# CUDA FOR GRAPHICS

Advanced Computer Graphics    29.2.2012

Ville Timonen
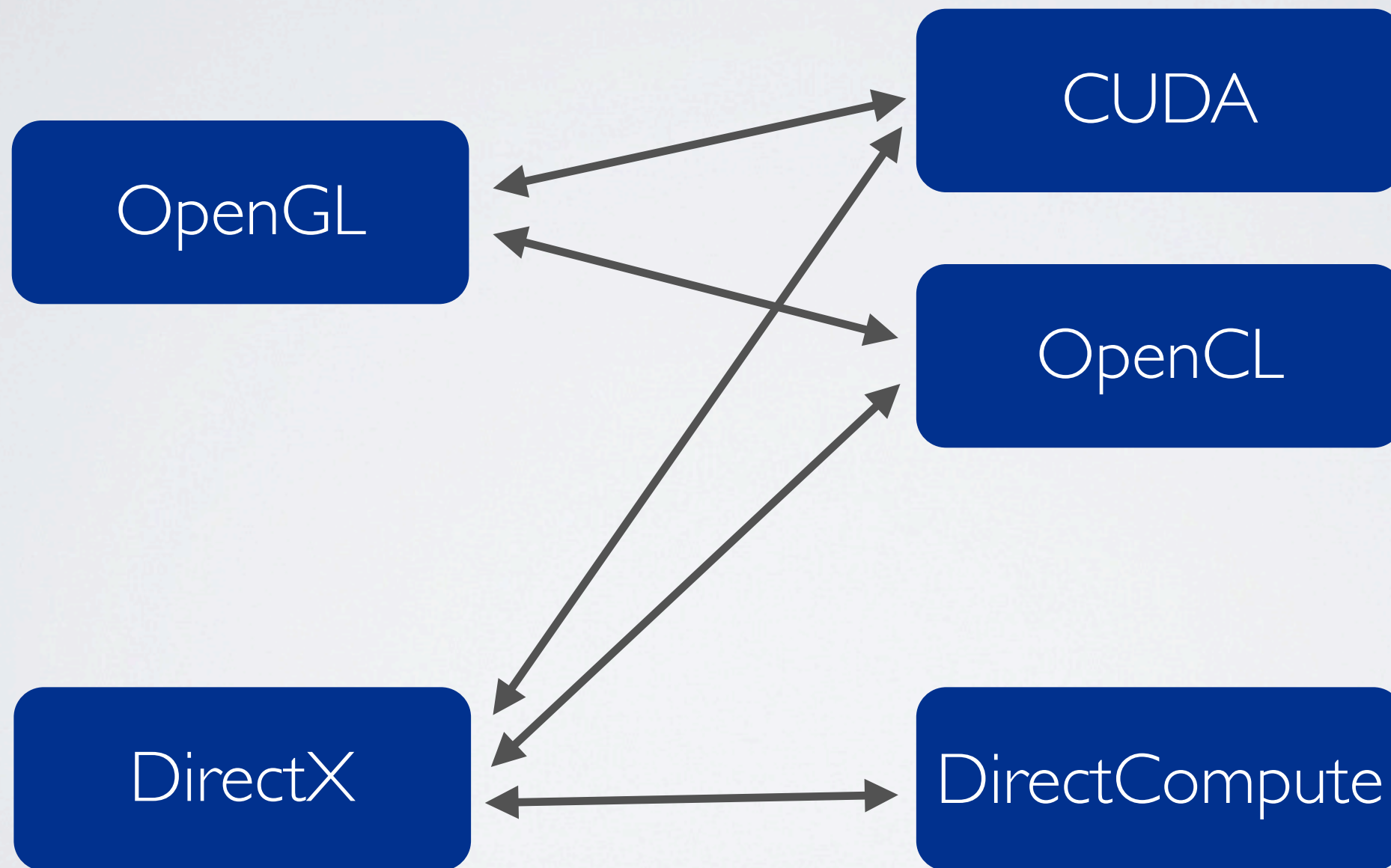
# CONTENTS

- CUDA in the big picture

- When to use it in graphics apps

- How to use it

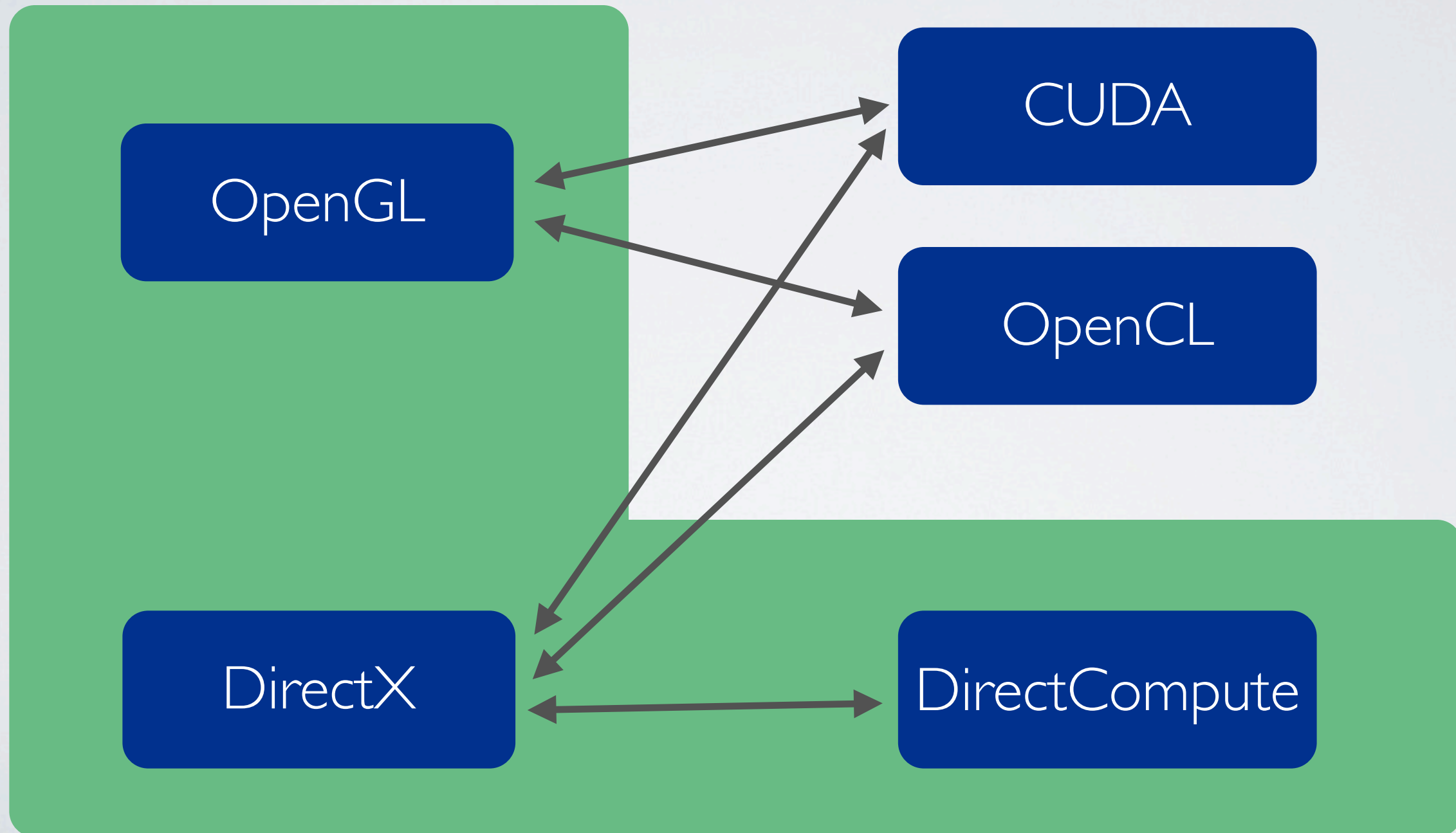- Example: summed area tables (SAT)

# CUDA IN THE BIG PICTURE

Graphics:

GPGPU:

OpenGL

DirectX

CUDA

OpenCL

DirectCompute

# WHEN TO USE CUDA

# WHEN TO USE CUDA

- GPGPU is less limiting, allowing e.g.:

    - Arbitrary memory access patterns

    - On-chip memory communication

- CUDA when:

    - Only targeting NVidia hardware

    - Need advanced hardware features:

        - L1 config, vote functions, function pointers, etc

# WHEN TO USE CUDA

- But choose it only when you really have to

- Don't underestimate optimized OpenGL operations

    - Driver writers know what they are doing

    - You will lose if you try to reinvent the wheel in CUDA

- For optimal performance, you need to know the target HW

    - If you didn't care for performance, you would do it in CPU

# HOW TO USE CUDA

1   Initialize OpenGL

2   Initialize CUDA telling to share an OpenGL context

3   Pass data OpenGL -> CUDA

4   Perform calculations in CUDA

5   Pass results CUDA -> OpenGL

- Main data resources, sharable from OpenGL:

  - Linear allocations (OGL: buffers)

  - CUDA arrays (OGL: textures)

- Execution in kernels

  - Grouping into thread blocks

  - Results into linear allocations
    (can be copied into textures later on)

# HOW TO USE CUDA

- You can choose either the runtime API or the driver API

| Runtime API | Driver API |
|---|---|
| • Mix'n match GPU and CPU functions in same source files<br>• Compile with nvcc into objects<br>• Link into a complete binary<br><br>*Pick me, I'm easy!* | • Compile GPU functions into PTX with nvcc<br>• Compile CPU code separately (e.g. with gcc/g++)<br>• Use CUDA as a normal library, upload the PTX file at runtime<br>• Similar to uploading shader sources in OpenGL |

# EXAMPLE: SUMMED AREA TABLES

# SUMMED AREA TABLES

- Motivation: need to take an average over a region of pixels

  - Generation of texture mip-map levels

  - Fast blur filters (semi-glossy reflections, defocus blur)

Each element $s_{mn}$ of a summed-area table $S$ contains the sum of all elements above and to the left of the original table/texture $T$ [Crow84]

$$s_{mn} = \sum_{i=1}^{m} \sum_{j=1}^{n} t_{ij}$$

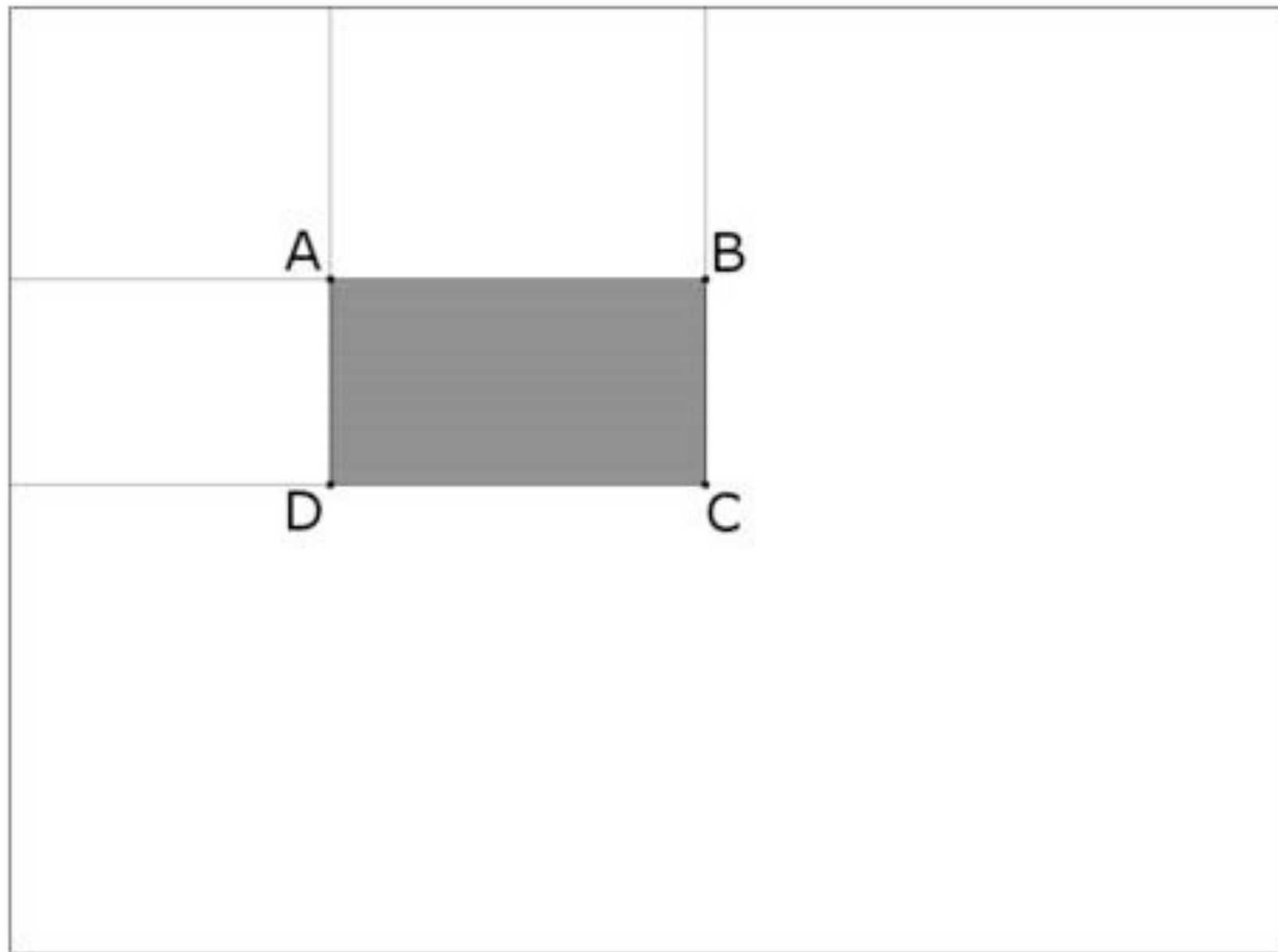|     | 1 | 2 | 3 | 4 |
| --- | --- | --- | --- | --- |
| 1 | 2 | 3 | 2 | 1 |
| 2 | 3 | 0 | 1 | 2 |
| 3 | 1 | 3 | 1 | 0 |
| 4 | 1 | 4 | 2 | 2 |

