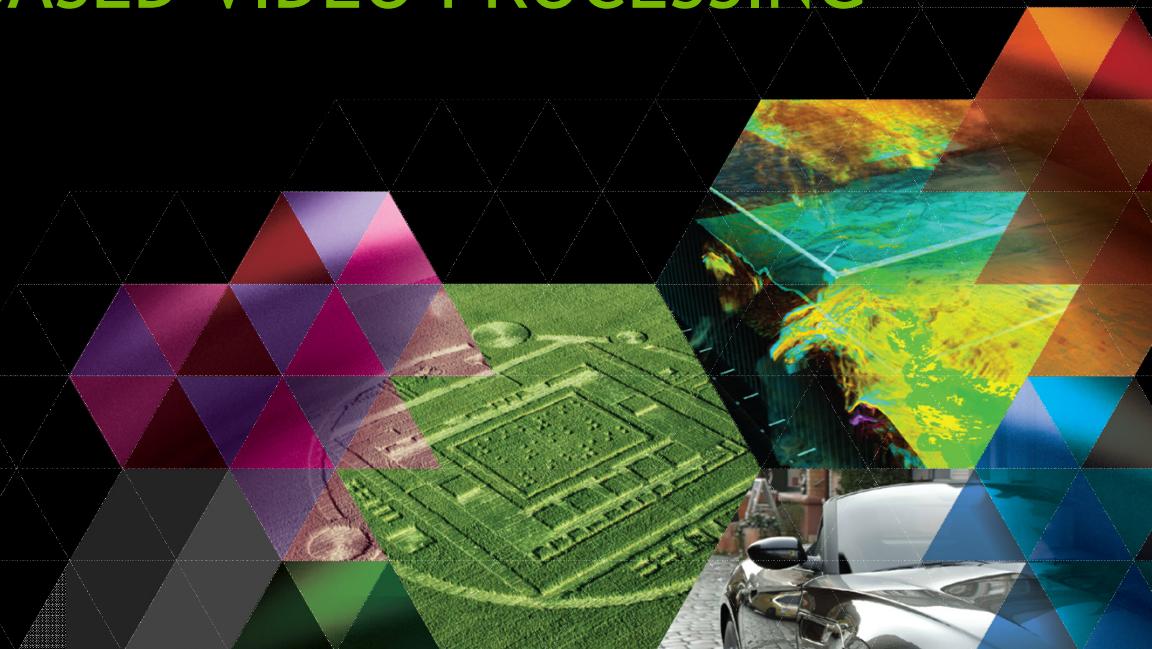


TOPICS IN GPU-BASED VIDEO PROCESSING

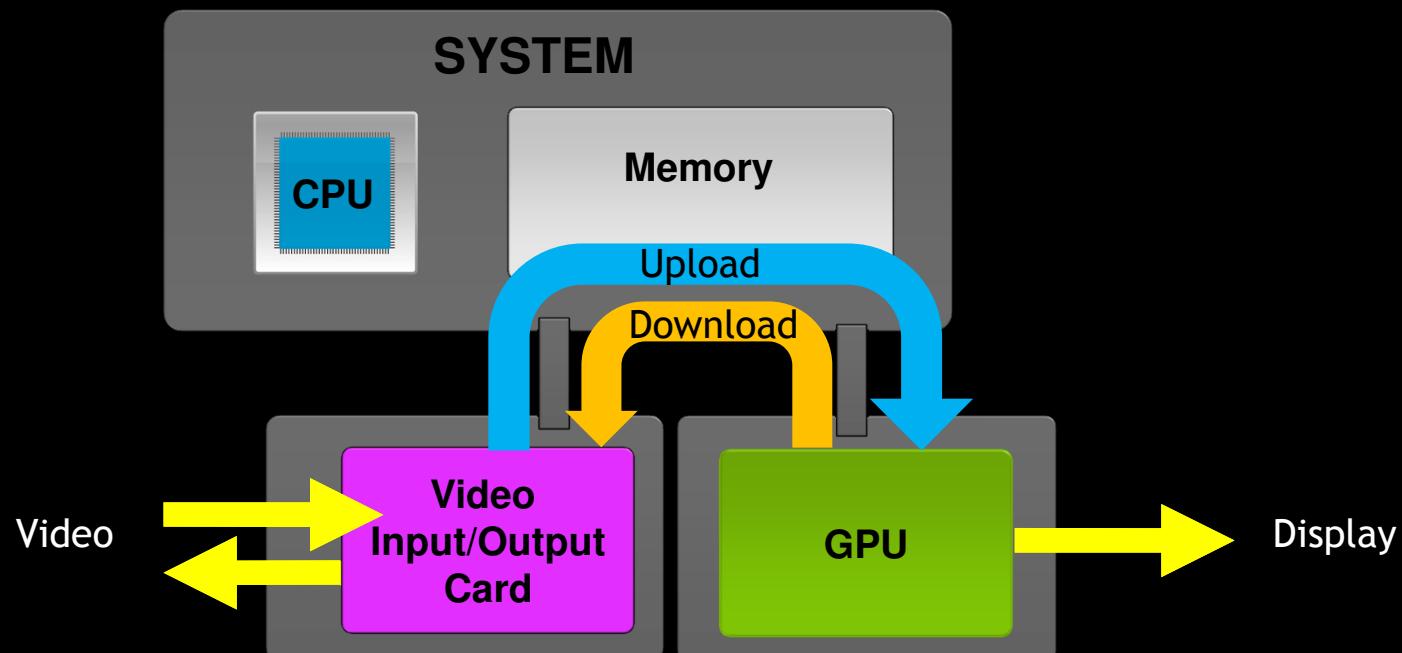
Thomas J. True
ttrue@nvidia.com



AGENDA

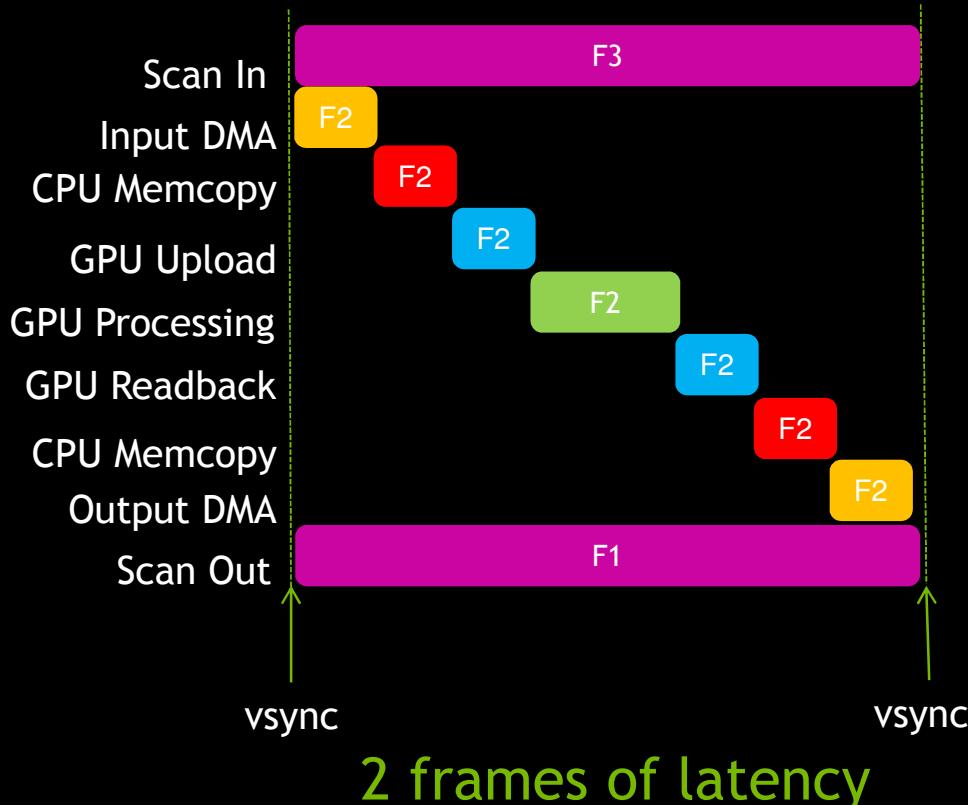
1. *GPU-Based Video Processing Pipeline*
2. *Optimized Data Transfers - Overlap Copy and Compute*
3. *Minimizing the Impact of Compute/Graphics Interop*
4. *OpenGL Compute*
5. *Compute Kernel Design and Optimization*
6. *High Bit Depth Considerations*
7. *High Resolution Display Considerations*

