, cudaMemcpyDeviceToHost, 0);
5      cudaStreamSynchronize(0);
6
7      // Map the render-GPU's resource and upload the host buffer...
8      // (all commands must be asynchronous.)
9      cudaSetDevice(renderGPU);
10     cudaGraphicsMapResources(1, &vboRes, 0);
11     cudaGraphicsResourceGetMappedPointer((void**)&vboPtr, &size, vboRes);
12     cudaMemcpyAsync(vboPtr, h_data, size, cudaMemcpyHostToDevice, 0);
13     cudaGraphicsUnmapResources(1, &vboRes, 0);
14     cudaEventRecord(uploadFinished, 0);
15     cudaSetDevice(computeGPU);
16
17     runGL(vboId);
18 }
```


Timeline: Multi-GPU

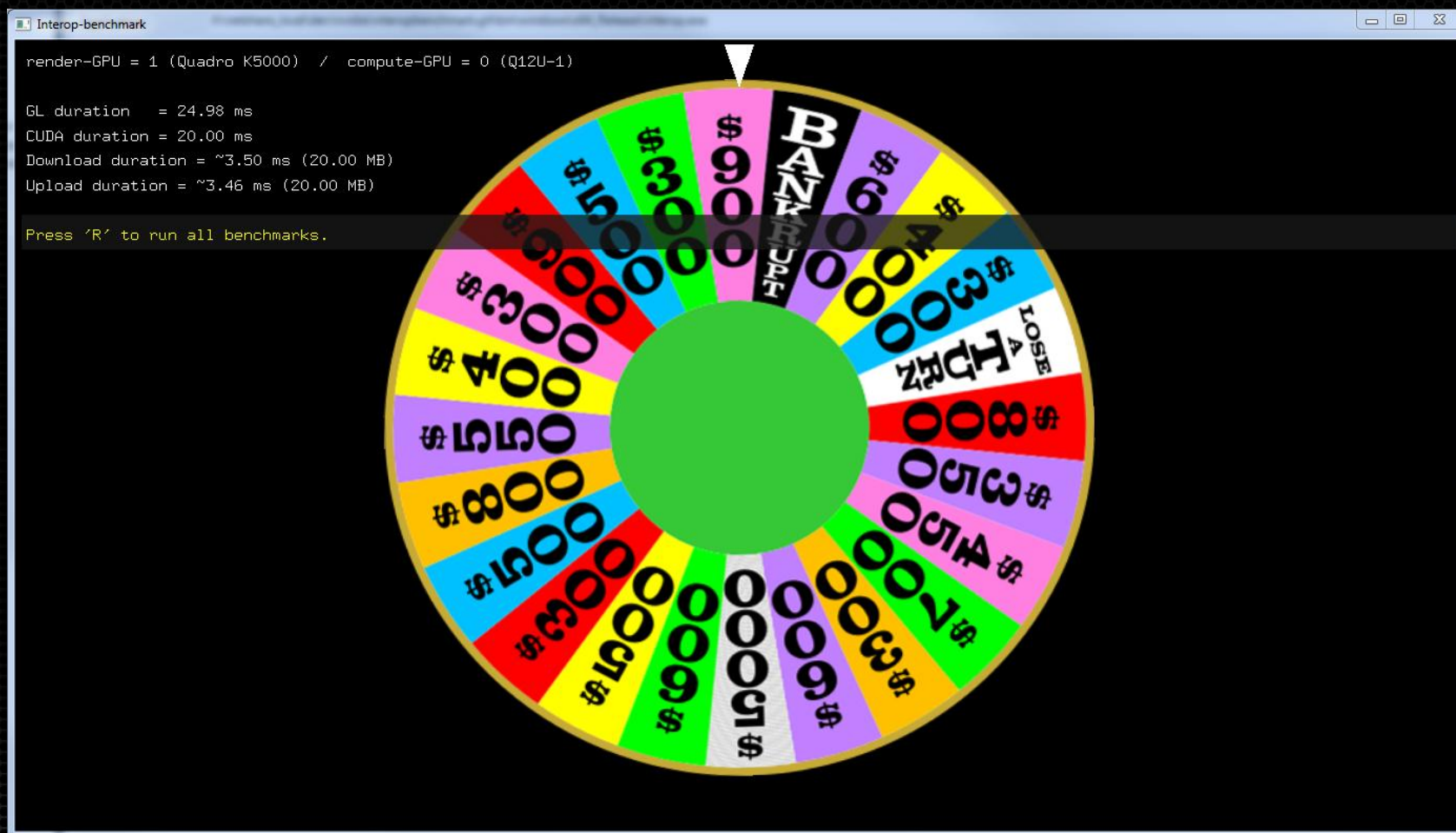
- Manual-interop, WRITE-DISCARD, **flipping CUDA**



Code Sample - Manual Interop (flipping)

```
0  int read = 1, write = 0;
1
2  while (!done) {
3      // Compute the data in a temp buffer, and copy to a host buffer...
4      cudaStreamWaitEvent(custream[write], kernelFinished[read]);
5      runCUDA(d_data[write], custream[write]);
6      cudaEventRecord(kernelFinished[write], custream[write]);
7      cudaStreamWaitEvent(custream[write], uploadFinished[read]);
8      cudaMemcpyAsync(h_data[write], d_data[write], vboSize, cudaMemcpyDeviceToHost, custream[write]);
9      cudaEventRecord(downloadFinished[write], custream[write]);
10
11     // Map the renderGPU's resource and upload the host buffer...
12     cudaSetDevice(renderGPU);
13     cudaGraphicsMapResources(1, &vboRes, glstream);
14     cudaGraphicsResourceGetMappedPointer((void*)&vboPtr, &size, vboRes);
15     cudaStreamWaitEvent(glstream, downloadFinished[read]);
16     cudaMemcpyAsync(vboPtr, h_data[read], size, cudaMemcpyHostToDevice, glstream);
17     cudaGraphicsUnmapResources(1, &vboRes, glstream);
18     cudaEventRecord(uploadFinished[read], glstream);
19     cudaStreamSynchronize(glstream); // Sync for easier analysis!
20     cudaSetDevice(computeGPU);
21
22     runGL(vboId);
23     swap(&read, &write);
24 }
```


Demo



Demo - results

- runCUDA (20ms)
- runGL (10ms)
- copy (10ms)
- Single-GPU
 - Driver-interop = 30ms
- Multi-GPU
 - Driver-interop = 36ms
 - Async Manual-interop = 32ms
 - Flipped Manual-interop = 22ms

Too large data size
makes multi-GPU
interop worse.

Overlapping the
download helps
us break even.

But using streams
and flipping is a
significant win!

Some final thoughts on Interoperability.

- Similar considerations when OpenGL produces and CUDA consumes
 - Use `cudaGraphicsMapFlagsReadOnly`
- Avoid synchronizing GPUs!
 - Watch out for Windows's WDDM implicit synchronization on unmap.
- CUDA-OpenGL interoperability can perform slower if OpenGL context spans multiple GPUs.
- Context-switch performance varies with system configuration and OS.

Resources

- CUDA samples/documentation
 - <http://developer.nvidia.com/cuda-downloads>
- Interop-benchmark demo
 - http://bitbucket.org/wbraithwaite_nvidia/interopbenchmark
- OpenGL Insights, Patrick Cozzi, Christophe Riccio, 2012. ISBN 1439893764.
 - www.openglinsights.com
- Thursday 12pm, Multi-GPU programming for Visual Computing (Part 2)