Original

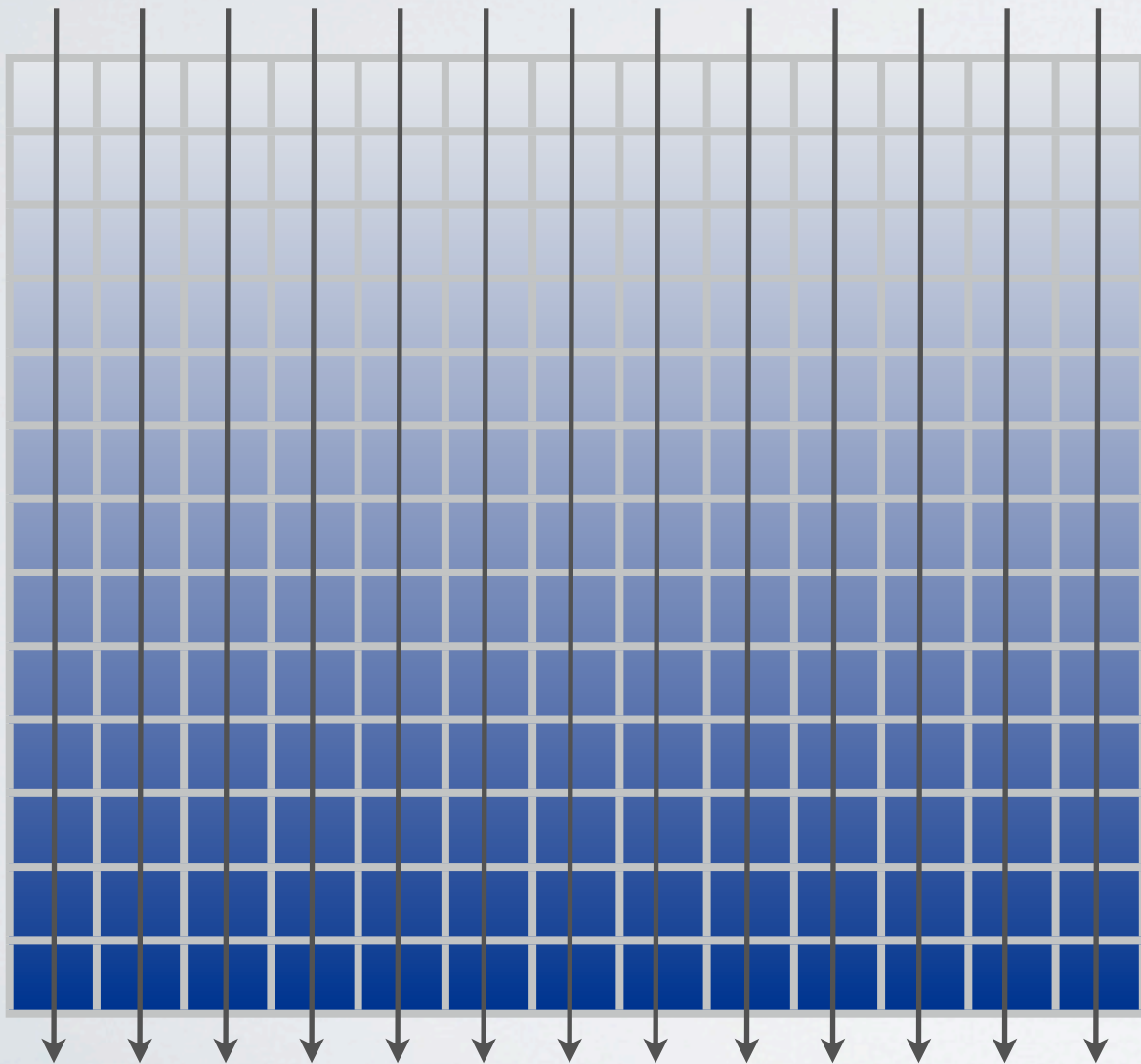| | | | |
| --- | --- | --- | --- |
| 2 | 5 | 7 | 8 |
| 5 | 8 | 11 | 14 |
| 6 | 12 | 16 | 19 |
| 7 | 17 | 23 | 28 |