GPU-BASED VIDEO PROCESSING PIPELINE



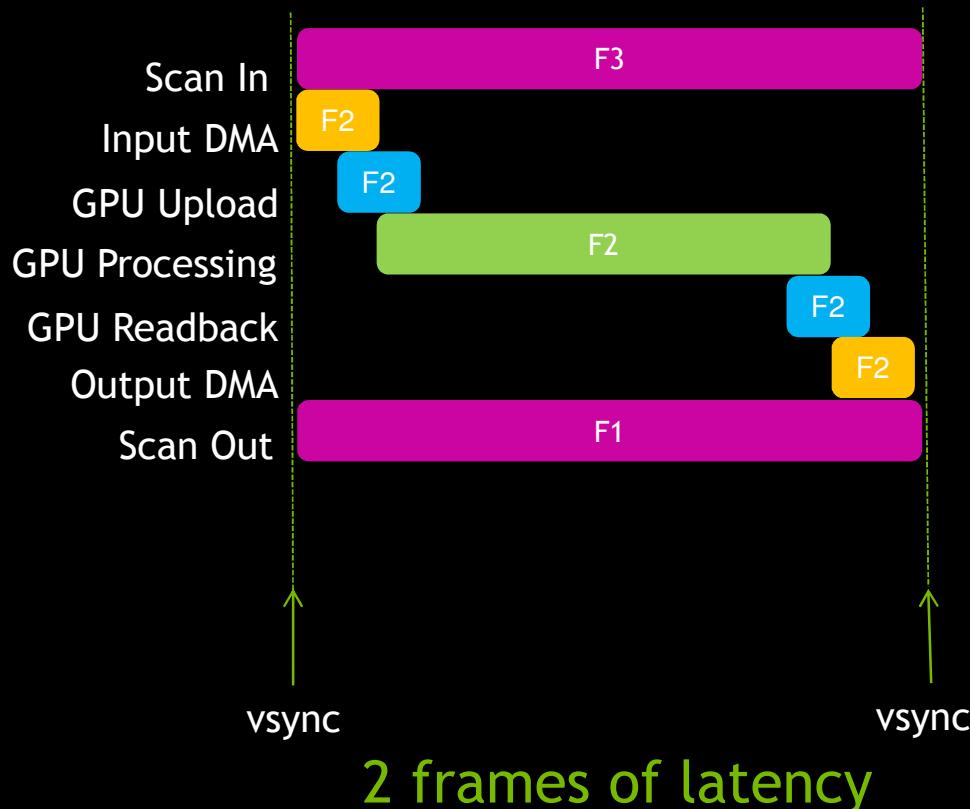
Video I/O With GPU Processing

INEFFICIENT GPU-BASED VIDEO PROCESSING



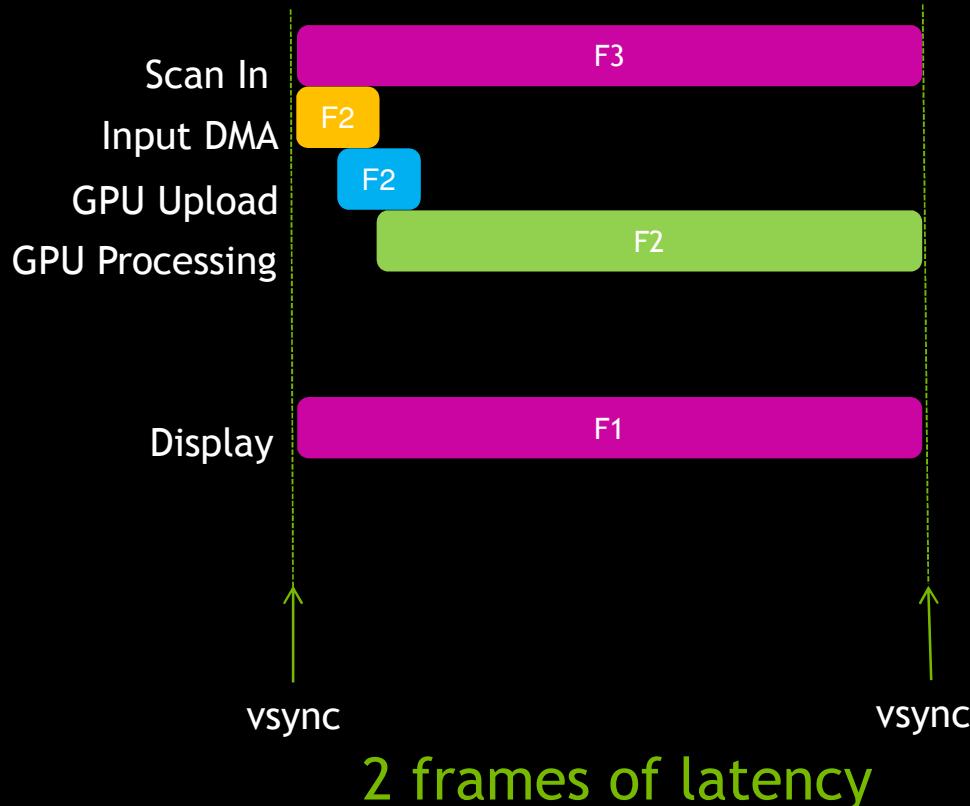
- I/O <-> GPU through GPU unshareable system memory buffer
- GPU synchronous transfers

OPTIMAL GPU-BASED VIDEO PROCESSING



- Pinned system memory
- Overlap I/O DMA transfers with GPU transfers and compute
- GPU Asynchronous
- Partial frame support

OPTIMAL GPU-BASED VIDEO PROCESSING



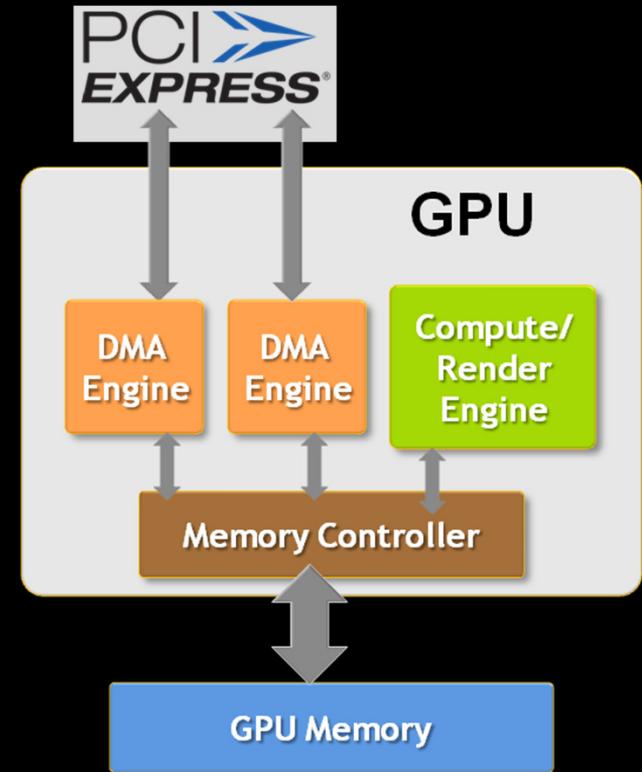
- Pinned system memory
- Overlap I/O DMA transfers with GPU transfers and compute
- GPU Asynchronous
- Partial frame support

OPTIMAL DATA MOVEMENT

- Minimize frame time consumed by GPU upload and download
- Overlap data transfers with compute / rendering

DUAL COPY ENGINES

- Allows copy-to-device + compute + copy-to-host to overlap simultaneously
- Compute/CUDA/OpenCL
 - Using async API's for memcpy with CUDA streams
 - Memory is allocated as page-locked
- Graphics/OpenGL
 - Using PBO's in multiple threads



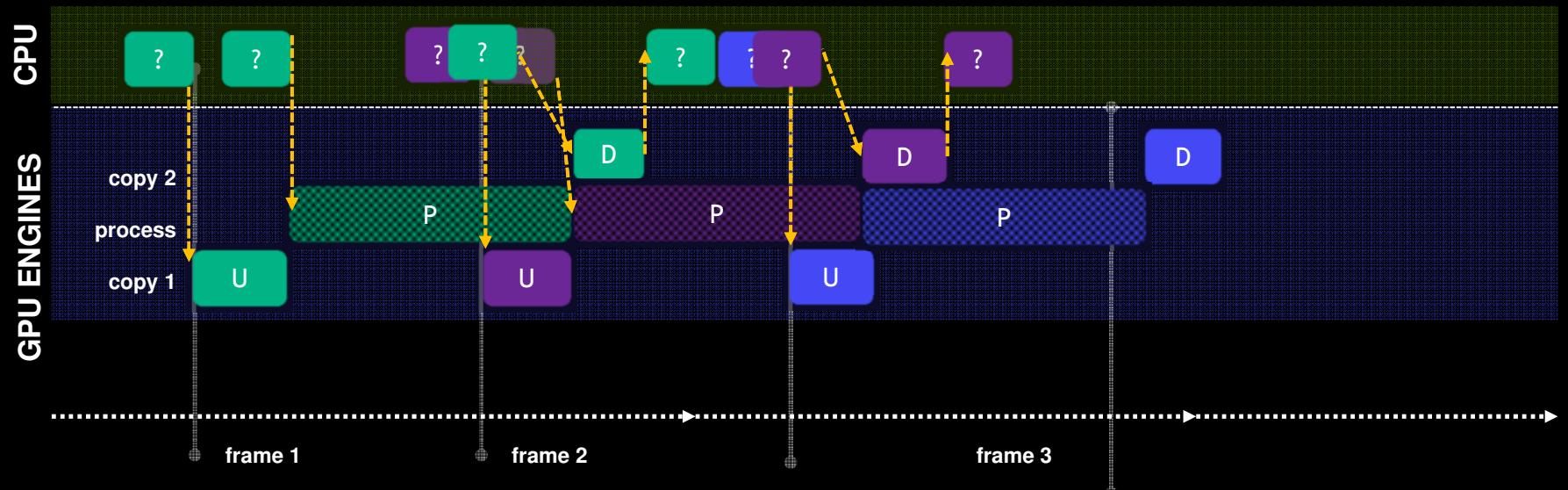
PINNED HOST MEMORY ALLOCATION IN CUDA

```
void *pmem  
cudaHostAlloc(pmем, size, cudaHostAllocDefault);
```

```
void *pmem  
pmем = malloc(size);  
  
cudaHostRegister(pmем, size, flags);
```

GOAL: OVERLAP I/O WITH PROCESSING

- How can we design the application so that GPU processing can last almost the entire frame time and we don't drop frames?



VIDEO CAPTURE, PROCESS AND PLAY - CUDA

```
while (!done) {
    // Start download of processed frame
    cudaMemcpyAsync(pCPUbuf [(bufNum+1)%3], pGPUbuf [bufNum],
                   size, cudaMemcpyDeviceToHost, downloadStream) ;

    // Start upload of captured frame
    cudaMemcpyAsync(pGPUbuf [bufNum+2], pCPUbuf [(bufNum+2)%3],
                   size, cudaMemcpyHostToDevice, uploadStream) ;

    // GPU-based video processing
    myKernel1<<<gridSz, BlockSz, 0, computeStream>>>(pGPUbuf [(bufNum+1)%3]...) ;
    myKernel2<<<gridSz, BlockSz, 0, computeStream>>>(pGPUbuf [(bufNum+1)%3]...) ;

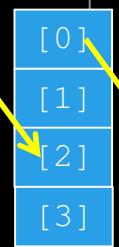
    // Execute other asynchronous GPU stuff

    // Capture next video frame
    GrabMyFrame(pCPUbuf [bufNum]) ;
    // Do other CPU stuff
    cudaThreadSynchronize();
    bufNum++; bufNum %=3;
}
```

OVERLAP CAPTURE, PROCESS AND PLAY - OPENGL

Capture Thread

```
// Wait for signal to start upload  
CPUWait(startUploadValid[2]);  
glWaitSync(startUpload[2]);  
  
// Bind texture object  
BindTexture(capTex[2]);  
  
// Upload  
glTexSubImage(texID...);  
  
// Signal upload complete  
GLOsync endUpload[2] = glFenceSync(...);  
Signal(endUploadValid[2]);
```



Render Thread

```
// Wait for download to complete  
CPUWait(endDownloadValid[3]);  
glWaitSync(endDownload[3]);  
  
// Wait for upload to complete  
CPUWait(endUploadValid[0]);  
glWaitSync(endUpload[0]);  
  
// Bind render target  
glFramebufferTexture(playTex[3]);  
  
// Bind video capture source texture  
BindTexture(capTex[0]);  
  
// Draw  
  
// Signal next upload  
startUpload[0] = glFenceSync(...);  
Signal(startUploadValid[0]);  
// Signal next download  
startDownload[3] = glFenceSync(...);  
Signal(startDownloadValid[3]);
```



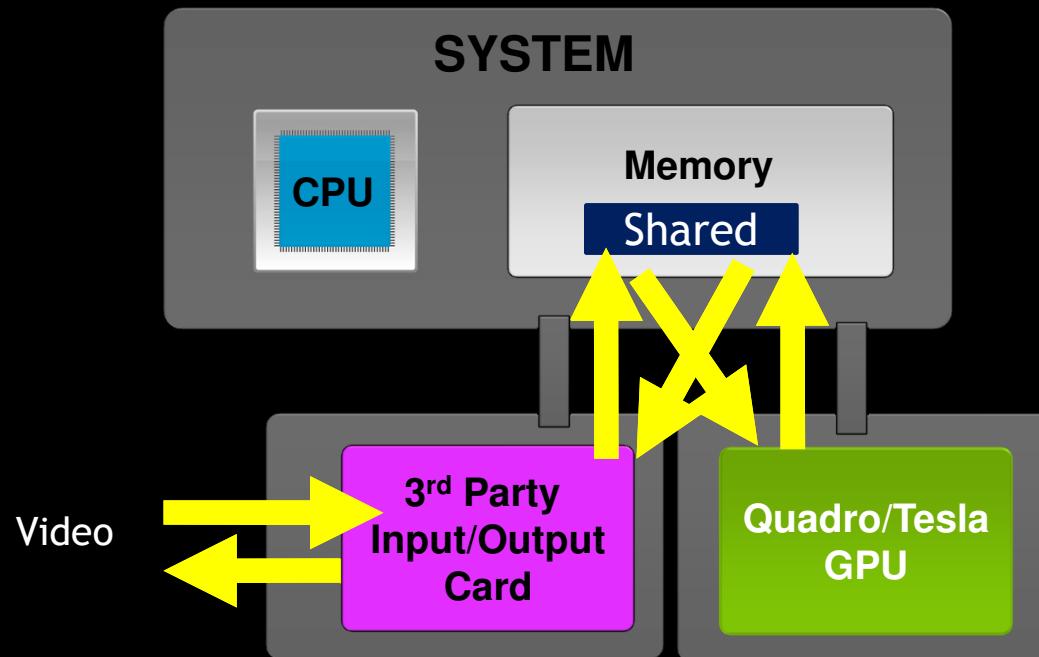
Playout Thread

```
// Playout thread  
CPUWait(startDownloadValid[2]);  
glWaitSync(startDownload[2]);  
  
// Readback  
glGetTexImage(playTex[2]);  
  
// Read pixels to PBO  
  
// Signal download complete  
endDownload[2] = glFenceSync(...);  
Signal(endDownloadValid[2]);
```



