**GPU** TECHNOLOGY CONFERENCE

# Mixing graphics and compute with multiple GPUs

# Agenda

- Compute and Graphics Interoperability
- Interoperability at a system level
- Application design considerations

# Putting Graphics & Compute together

- Compute and Visualize the same data

# Compute/Graphics interoperability

- Set of compute API functions
  - Graphics sets up the objects
  - Register/Unregister the objects with compute context
  - Mapping/Unmapping of the objects to/from the compute context every frame

**Application**

| CUDA | | OpenGL/DX |
|---|---|---|
| CUDA Linear Memory | ⟷ | Buffer |
| CUDA Array | ⟷ | Texture |

# Simple OpenGL-CUDA interop sample: Setup and Register of Buffer Objects

```
GLuint imagePBO;

cudaGraphicsResource_t    cudaResourceBuf;

//OpenGL buffer creation

glGenBuffers(1, &imagePBO);

glBindBuffer(GL_PIXEL_UNPACK_BUFFER_ARB, imagePBO);

glBufferData(GL_PIXEL_UNPACK_BUFFER_ARB, size, NULL, GL_DYNAMIC_DRAW);

glBindBuffer(GL_PIXEL_UNPACK_BUFFER_ARB,0);

//Registration with CUDA

cudaGraphicsGLRegisterBuffer(&cudaResourceBuf, imagePBO,
    cudaGraphicsRegisterFlagsNone);
```
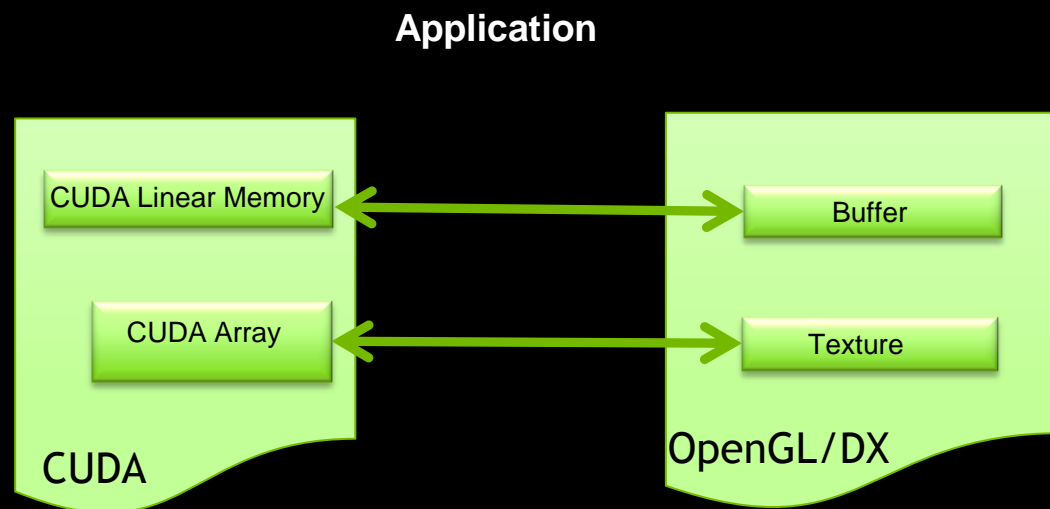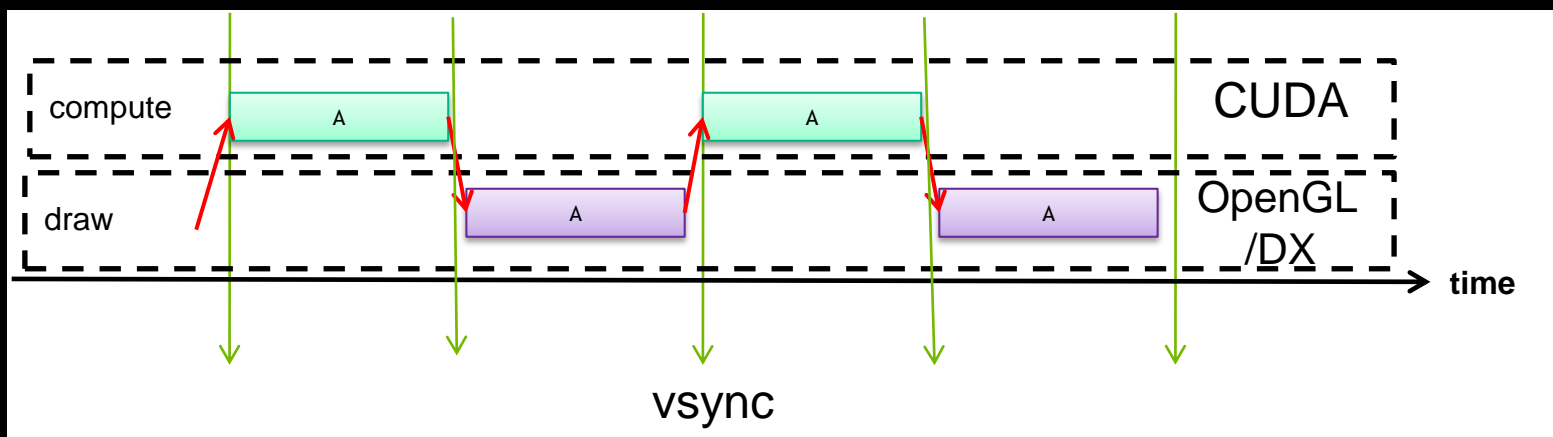
# Simple OpenGL-CUDA interop sample: Setup and Register of Texture Objects

```
GLuint imageTex;

cudaGraphicsResource_t    cudaResourceTex;

//OpenGL texture creation

glGenTextures(1, &imageTex);

glBindTexture(GL_TEXTURE_2D, imageTex);

//set texture parameters here

glTexImage2D(GL_TEXTURE_2D,0, GL_RGBA8UI_EXT, width, height, 0,
    GL_RGBA_INTEGER_EXT, GL_UNSIGNED_BYTE, NULL);

glBindTexture(GL_TEXTURE_2D, 0);

//Registration with CUDA

cudaGraphicsGLRegisterImage (&cudaResourceTex, imageTex, GL_TEXTURE_2D,
    cudaGraphicsRegisterFlagsNone);
```