Summed-area table

Sum of the region: c-b-d+a

# SUMMED AREA TABLES

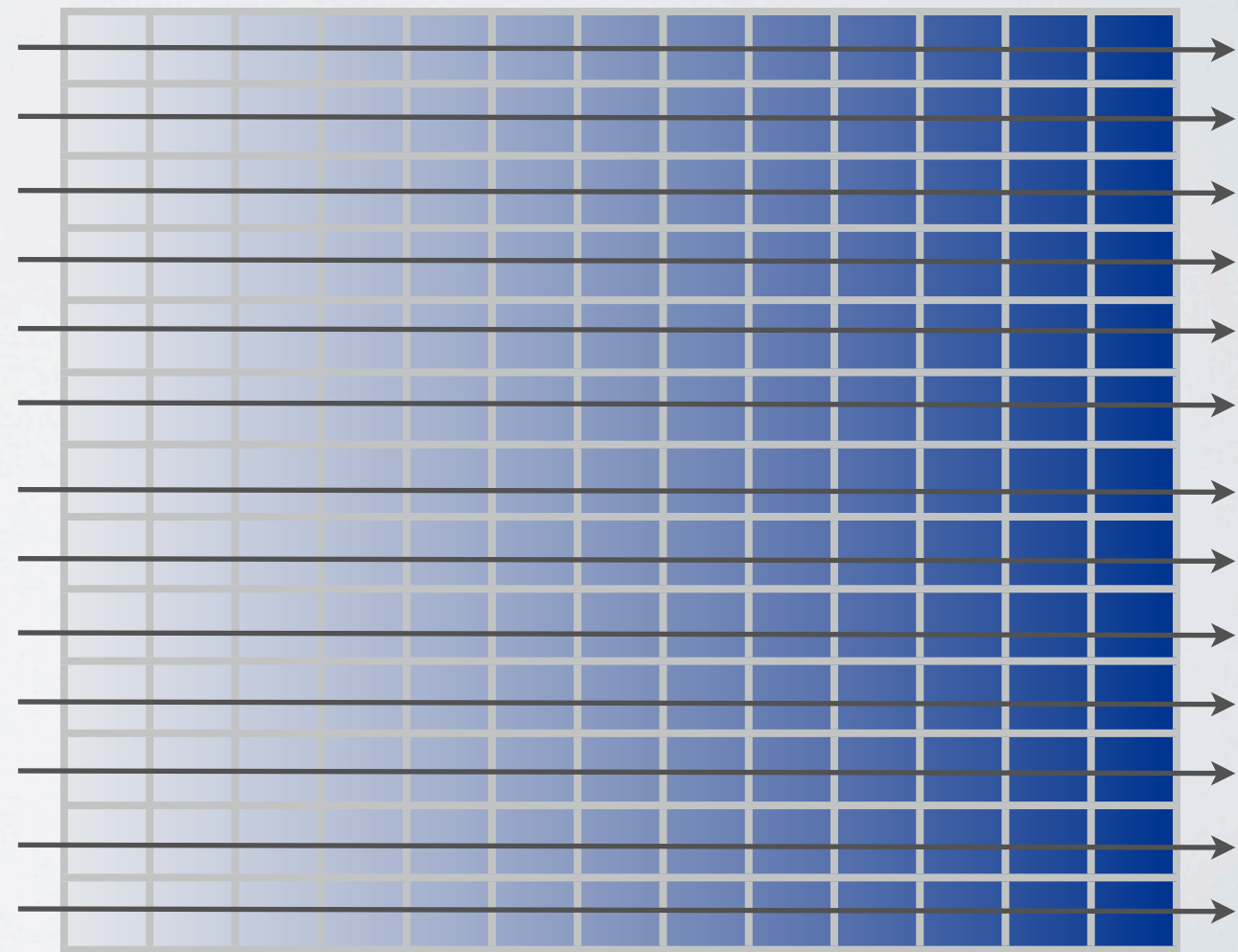Vertical sweep          +          Horizontal sweep