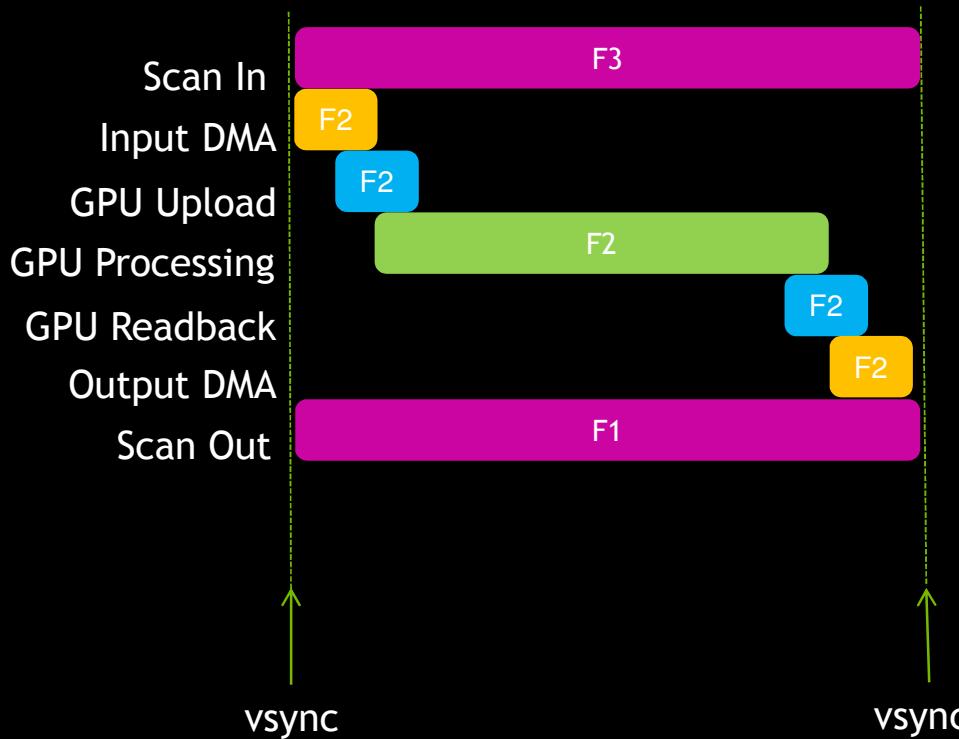
http://www.nvidia.com/docs/IO/40049/Dual_copy_engines.pdf

GPU DIRECT FOR VIDEO



Video Transfers through a Shareable System Memory Buffer!

OVERLAP COPY AND COMPUTE - GPUDIRECT FOR VIDEO



- Unified data transfer API for all Graphics and Compute APIs' objects
 - Video oriented
 - Efficient synchronization
 - GPU shareable system memory
 - Sub-field transfers
 - GPU Asynchronous transfers

Requires 3rd Party Video I/O SDK support GPUDirect for Video

GPU DIRECT FOR VIDEO APP STRUCTURE

- This structure can be used for any GPU API of your choice: CUDA, OpenGL, DirectX

Capture Thread

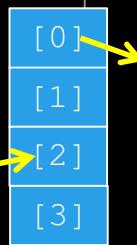
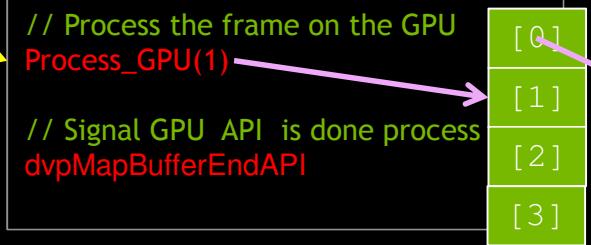
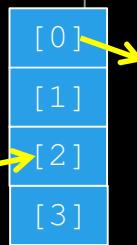
```
3rdPartyInitiateTransfer(2)  
// Wait for 3rd party to be done with frame  
dvpSyncObjClientWaitComplete(2)  
  
// Wait for GPU API to be done  
dvpMapBufferWaitDVP  
  
// transfer to GPU  
dvpMemcpy(2)  
  
// Signal GPU transfer is complete  
dvpMapBufferEndDVP
```

Render Thread

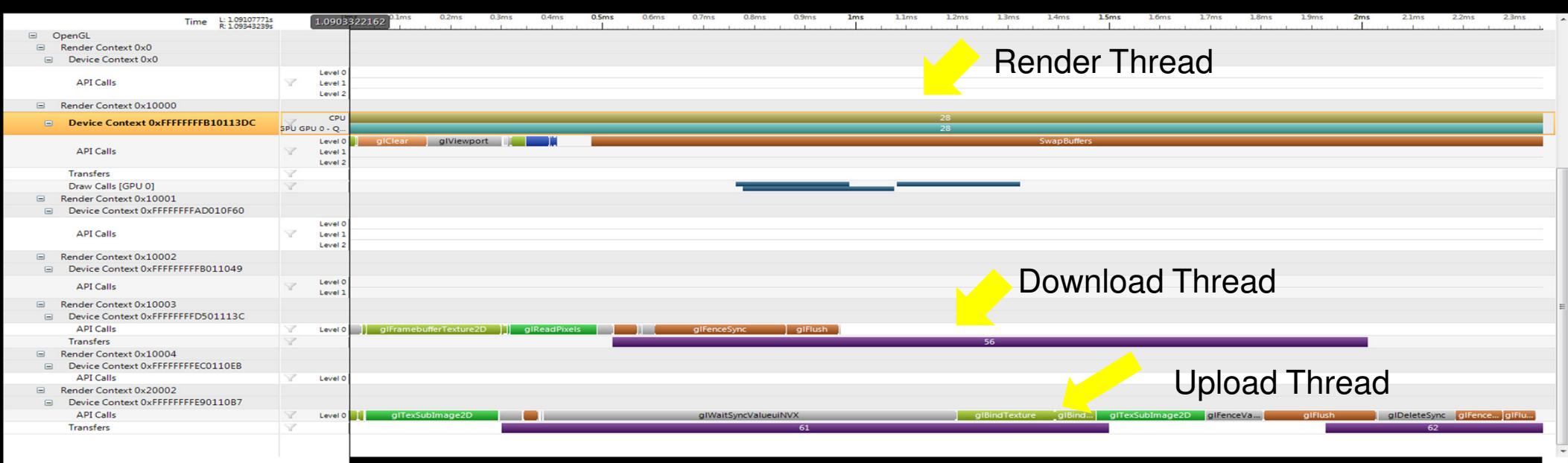
```
// Wait for GPU transfers to be done  
dvpMapBufferWaitAPI  
  
// Process the frame on the GPU  
Process_GPU(1)  
  
// Signal GPU API is done process  
dvpMapBufferEndAPI
```

Playout Thread

```
// Wait for GPU API to be done  
dvpMapBufferWaitDVP  
  
// Transfer from GPU  
dvpMemcpy(0)  
  
// Signal GPU transfer is complete  
dvpMapBufferEndDVP  
  
// Wait for signal to start transfer to 3rd party  
3rdPartyWaitForGPUComplete_CPU(0)  
3rdPartyInitiateTransfer(0)
```



GPU DIRECT FOR VIDEO



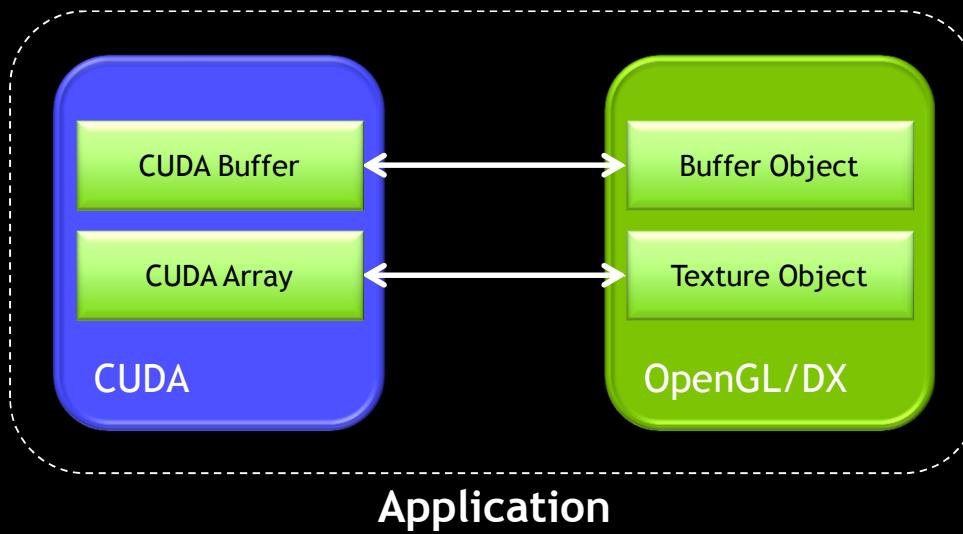
Parallel Nsight demonstrates the overlap of upload, render and download.