# Simple OpenGL-CUDA interop sample

```
unsigned char *memPtr;
cudaArray *arrayPtr;
while (!done) {
    cudaGraphicsMapResources(1, &cudaResourceTex, cudaStream);
    cudaGraphicsMapResources(1, &cudaResourceBuf, cudaStream);
    cudaGraphicsSubResourceGetMappedArray(&cudaArray, cudaResourceTex, 0, 0);
    cudaGraphicsResourceGetMappedPointer((void **)&memPtr, &size, cudaResourceBuf);
    doWorkInCUDA(cudaArray, memPtr, cudaStream); //asynchronous
    cudaGraphicsUnmapResources(1, & cudaResourceTex, cudaStream);
    cudaGraphicsUnmapResources(1, & cudaResourceBuf, cudaStream);
    doWorkInGL(imagePBO, imageTex); //asynchronous
}
```
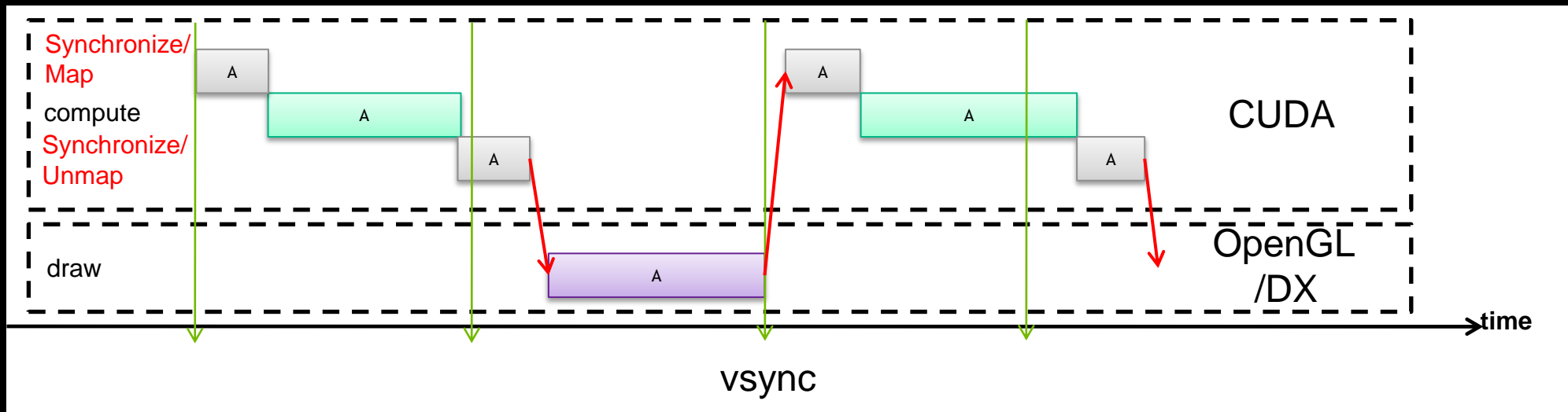
# Interoperability behavior:single GPU