WHY NOT PEER-TO-PEER

- Same performance and latency for video size frames
- Supportability for ISVs, 3rd party I/O vendors and NVIDIA
 - The upgrade from Fermi to Kepler GPUs required NO code changes to 3rd party I/O SDKs
- IOH-to-IOH communication issues
 - Limited PCIE slot options due to lane allocations
- Support Graphics APIs as well as CUDA
- Multi-GPU Support with R313 display driver
- Multi-OS support, Windows and Linux

COMPUTE/GRAFICS INTEROP

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



CUDA GL INTEROP SAMPLE

- Setup and register OpenGL buffer objects

CUDA GL INTEROP SAMPLE

- Mapping between contexts:

```
cudaArray* texArray;  
  
while (!done)  
{  
    cudaGraphicsMapResources(1, &texRes);  
  
    cudaGraphicsSubResourceGetMappedArray(&texArray, texRes, 0, 0);  
  
    doCUDA(texArray);  
  
    cudaGraphicsUnmapResources(1, &texRes);  
  
    doGL(texId);  
}
```

CUDA GL INTEROP SAMPLE

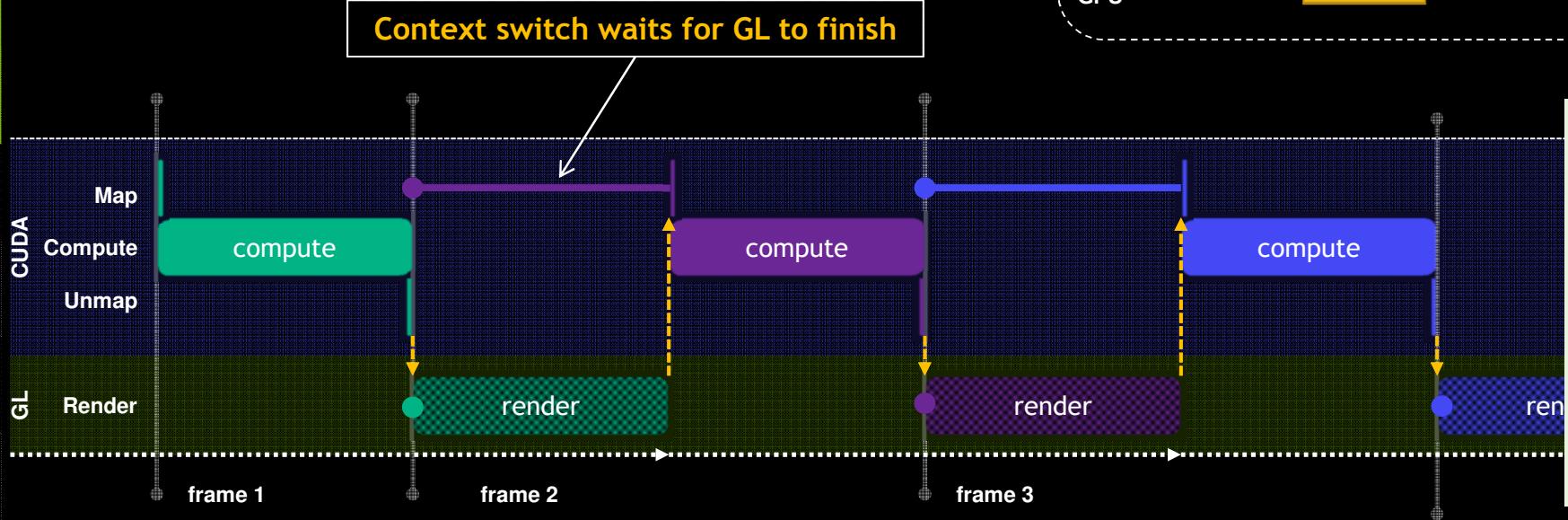
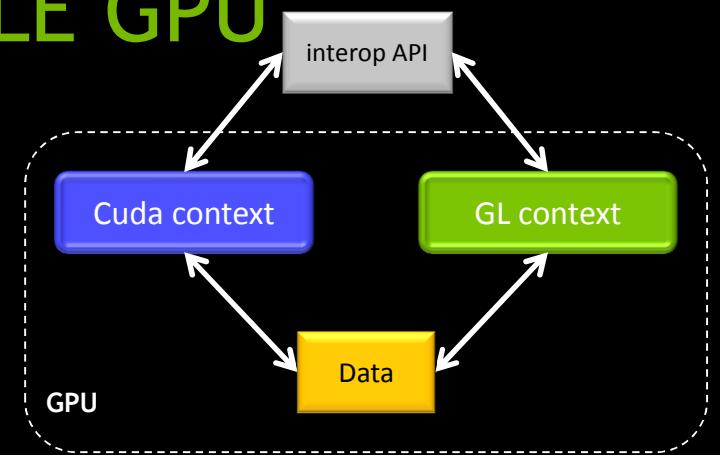
- Sync points:

```
cudaArray* texArray;  
  
while (!done)  
{  
    cudaGraphicsMapResources(1, &texRes);  
  
    cudaGraphicsSubResourceGetMappedArray(&texArray, texRes, 0, 0);  
    doCUDA(texArray);  
  
    cudaGraphicsUnmapResources(1, &texRes);  
  
    doGL(texId);  
}
```

Context-switching happens here.

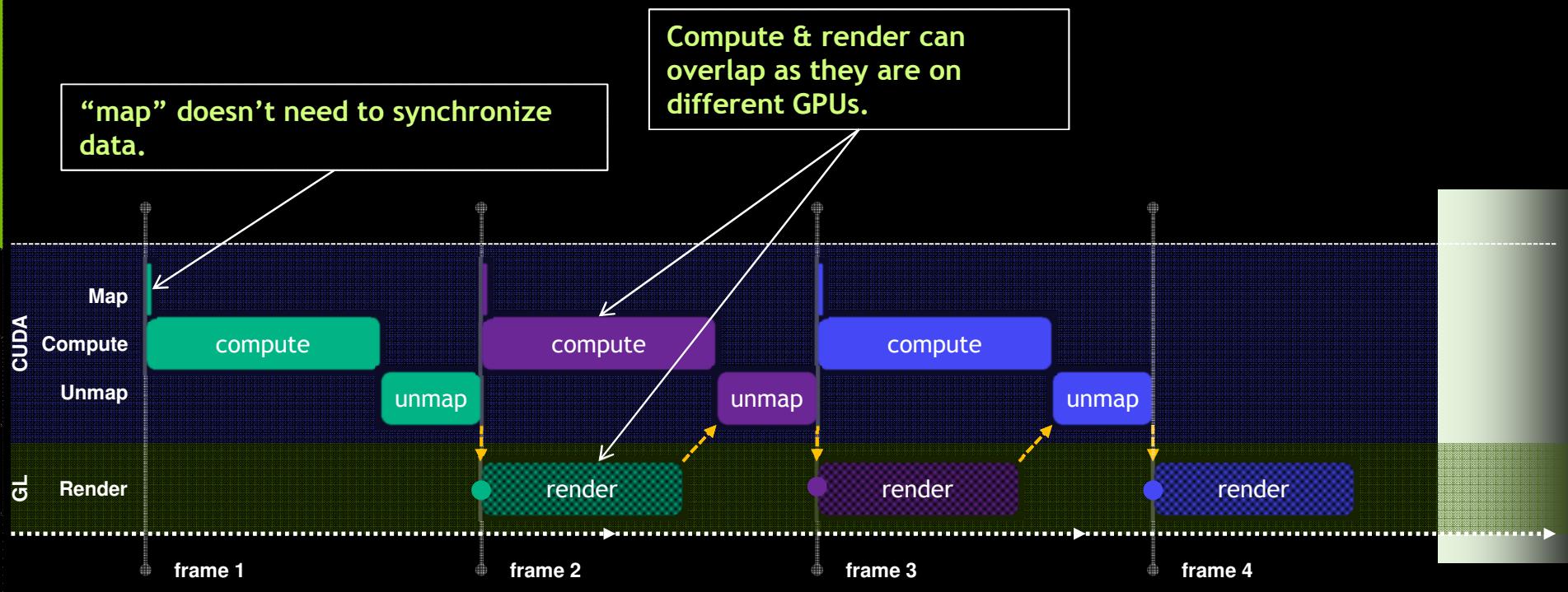
INTEROP BEHAVIOR: SINGLE GPU

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



INTEROP BEHAVIOR: MULTI-GPU

- CUDA-based Image Processing + OpenGL Display...



CODE SAMPLE - WRITE-DISCARD

```
cudaArray* texArray;

cudaGraphicsResourceSetMapFlags(texRes, cudaGraphicsMapFlagsWriteDiscard);

while (!done)
{
    cudaGraphicsMapResources(1, &texRes);

    cudaGraphicsSubResourceGetMappedArray(&texArray, texRes, 0, 0);
    doCUDA(texArray);

    cudaGraphicsUnmapResources(1, &texRes);

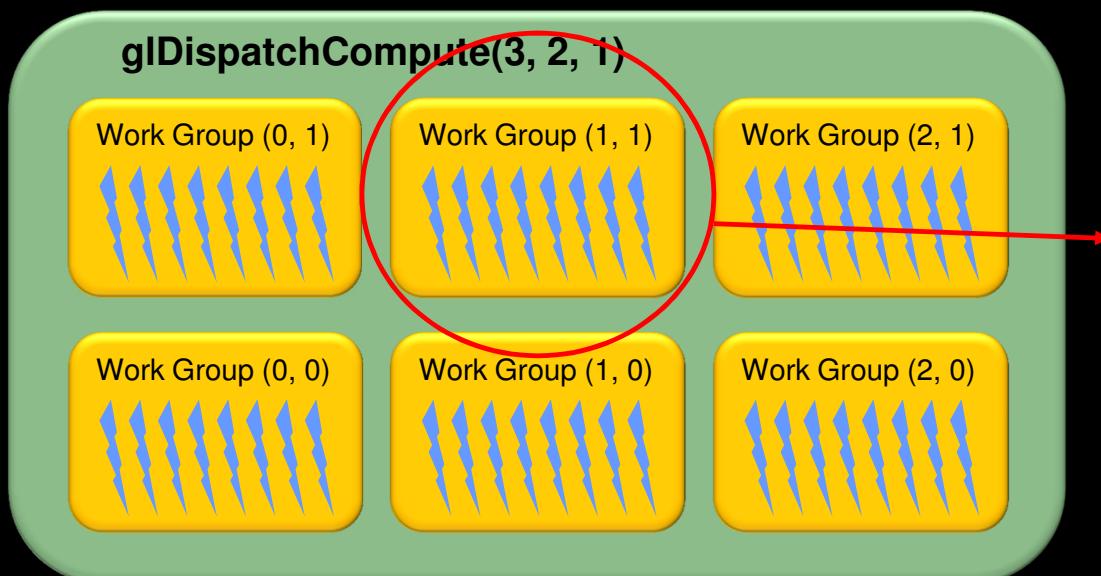
    doGL(texId);
}
```

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

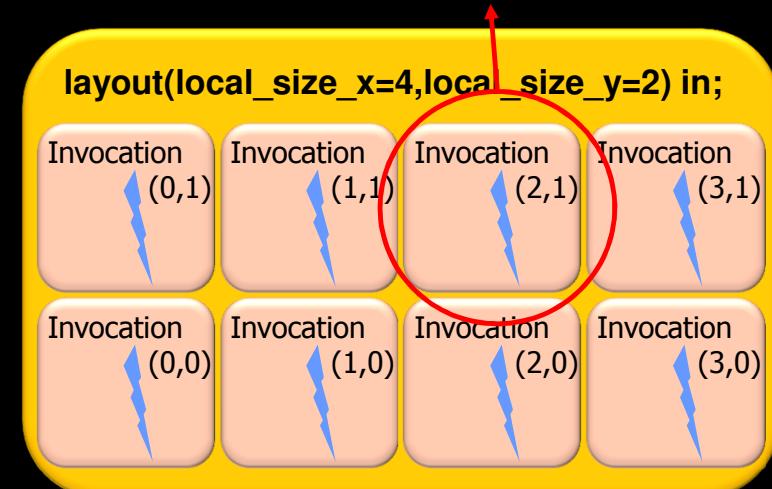
AVOID INTEROP WITH COMPUTE SHADERS

- OpenGL:
 - GL_ARB_compute_shader
 - Consistent GLSL API with familiar GLSL structures, textures, buffers, etc. (Requires OpenGL 4.3)
- DirectX:
 - DX11 DirectCompute
- No context-switching required
- No extra drivers or libraries needed.
- Less overhead on dispatching work
- No device setup code required

COMPUTE PROGRAMMING MODEL

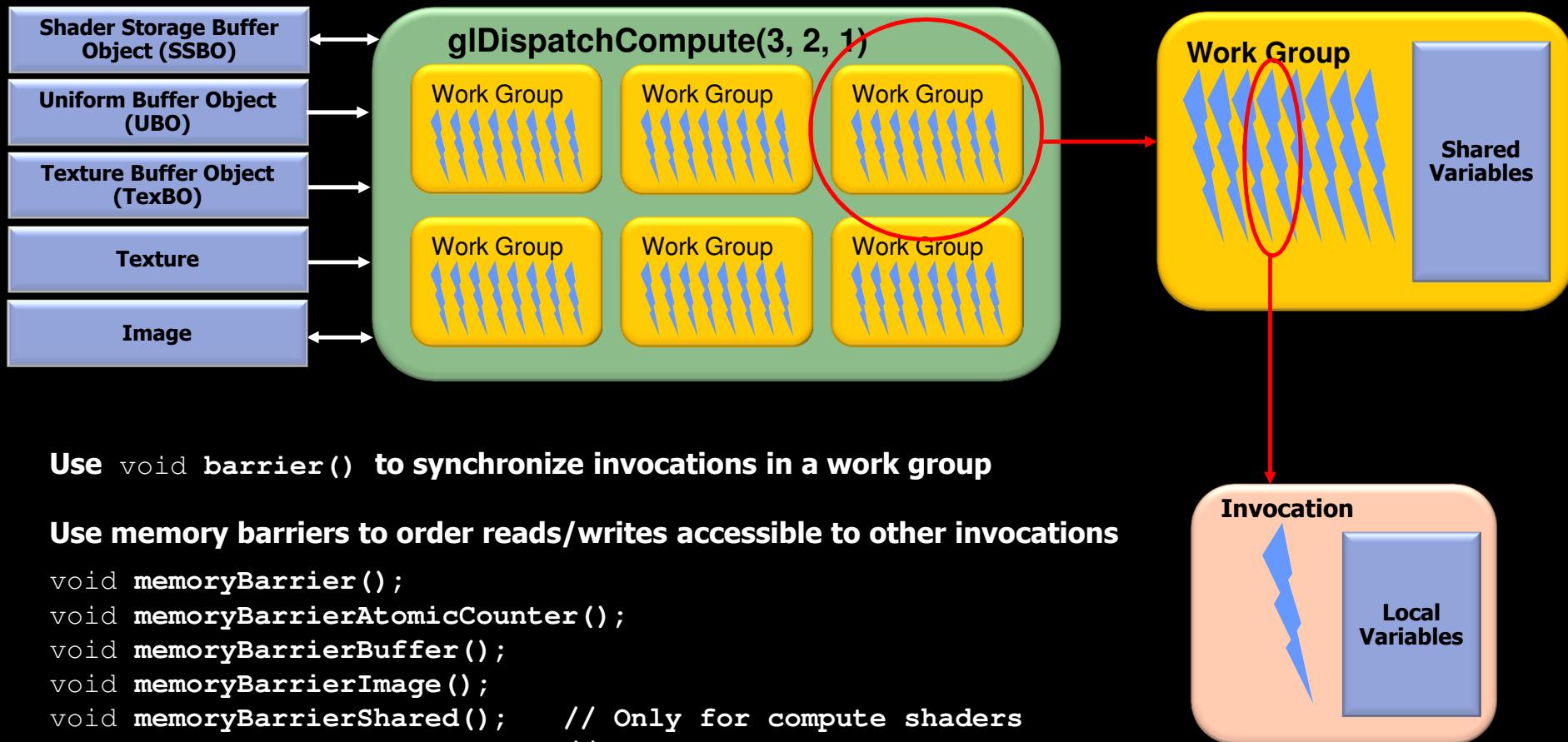


```
gl_WorkGroupSize = (4, 2, 0)
gl_WorkGroupID = (1, 1, 0)
gl_LocalInvocationID = (2, 1, 0)
gl_GlobalInvocationID = (6, 3, 0)
```



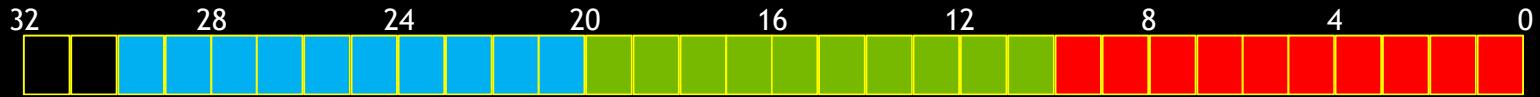
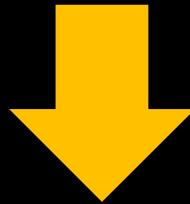
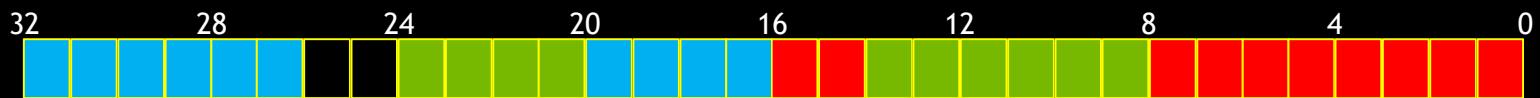
```
in uvec3 gl_NumWorkGroups; // Number of workgroups dispatched
const uvec3 gl_WorkGroupSize; // Size of each work group for current shader
in uvec3 gl_WorkGroupID; // Index of current work group being executed
in uvec3 gl_LocalInvocationID; // index of current invocation in a work group
in uvec3 gl_GlobalInvocationID; // Unique ID across all work groups
```

COMPUTE MEMORY HIERARCHY



GLCOMPUTE EXAMPLE

Problem: 10-Bit Big Endian to Little Endian Conversion



GLCOMPUTE EXAMPLE

```
#version 430

// Needed to use integer textures
#extension GL_EXT_gpu_shader4 : enable
#extension GL_ARB_compute_shader : enable

// Define local size of compute kernel
layout(local_size_x = 24,
       local_size_y = 24,
       local_size_z = 1) in;

// Input = single-component 32-bit unsigned integer texture
layout(r32ui, location = 0) uniform uimage2D inTex;

// Output = four component R10G10B10A2 unsigned integer texture
layout(rgb10_a2ui, location = 1) uniform uimage2D outTex;
```

GLCOMPUTE EXAMPLE

```
void main() {
    ivec2 outPos = ivec2(gl_GlobalInvocationID.xy);
    ivec2 inPos = ivec2(gl_GlobalInvocationID.xy);

    uint texel;
    uvec4 color2;

    // Components are in the upper 30 bits.
    texel = imageLoad(inTex, inPos);

    // Swizzle from Big Endian to Little Endian
    color2.r = ((texel & 0xFC000000)>>26) + ((texel & 0xF0000)>>10);
    color2.g = ((texel & 0x3F00)>>4) + ((texel & 0xF00000)>>20);
    color2.b = ((texel & 0xFF)<<2) + ((texel & 0xC000)>>14);
    color2.a = 0;\n" " \n"

    // Store 4-component RGB result as uint
    imageStore(outTex, outPos, color2);
};
```

GLCOMPUTE EXAMPLE

```
// Enable the compute shader
glUseProgram(computeProgram);

// Bind input image texture
glBindImageTexture(0, inGpuTexture, 0, GL_FALSE, 0, GL_READ_ONLY, GL_R32UI);

// Bind output image texture
glBindImageTexture(1, scratchGpuTexture, 0, GL_FALSE, 0, GL_WRITE_ONLY, GL_RGB10_A2UI);

// Launch the thread groups
glDispatchCompute(width/24, height/24, 1);
```

OPENGL FRAGMENT SHADER EQUIVALENT

```
// Needed to use integer textures:  
#extension GL_EXT_gpu_shader4 : enable  
uniform usampler2D tex;  
  
uint colorLookup(vec2 coord) {  
    //Get Texel  
    return(texture2D(tex, coord));  
}\\"/>  
  
void main(out vec4 color2) {  
    vec2 coord = vec2(gl_TexCoord[0].s, gl_TexCoord[0].t);  
  
    uint color1 = colorLookup(coord);  
    float scale = 1.0f / 1023.0f;  
    float scalea = 1.0f / 3.0f;\n"  
    color2.r = float(((color1 & 0xFF)<<2) + ((color1 & 0xC000)>>14)) * scale;  
    color2.g = float(((color1 & 0x3F00)>>4) + ((color1 & 0xF00000)>>20)) * scale;  
    color2.b = float(((color1 & 0xFC000000)>>26)+ ((color1 & 0xF0000)>>10)) * scale;  
    color2.a = 0.0f;  
}
```

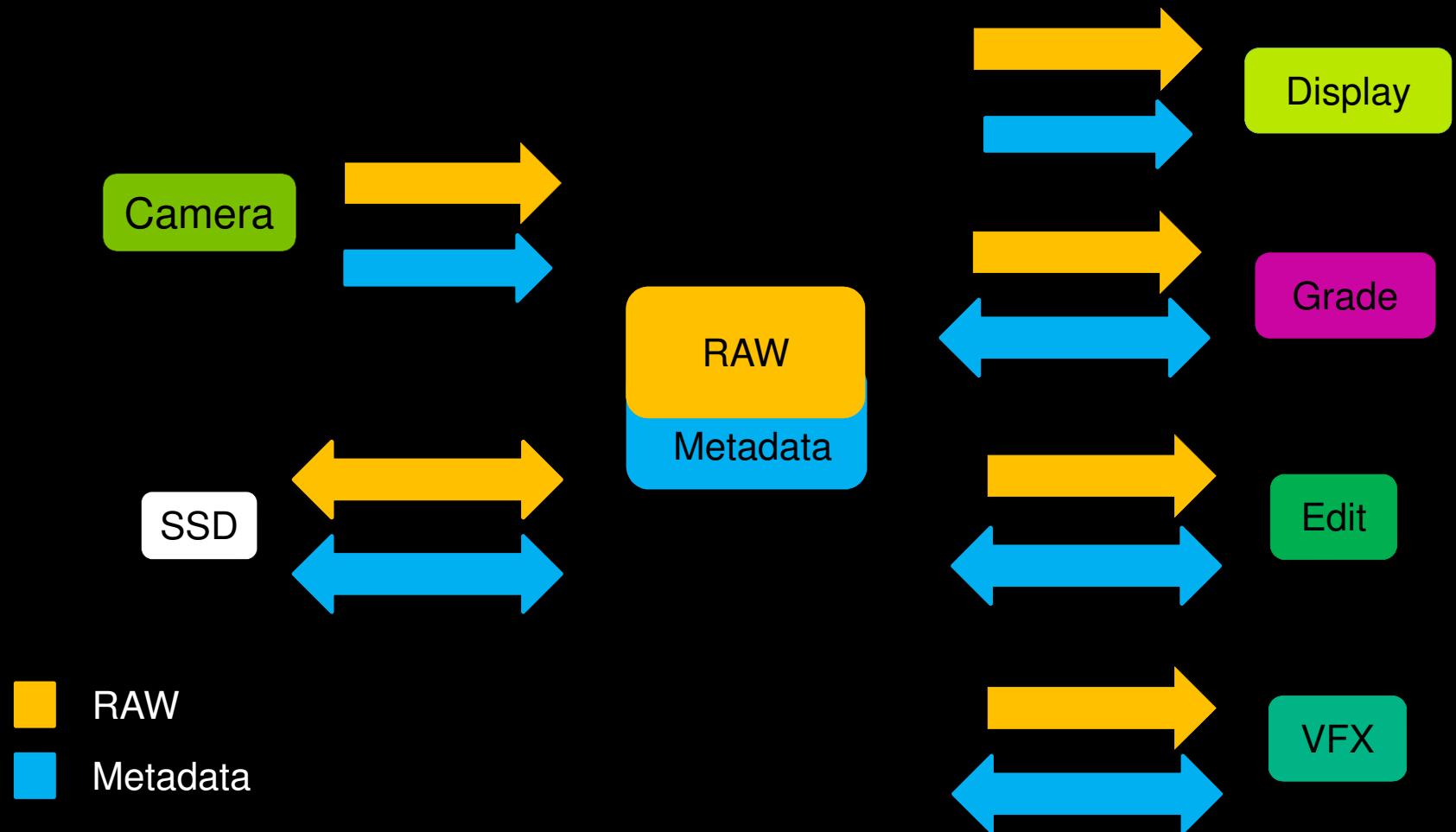
COMPUTE OR FRAGMENT SHADER?