- The resource is shared
- Tasks are serialized

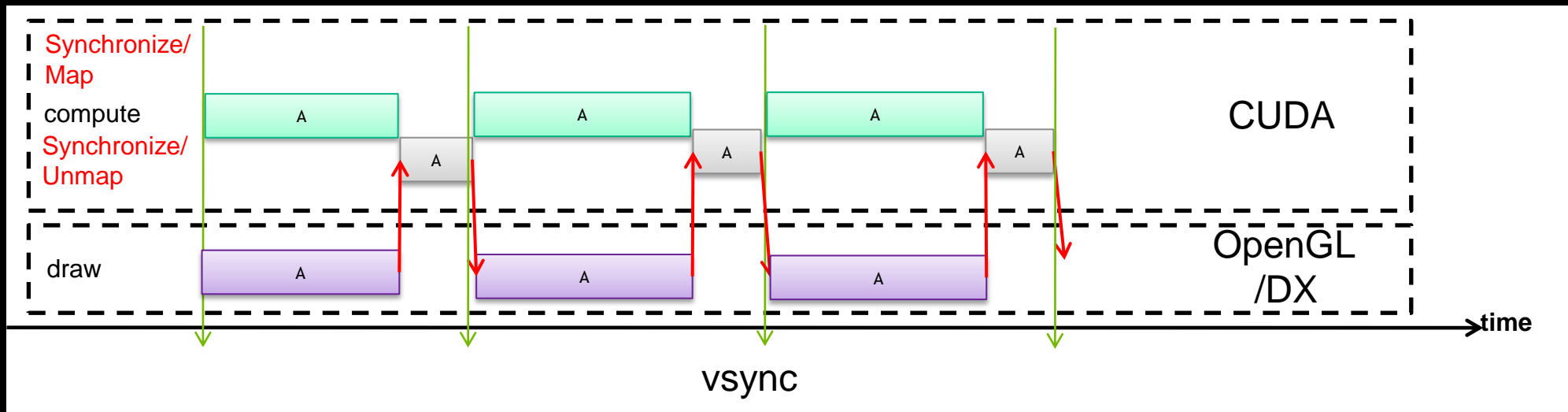# Interoperability behavior: multiple GPUs

- Each context owns a copy of the resource
- Tasks are serialized
- map/unmap might do a host side synchronization

# Interoperability behavior: multiple GPUs Improvements

- If one of the APIs is a producer and another is a consumer then the tasks can overlap.
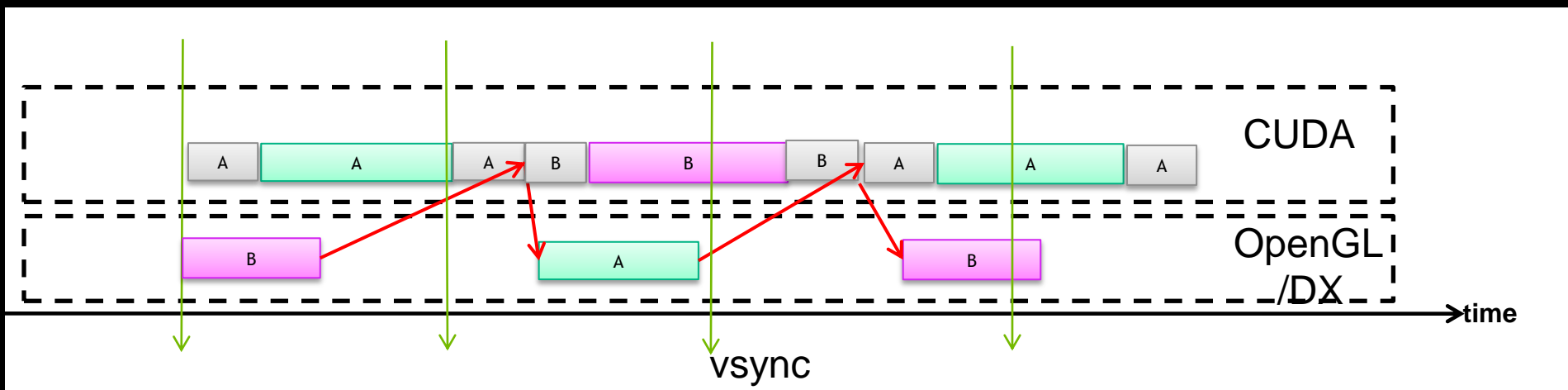
CUDA as a producer example:



| | CUDA |
| --- | --- |
| Synchronize/Map compute Synchronize/Unmap | |
| draw | OpenGL /DX |

vsync

# Simple OpenGL-CUDA interop sample

- Use mapping hint with cudaGraphicsResourceSetMapFlags()

cudaGraphicsMapFlagsReadOnly/cudaGraphicsMapFlagsWriteDiscard:

```
unsigned char *memPtr;
cudaGraphicsResourceSetMapFlags(cudaResourceBuf, cudaGraphicsMapFlagsWriteDiscard)
while (!done) {
    cudaGraphicsMapResources(1, &cudaResourceBuf, cudaStream);
    cudaGraphicsResourceGetMappedPointer((void **)&memPtr, &size, cudaResourceBuf);
    doWorkInCUDA(memPtr, cudaStream);//asynchronous
    cudaGraphicsUnmapResources(1, &cudaResourceBuf, cudaStream);
    doWorkInGL(imagePBO);  //asynchronous
}
```

# Interoperability behavior: multiple GPUs Improvements

- If the graphics and compute are interdependant use ping-pong buffers for task overlap

# Simple OpenGL-CUDA interop sample

- **ping-pong buffers:**

```
unsigned char *memPtr;

int count = 0;

while (!done) {

    cudaResourceBuf = (count%2) ? cudaResourceBuf1 : cudaResourceBuf2;

    imagePBO = (count%2) ? imagePBO1 : imagePBO2;

    cudaGraphicsMapResources(1, &cudaResourceBuf, cudaStream);

    cudaGraphicsResourceGetMappedPointer((void **)&memPtr, &size, cudaResourceBuf);

    doWorkInCUDA(memPtr, cudaStream); //asynchronous

    cudaGraphicsUnmapResources(1, & cudaResourceBuf, cudaStream);

    doWorkInGL(imagePBO); //asynchronous

    count++;

}
```
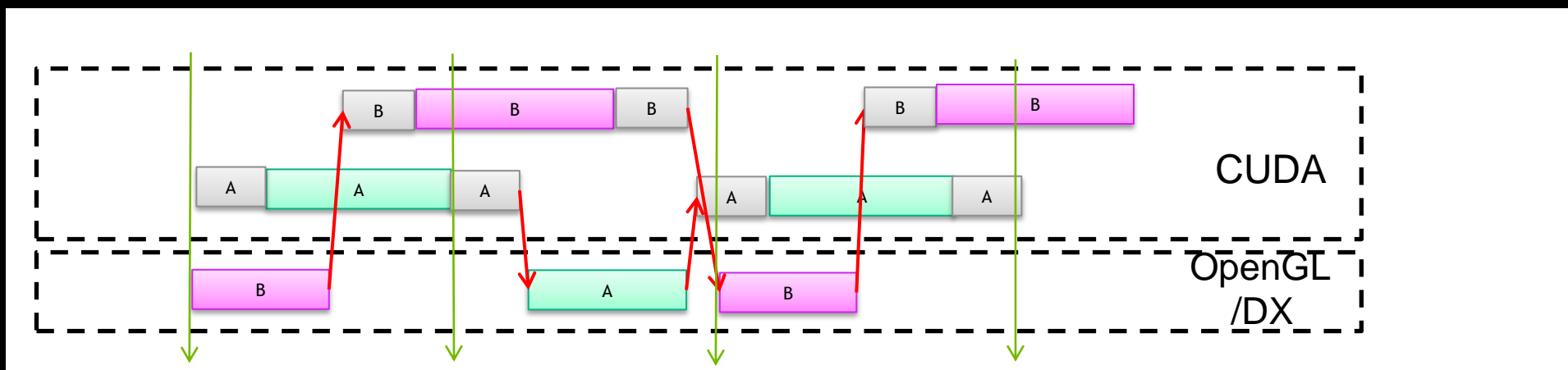
# Interoperability behavior: multiple GPUs Improvements

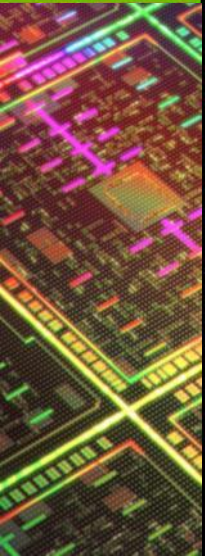- Utilize the copy engines to hide the data transfer latency cost.

# Simple OpenGL-CUDA interop sample

- ▪ Use streams

```
unsigned char *memPtr;

int count = 0;

while (!done) {

    cudaResourceBuf = (count%2) ? cudaResourceBuf1 : cudaResourceBuf2;

    imagePBO = (count%2) ? imagePBO1 : imagePBO2;

    cudaStream= (count%2) ? cudaStream1 : cudaStream2;

    cudaGraphicsMapResources(1, &cudaResourceBuf, cudaStream);

    cudaGraphicsResourceGetMappedPointer((void **)&memPtr, &size, cudaResourceBuf);

    doWorkInCUDA(memPtr, cudaStream);//asynchronous

    cudaGraphicsUnmapResources(1, & cudaResourceBuf, cudaStream);

    doWorkInGL(imagePBO); //asynchronous

    count++;

}
```
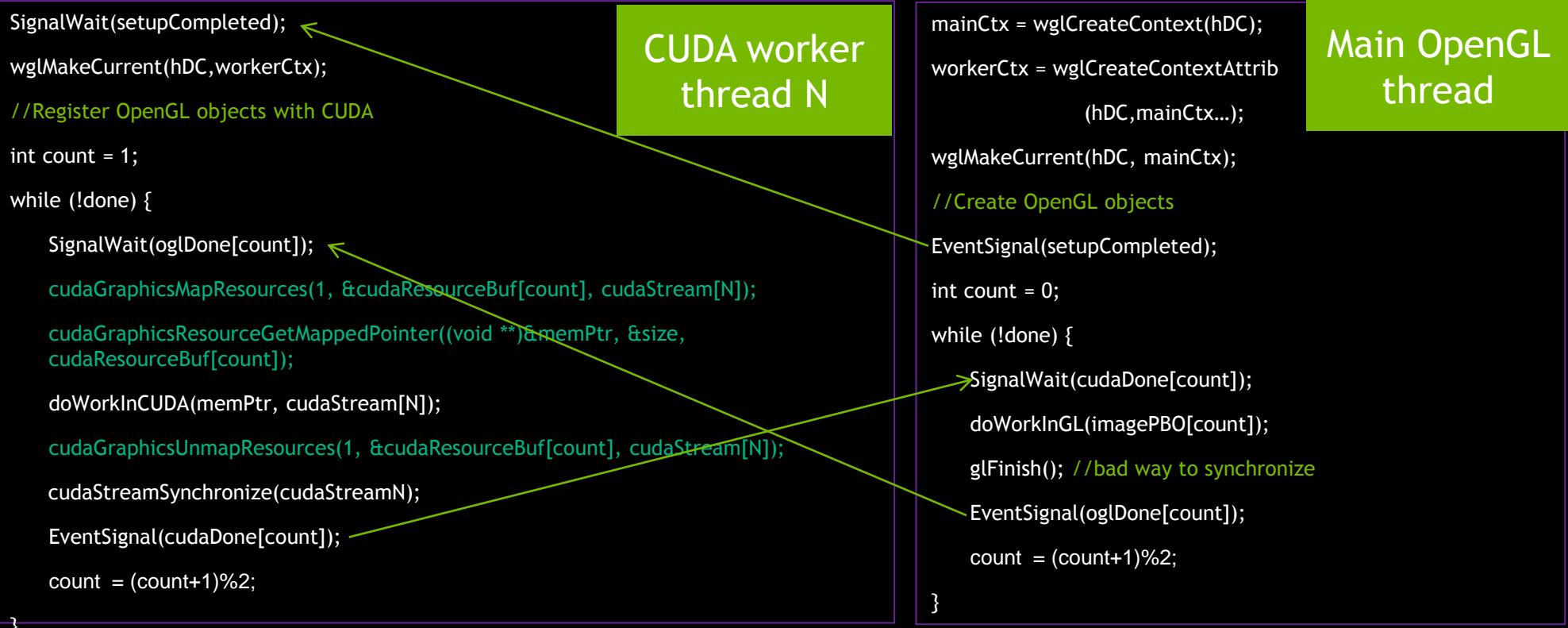
# Simple OpenGL-CUDA interop sample: Quick Summary of Techniques