- Fragment Shader
 - Granularity determined by the rasterizer
- Compute Shader
 - Granularity determined by the specified work group size

```
glDispatchCompute(GetWidth()/36, GetHeight(), 1);
```

- Provide more explicit control
- Each work group will execute on the same compute unit
- Shared memory

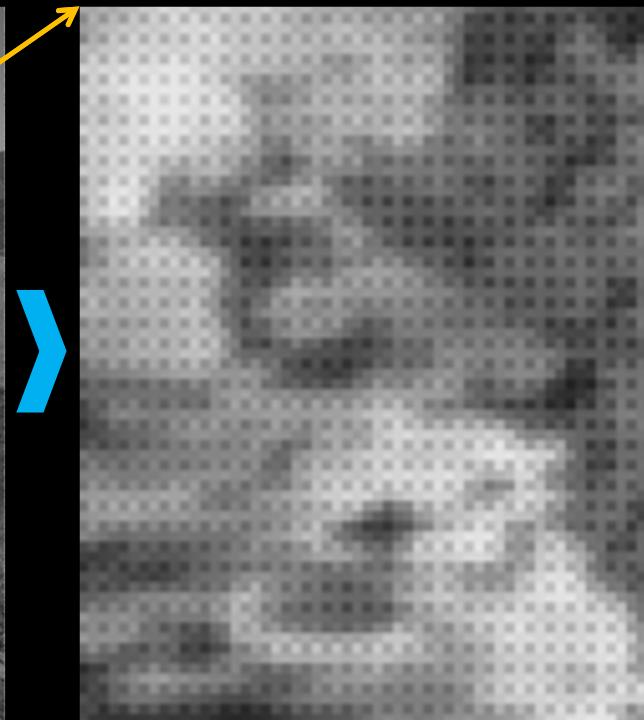
Example: 4K Debayering



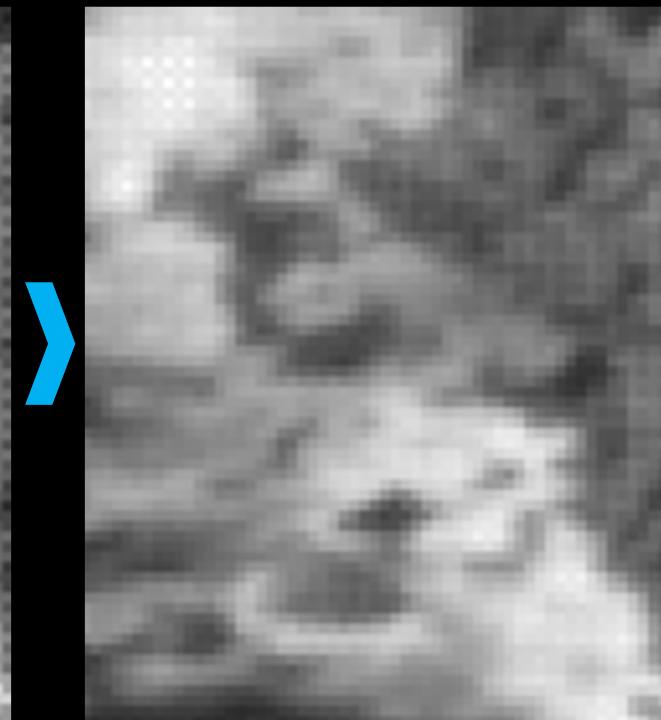
Final 4K Image Rendering



Original Image



Raw sensor data



**Linearization
Black Level Adjustment**

Images Courtesy of Jack Holm

Final 4K Image Rendering



De-Bayer and create
raw device RGB image



White Balance and
Color Space Conversion to
working color space



Tone Mapping
Gamma
Correction

Images Courtesy of Jack Holm

DEBAYER USING BILINEAR INTERPOLATION

R _{0,0}	G _{1,0}	R _{2,0}	G _{3,0}	R _{4,0}	G _{5,0}
G _{0,1}	B _{1,1}	G _{2,1}	B _{3,1}	G _{4,1}	B _{5,1}
R _{0,2}	G _{1,2}	R _{2,2}	G _{3,2}	R _{4,2}	G _{5,2}
G _{0,3}	B _{1,3}	G _{2,3}	B _{3,3}	G _{4,3}	B _{5,3}
R _{0,4}	G _{1,4}	R _{2,4}	G _{3,4}	R _{4,4}	G _{5,4}
G _{0,5}	B _{1,5}	G _{2,5}	B _{3,5}	G _{4,5}	B _{5,5}

Red Pixel

$$R_{2,2} = R_{2,2}$$

$$G_{2,2} = \frac{G_{2,1} + G_{1,2} + G_{3,2} + G_{2,3}}{4}$$

$$B_{2,2} = \frac{B_{1,1} + B_{3,2} + B_{1,3} + B_{3,3}}{4}$$

Blue Pixel

$$R_{1,1} = \frac{R_{0,0} + R_{2,0} + R_{0,2} + R_{2,2}}{4}$$

$$G_{1,1} = \frac{G_{1,0} + G_{0,1} + G_{2,1} + G_{1,2}}{4}$$

$$B_{1,1} = B_{1,1}$$

Green Pixel

$$R_{2,1} = \frac{R_{2,0} + R_{2,2}}{2}$$

$$G_{2,1} = G_{2,1}$$

$$B_{2,1} = \frac{B_{1,1} + B_{3,1}}{2}$$

DETERMINE UNIT OF WORK TO AVOID THREAD DIVERGENCE

Each thread computes color components for 2x2 quad:

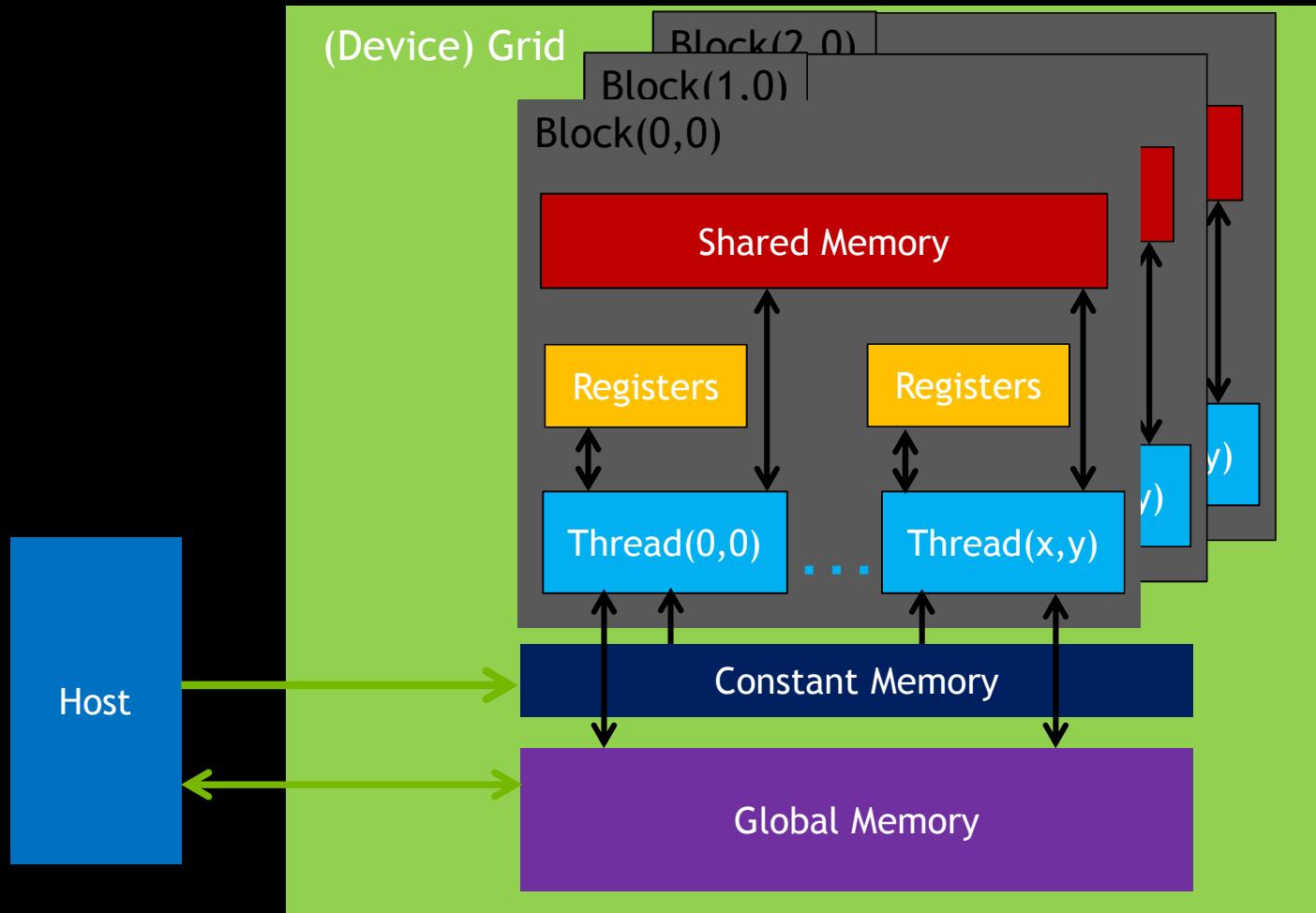
R _{0,0}	G _{1,0}	R _{2,0}	G _{3,0}	R _{4,0}	G _{5,0}
G _{0,1}	B _{1,1}	G _{2,1}	B _{3,1}	G _{4,1}	B _{5,1}
R _{0,2}	G _{1,2}	R _{2,2}	G _{3,2}	R _{4,2}	G _{5,2}
G _{0,3}	B _{1,3}	G _{2,3}	B _{3,3}	G _{4,3}	B _{5,3}
R _{0,4}	G _{1,4}	R _{2,4}	G _{3,4}	R _{4,4}	G _{5,4}
G _{0,5}	B _{1,5}	G _{2,5}	B _{3,5}	G _{4,5}	B _{5,5}

$$R_{2,2} = R_{2,2}$$

$$G_{2,2} = \frac{G_{2,1} + G_{1,2} + G_{3,2} + G_{2,3}}{4}$$

$$B_{2,2} = \frac{B_{1,1} + B_{3,2} + B_{1,3} + B_{3,3}}{4}$$

OPTIMIZE / COALESCE MEMORY ACCESS



TEXTURE MEMORY

- Separate dedicated texture cache
 - Only 1 device memory read on a miss
 - Out-of-bounds index handling (clamp or wrap)
 - Interpolation (linear, bilinear, trilinear or none)
 - Format conversion ({char, short, int} -> float)
- 
- Free

Memory Coalescing

Use memory access patterns that enable the GPU to coalesce groups of reads or writes of multiple data items into one operation.

TEXTURE MEMORY

```
// Allocate CUDA device memory for the source image on the GPU
unsigned short * pGPU_Source_Image;
size_t SourcePitch;
err = cudaMallocPitch((void**)&pGPU_Source_Image, &SourcePitch,
                     Image_Width*sizeof(unsigned short),
                     Image_Height);

// Copy the source image to the GPU
err = cudaMemcpy2DAsync(pGPU_Source_Image, SourcePitch,
                      pCPU_Source_Image,
                      Image_Width*sizeof(unsigned short),
                      Image_Width*sizeof(unsigned short),
                      Image_Height, cudaMemcpyHostToDevice, 1);
```

TEXTURE MEMORY

HSV Table	HSV mapping table
Exposure Curve	1D array specifying the exposure curve
Tone Map	1D array specifying the tone/look map
Gamma Table	1D array specifying the gamma ramp

CONSTANT MEMORY

- 64KB/device
 - Cached
 - As fast as reading from a register

CONSTANT MEMORY

Coefficient Table	Weight multipliers that represent the contribution of each of the 2 x 2 single color source pixels to each color channel component in the 2 x 2 RGB output pixels.
Offset Table	Precomputed offsets to each of the source pixels within the 2 x 2 source pixel quad.
White Point	Vector containing the white point
Camera CSC	3x3 matrix specifying the camera RGB colorspace conversion
Final CSC	3x3 matrix specifying the final output colorspace conversion

OPTIMIZE CGMA - COMPUTE TO GMEM ACCESS

$$R_{2,2} = R_{2,2} \quad G_{2,2} = \frac{G_{2,1} + G_{1,2} + G_{3,2} + G_{2,3}}{4} \quad B_{2,2} = \frac{B_{1,1} + B_{3,2} + B_{1,3} + B_{3,3}}{4}$$

No Shared Memory

3 values calculated for each 9 GMEM reads

.33:1

Quadro K5000: 173 GB/sec GMEM Bandwdith

$173 * 0.33 = 57 \text{ GB/sec} = 19 \text{ gigapixels}$

Shared Memory

Block Size: 32 x 21

SMEM Size: 34 x 22 (includes apron pixels)

672 values : 748 GMEM reads

.90:1

Quadro K5000: 173 GB/sec GMEM Bandwdith

$173 * 0.90 = 156 \text{ GB/sec} = 51 \text{ gigapixels}$

Use SMEM to improve CGMA!

USE SHARED MEMORY

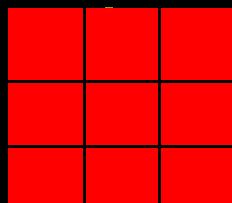
- Faster than global device memory.
- Structure usage to avoid bank conflicts

```
// Shared memory to hold the image tile we will read in from the raw image.  
__shared__ float tile_R[BILINEAR_TILE_H][BILINEAR_TILE_W];  
__shared__ float tile_G1[BILINEAR_TILE_H][BILINEAR_TILE_W];  
__shared__ float tile_B[BILINEAR_TILE_H][BILINEAR_TILE_W];  
__shared__ float tile_G2[BILINEAR_TILE_H][BILINEAR_TILE_W];
```

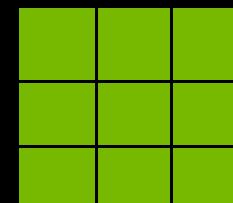
AVOID SMEM BANK CONFLICTS

Use four different color planes to store raw source pixels:

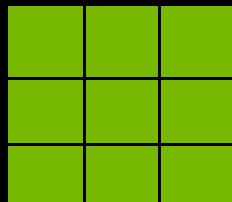
R _{0,0}	G _{1,0}	R _{2,0}	G _{3,0}	R _{4,0}	G _{5,0}
G _{0,1}	B _{1,1}	G _{2,1}	B _{3,1}	G _{4,1}	B _{5,1}
R _{0,2}	G _{1,2}	R _{2,2}	G _{3,2}	R _{4,2}	G _{5,2}
G _{0,3}	B _{1,3}	G _{2,3}	B _{3,3}	G _{4,3}	B _{5,3}
R _{0,4}	G _{1,4}	R _{2,4}	G _{3,4}	R _{4,4}	G _{5,4}
G _{0,5}	B _{1,5}	G _{2,5}	B _{3,5}	G _{4,5}	B _{5,5}



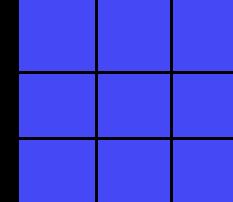
Red SMEM Tile



Green 1 SMEM Tile



Green 2 SMEM Tile



Blue SMEM Tile

MAXIMIZE OCCUPANCY

$$\text{occupancy} = \frac{\# \text{ of executing threads}}{\# \text{ of possible executing threads}}$$

- Function of:
 - Block size
 - Shared memory
 - Registers

MAXIMIZE OCCUPANCY

CUDA GPU Occupancy Calculator (version 1) [Autosaved] [Compatibility Mode]

Just follow steps 1, 2, and 3 below! (or click here for help)

1.) Select Compute Capability (click): 3.0 (Help)
 2.) Select Shared Memory Size Config (bytes): 49152

2.) Enter your resource usage:
 Threads Per Block: 672 (Help)
 Registers Per Thread: 11
 Shared Memory Per Block (bytes): 11968

(Don't edit anything below this line)

3.) GPU Occupancy Data is displayed here and in the graphs:
 Active Threads per Multiprocessor: 2016 (Help)
 Active Warps per Multiprocessor: 63
 Active Thread Blocks per Multiprocessor: 3
 Occupancy of each Multiprocessor: 98%

Physical Limits for GPU Compute Capability: 3.0
 Threads per Warp: 32
 Warps per Multiprocessor: 64
 Threads per Multiprocessor: 2048
 Thread Blocks per Multiprocessor: 16
 Total # of 32-bit registers per Multiprocessor: 65536
 Register allocation unit size: 256
 Register allocation granularity: warp
 Registers per Thread: 63
 Shared Memory per Multiprocessor (bytes): 49152
 Shared Memory Allocation unit size: 256
 Warp allocation granularity: 4
 Maximum Thread Block Size: 1024

Allocated Resources = Allocatable
 Warps (Threads Per Block / Threads Per Warp): 21 (Warp limit per SM due to per-warp reg count): 64 (Blocks/SM * Warps/Block = Warps/SM): 3
 Registers (Warp limit per SM due to per-warp reg count): 21 (Registers Per Thread): 128 (Blocks/SM * Registers/Block = Registers/SM): 6
 Shared Memory (Bytes): 12032 (Shared Memory Per Block (bytes)): 49152 (Blocks/SM * Shared Memory/Block = Shared Memory/SM): 4

Note: SM is an abbreviation for (Streaming) Multiprocessor

Maximum Thread Blocks Per Multiprocessor: Blocks/SM * Warps/Block = Warps/SM
 Limited by Max Warps or Max Blocks per Multiprocessor: 3 21 63
 Limited by Registers per Multiprocessor: 6
 Limited by Shared Memory per Multiprocessor: 4

Note: Occupancy limiter is shown in orange
 Physical Max Warps/SM = 64
 Occupancy = 63 / 64 = 98%

Click Here for detailed instructions on how to use this occupancy calculator.
 For more information on NVIDIA CUDA, visit <http://developer.nvidia.com/cuda>

Your chosen resource usage is indicated by the red triangle on the graphs. The other data points represent the range of possible block sizes, register counts, and shared memory allocation.

Impact of Varying Block Size

Threads Per Block	Multiprocessor Warp Occupancy (# warps)
64	4
128	8
192	12
256	16
320	20
384	24
448	28
512	32
576	36
640	40
704	44
768	48
832	52
896	56
960	60
1024	64

Impact of Varying Register Count Per Thread

Registers Per Thread	Multiprocessor Warp Occupancy (# warps)
8	64
16	40
24	16
32	0
40	0
48	0
56	0
64	0
72	0
80	0
88	0
96	0
104	0
112	0
120	0
128	0
136	0
144	0
152	0
160	0
168	0
176	0
184	0
192	0
200	0
208	0
216	0
224	0
232	0
240	0
248	0
256	0

KNOW YOUR SYSTEM: DEVICEQUERY.EXE

```
Device 0: "Quadro K5100M"
  CUDA Driver Version / Runtime Version      6.0 / 5.5
  CUDA Capability Major/Minor version number: 3.0
  Total amount of global memory:             8192 MBytes (8589934592 bytes)
  ( 8) Multiprocessors, (192) CUDA Cores/MP:
  GPU Clock rate:                          771 MHz (0.77 GHz)
  Memory Clock rate:                      1800 Mhz
  Memory Bus Width:                       256-bit
  L2 Cache Size:                           524288 bytes
  Maximum Texture Dimension Size (x,y,z): 1D=(65536), 2D=(65536, 65536),
                                             3D=(4096, 4096, 4096)
  Maximum Layered 1D Texture Size, (num) layers 1D=(16384), 2048 layers
  Maximum Layered 2D Texture Size, (num) layers 2D=(16384, 16384), 2048 layers
  Total amount of constant memory:           65536 bytes
  Total amount of shared memory per block:   49152 bytes
  Total number of registers available per block: 65536
  Warp size:                                32
  Maximum number of threads per multiprocessor: 2048
  Maximum number of threads per block:        1024
  Max dimension size of a thread block (x,y,z): (1024, 1024, 64)
  Max dimension size of a grid size (x,y,z):  (2147483647, 65535, 65535)
  Maximum memory pitch:                     2147483647 bytes
  Texture alignment:                        512 bytes
```

C:\NVIDIA\CUDA\CUDA Samples\Bin\win64\Release

KNOW YOUR SYSTEM: DEVICEQUERY.EXE

Concurrent copy and kernel execution:	Yes with 2 copy engine(s)
Run time limit on kernels:	Yes
Integrated GPU sharing Host Memory:	No
Support host page-locked memory mapping:	Yes
Alignment requirement for Surfaces:	Yes
Device has ECC support:	Disabled
CUDA Device Driver Mode (TCC or WDDM):	WDDM (Windows Display Driver Model)
Device supports Unified Addressing (UVA):	Yes
Device PCI Bus ID / PCI location ID:	1 / 0
Compute Mode: < Default (multiple host threads can use ::cudaSetDevice() with device simultaneously) >	
deviceQuery, CUDA Driver = CUDART, CUDA Driver Version = 6.0, CUDA Runtime Version = 5.5, NumDevs = 1, Device0 = Quadro K5100M	
Result = PASS	

C:\NVIDIA\CUDA\CUDA Samples\Bin\win64\Release

KNOW YOUR SYSTEM: BANDWIDTHTEST.EXE

Device 0: Quadro K5000
Quick Mode

Host to Device Bandwidth, 1 Device(s)

PINNED Memory Transfers

Transfer Size (Bytes)	Bandwidth (MB/s)
33554432	5937.4

Device to Host Bandwidth, 1 Device(s)

PINNED Memory Transfers

Transfer Size (Bytes)	Bandwidth (MB/s)
33554432	3409.9

Device to Device Bandwidth, 1 Device(s)