- Mapping hints

- Ping-pong buffers

- Streams

- This is enough if:
  - You can use the map/unmap hints
  - map/unmap is CPU asynchronous
  - You are afraid of multiple threads
  - You developed the whole application

# Application Example:pseudocode

- Multithreaded OpenGL centric application: Autodesk Maya with a CUDA plug-in

**CUDA worker thread N**

```
SignalWait(setupCompleted);
wglMakeCurrent(hDC,workerCtx);
//Register OpenGL objects with CUDA

int count = 1;
while (!done) {
    SignalWait(oglDone[count]);
    cudaGraphicsMapResources(1, &cudaResourceBuf[count], cudaStream[N]);
    cudaGraphicsResourceGetMappedPointer((void **)&memPtr, &size,
    cudaResourceBuf[count]);
    doWorkInCUDA(memPtr, cudaStream[N]);
    cudaGraphicsUnmapResources(1, &cudaResourceBuf[count], cudaStream[N]);
    cudaStreamSynchronize(cudaStreamN);
    EventSignal(cudaDone[count]);
    count  = (count+1)%2;
}
```

**Main OpenGL thread**

```
mainCtx = wglCreateContext(hDC);
workerCtx = wglCreateContextAttrib
                (hDC,mainCtx…);
wglMakeCurrent(hDC, mainCtx);
//Create OpenGL objects
EventSignal(setupCompleted);
int count = 0;
while (!done) {
    SignalWait(cudaDone[count]);
    doWorkInGL(imagePBO[count]);
    glFinish(); //bad way to synchronize
    EventSignal(oglDone[count]);
    count  = (count+1)%2;
}
```

# Application Example:pseudocode

- Multithreaded CUDA centric application: Adobe Premiere Pro with an OpenGL plugin
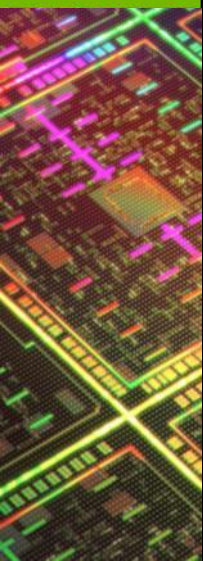
**Main CUDA thread**

```
int count = 1;
while (!done) {
    SignalWait(oglDone[count]);
    doWorkInCUDA(memPtr, NULL);
    EventSignal(cudaDone[count]);
    count  = (count+1)%2;
}
```

**OpenGL Worker thread**

```
mainCtx = wglCreateContext(hDC);
wglMakeCurrent(hDC, mainCtx);
//Register OpenGL objects with CUDA
int count = 0;
while (!done) {
    SignalWait(cudaDone[count]);
    cudaGraphicsUnmapResources(1, &cudaResourceBuf[count], NULL);
    doWorkInGL(imagePBO[count]);
    cudaGraphicsMapResources(1, &cudaResourceBuf[count], NULL);
    cudaGraphicsResourceGetMappedPointer((void **)&memPtr, &size,
    cudaResourceBuf[count]);
    EventSignal(oglDone[count]);
    count  = (count+1)%2;
}
```

# Application design considerations

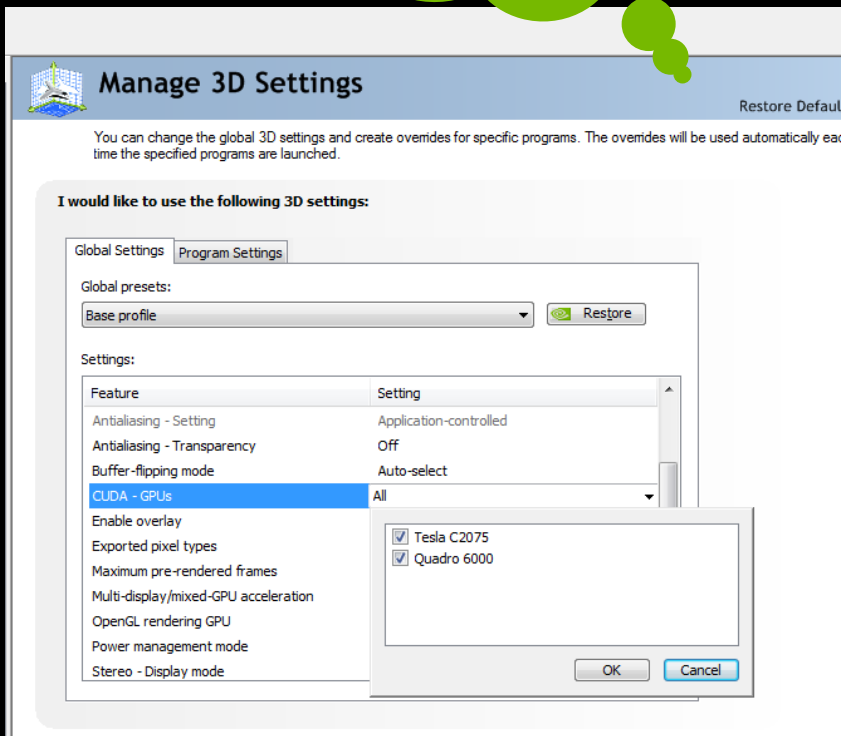- Use cudaD3D[9|10|11]GetDevices/cudaGLGetDevices to chose the right device to provision for multi-GPU environments.

- Avoid synchronized GPUs for CUDA!

- CUDA-OpenGL interop will perform slower if OpenGL context spans multiple GPU!
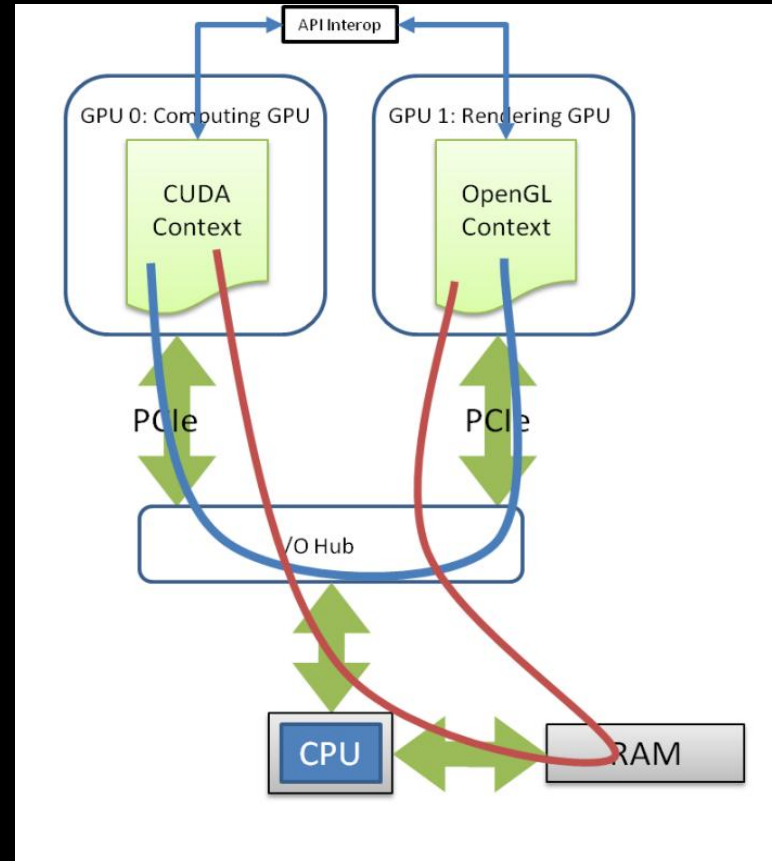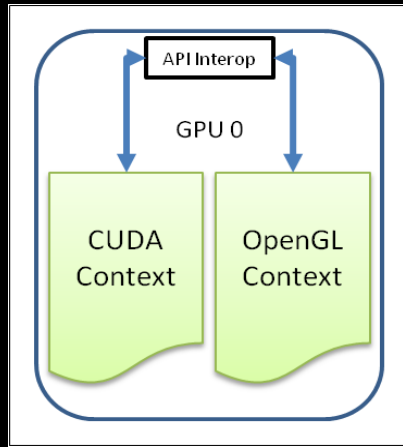
# Application design considerations