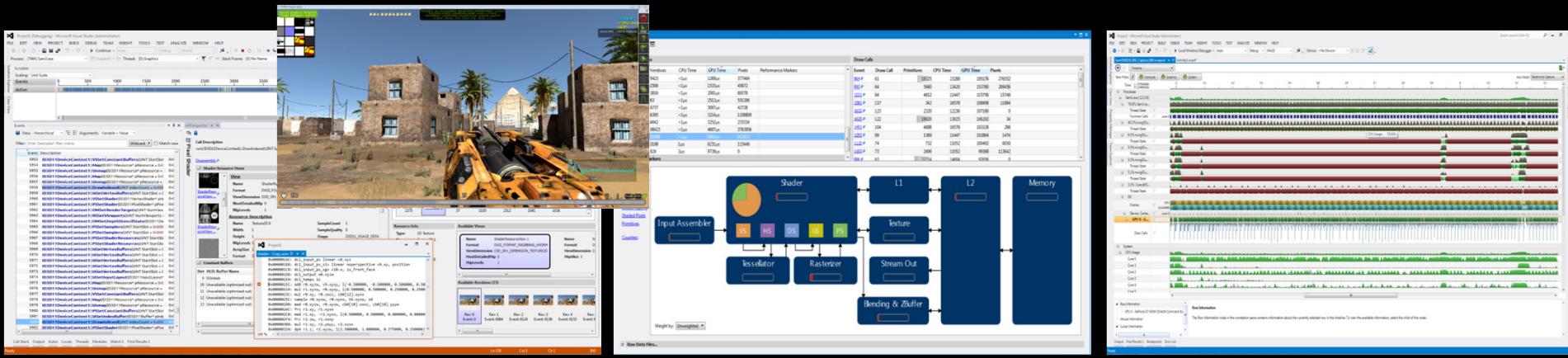
PINNED Memory Transfers

Transfer Size (Bytes)	Bandwidth (MB/s)
33554432	129900.5

C:\ProgramData\NVIDIA Corporation\CUDA Samples\v5.0\bin\

NVIDIA NSIGHT FOR GRAPHICS DEVELOPERS

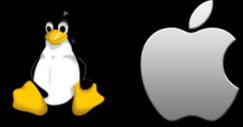
for OpenGL ~4.3+, Direct3D 9/11/11.1



- Frame debugger and profiler
- HLSL and GLSL shader debugger
- Application and system trace

NVIDIA® NSIGHT™ eclipse

Full featured IDE for developing CUDA-C apps on X86, ARM & POWER8.



```

// CUDA-aware Editor screenshot showing semantic highlighting of CUDA code.
// The code is annotated with various colors and icons to indicate different types of code elements.
// The interface includes a Project Explorer, a code editor with syntax highlighting, and a status bar at the bottom.

```

CUDA-Aware Editor

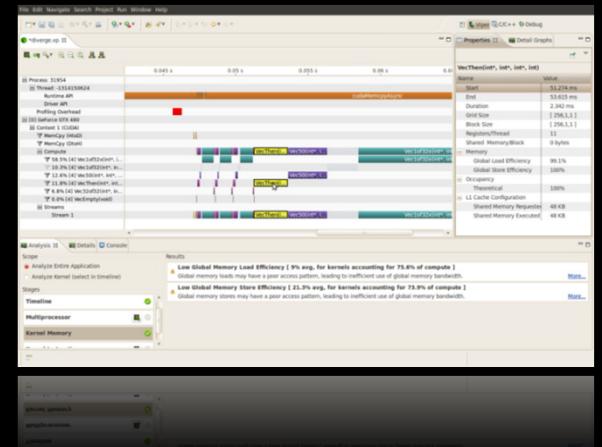
- Automated CPU to GPU code refactoring
- Semantic highlighting of CUDA code
- Integrated code samples & docs
- Cross-compilation for Linux target

The screenshot shows the Nsight Debugger interface with multiple panes:

- Debug View:** Shows a tree view of threads and their execution status (Running, Suspended, etc.).
- Variables View:** Displays variables across threads.
- CUDA Information View:** Provides details about CUDA operations like `cudaFindMax`.
- Registers View:** Shows register values for selected threads.
- Call Stack View:** Displays the call stack for the current thread.
- Output View:** Shows host-side logs and kernel output.
- Console View:** Displays CUDA error messages.

Nsight Debugger

- Simultaneously debug CPU and GPU code
- Inspect variables across CUDA threads
- Use breakpoints & single-step debugging
- Integrated CUDA memory checker



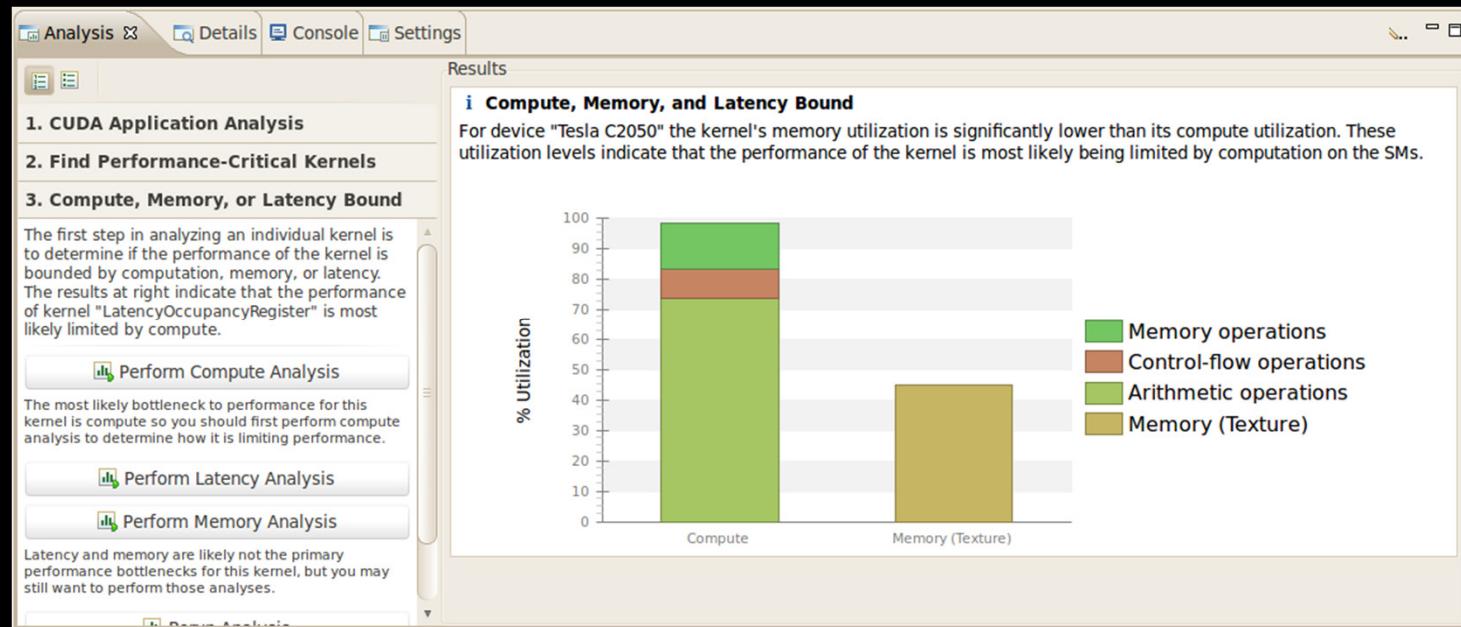
Nsight Profiler

- Quickly identifies performance issues
- Guided expert analysis
- Source line correlation

CUDA STANDALONE TOOLS

Visual Profiler

- Trace CUDA activities
- Performance instrumentation with source code correlation
- Guided Expert Analysis



NVPROF

- Generates execution summary
- Gather Performance events



CUDA-MEMCHECK

- Out of bounds memory access detection
- Detects Race Condition



* Android new in CUDA 6.0

CUDA-GDB

- Command line CUDA debugging
- Debug CPU and GPU code



HIGHER BIT DEPTH CONSIDERATIONS

Bit Depth	API	Scanout Surface	HDMI	DP	Full Screen	Window
10	DirectX	R10G10B10A2_UNORM R10G10B10_XR_BIAS_A2_UNORM	✓	✓	✓	X
10	OpenGL	A2RGB10	✓	✓	✓	✓ +
12	DirectX	R16G16B16A16_FLOAT	✓	✓ *	✓	X
12	OpenGL	X	X	X	X	X

*12bpc supported but not currently enabled due to unavailability of 12bpc DP devices on which to test it.

+Quadro only.

HIGHER RESOLUTION CONSIDERATIONS

- 4K, UltraHD, Cinema 4K -- many names for a display $\geq 3,840 \times 2,160$ pixels
 - Ultra HD: $3,840 \times 2,160$
 - Quad 1080p
 - Cinema 4K: $4,096 \times 2,160$
 - closer to std 35mm film aspect
- Need to pay attention to the frame-rate!
- Different panels/projectors have different connectors



HIGHER RESOLUTION CONSIDERATIONS

Number of Display Device Inputs:

Single Input

HDMI 1.4 (DVI-DL or DP1.1)

3840x2160@30, 4096x2160@24

DisplayPort 1.2

3,840x2,160@60

4,096x2,160@60 (K6000 only)

Multiple Inputs

2x DP1.1 or DVI-DL

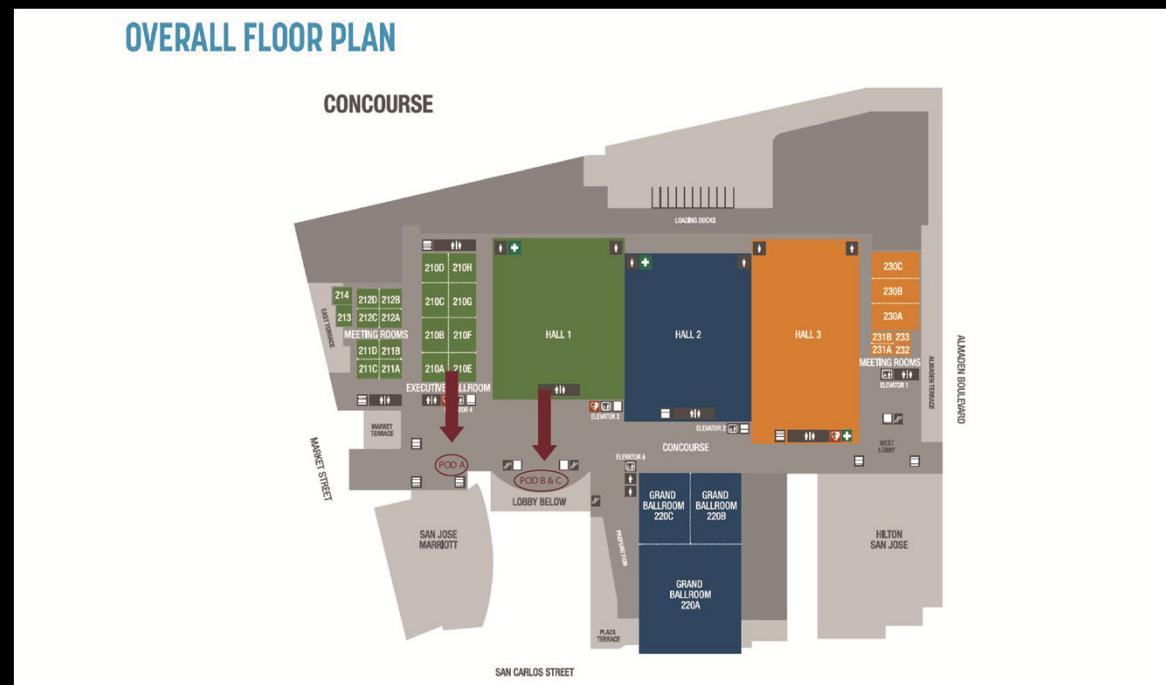
4,096x2,400@60

4x HDMI (3,840x2,400@60)

4x DP1.1 or DVI-DL (3,840x2,400@120)

GOT QUESTIONS? COME HANGOUT WITH US

GPU-Based Image and Video Processing Hangout
GTC Hangouts - Concourse Pod C
Wednesday 9-11AM



GPU TECHNOLOGY CONFERENCE



QUESTIONS?