- Allow users to specify the GPUs!
  - Typical heuristics:
    - TCC mode
    - GPU #
    - available memory
    - # of processing units
  - Affecting factors: OS, ECC, TCC mode
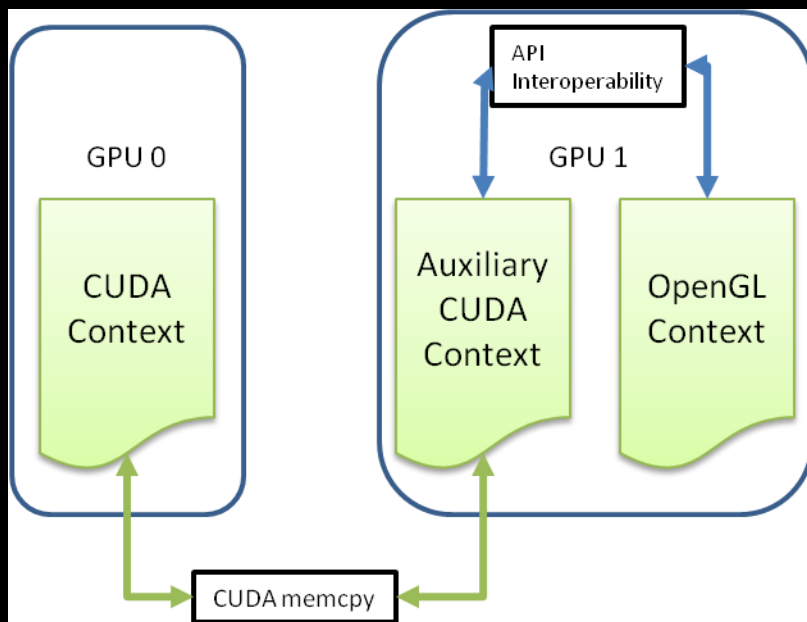
Don't make your users go here:

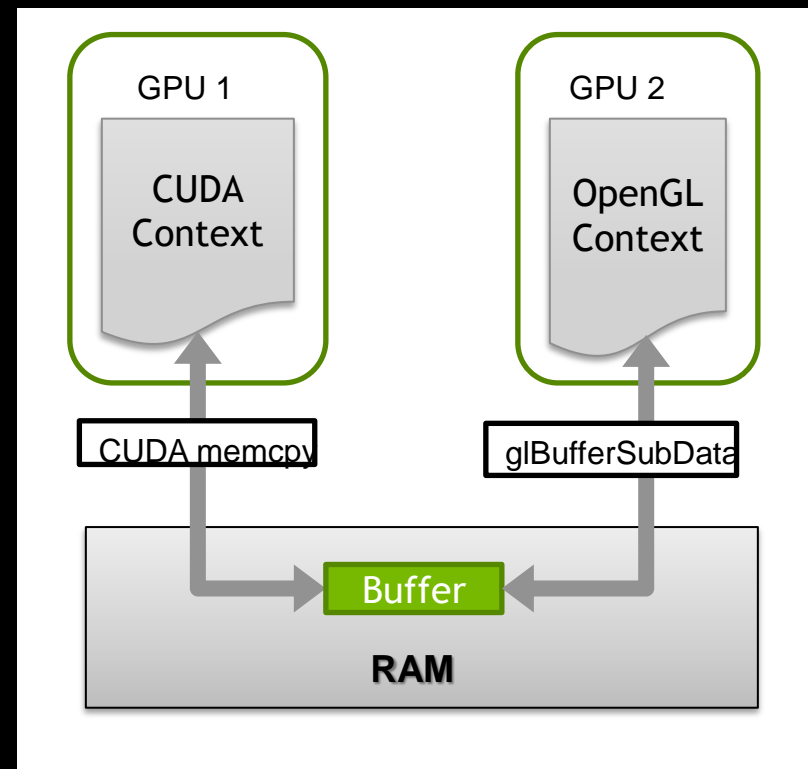# API Interop hides all the complexity

# If not cross-GPU API interop then what?
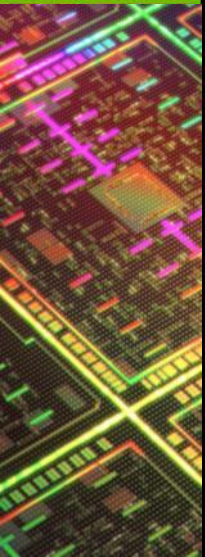
A)



B)

# Compute/Graphics interoperability: What's new?

- cudaD3D[9|10|11]SetDirect3DDevice/cudaGLSetGlDevice are no longer required

- All mipmap levels are shared

- Interoperate with OpenGL and DirectX at the same time

- Lots and lots of Windows WDDM improvements

# Conclusions/Resources

- The driver will do all the heavy lifting but..
- Scalability and final performance is up to the developer and..
- For fine grained control you might want to move data yourself.
- CUDA samples/documentation:

  http://developer.nvidia.com/cuda-downloads