# SUMMED AREA TABLES: IMPLEMENTATION

# SUMMED AREA TABLES

1. Initializations

2. Render a scene in OpenGL into a renderbuffer (via FBO)

3. Pass the renderbuffer + a result texture into CUDA

4. Perform sweeps in CUDA, write results in linear memory

5. Copy results into the result texture

6. Use the result texture in OpenGL in the usual fashion

## 1. Initializations

## 2. Render a scene in OpenGL into a renderbuffer (via FBO)

```
 1  createOpenGLWindow(...);
 2  GLuint renBuf, satTex;
 3
 4  // Initialize CUDA to be used with OpenGL
 5  cudaGLSetGLDevice(0);
 6
 7  // Register CUDA with the renderbuffer (renBuf)
 8  struct cudaGraphicsResource* renBuf_cuda;
 9  cudaGraphicsGLRegisterBuffer(&renBuf_cuda, renBuf,
10          cudaGraphicsRegisterFlagsReadOnly); // CUDA will not write
11
12  // Register CUDA with the result texture (satTex)
13  struct cudaGraphicsResource* satTex_cuda;
14  cudaGraphicsGLRegisterImage(&satTex_cuda, satTex, GL_TEXTURE_2D,
15          cudaGrsphicsRegisterFlagsWriteDiscard); // CUDA will overwrite, not read
16
17  // Allocate work data in CUDA (linear memory)
18  float4 *sweepData;
19  cudaMalloc(&sweepData, WIDTH*HEIGHT*sizeof(float4));
```

1. Initializations

**2. Render a scene in OpenGL into a renderbuffer**

3. Pass the renderbuffer + a result texture into CUDA

# SUMMED AREA TABLES

2. Render a scene in OpenGL into a renderbuffer (via FBO)

## 3. **Pass the renderbuffer + a result texture into CUDA**

4. Perform sweeps in CUDA, write results in linear memory

```
30  // Mapping tells OpenGL to flush changes and not to touch them until unmapped
31  cudaGraphicsMapResources(1, &renBuf_cuda);
32  cudaGraphicsMapResources(1, &satTex_cuda);
33
34  // Read in the CUDA pointers
35  float4 *renBufData; // Pointer to the buffer data
36  cudaGraphicsResourceGetMappedPointer(&renBufData,
37          WIDTH*HEIGHT*sizeof(float4), renBuf_cuda);
38
39  struct cudaArray *satTexArray; // Cast the imported texture as CUDA array
40  cudaGraphicsSubResourceGetMappedArray(&satTexArray,
41          satTex_cuda, 0, 0); // Tex ID (in case a cube map), mip-map level
```

3. Pass the renderbuffer + a result texture into CUDA

## 4. Perform sweeps in CUDA, write results in linear m..

5. Copy results into the result texture

```cpp
50  // In this example we ignore non-divisible-by-64 cases
51  sweepVert<<<WIDTH/64, 64>>>(renBufData, sweepData);
52  sweepHor <<<HEIGHT/64, 64>>>(sweepData);
53
54  __device__ void operator+=(float4 &a, const float4 b) {
55      a.x += b.x; a.y += b.y; a.z += b.z; a.w += b.w;
56  }
57
58  __global__ sweepVert(const float4 *renBufData, float4 *sweepData) {
59      unsigned int myId = blockIdx.x*blockDim.x + threadIdx.x;
60      float4 mySum = make_float4(0.0f, 0.0f, 0.0f, 0.0f);
61
62      for (int row = 0; row < HEIGHT; ++row) {
63          mySum += renBufData[row*WIDTH + myId];
64          sweepData[row*WIDTH + myId] = mySum;
65      }
66  }
67
68  __global__ sweepHor(float4 *sweepData) {
69      unsigned int myId = blockIdx.x*blockDim.x + threadIdx.x;
70      float4 mySum = make_float4(0.0f, 0.0f, 0.0f, 0.0f);
71
72      for (int col = 0; col < WIDTH; ++col) {
73          mySum += sweepData[myId*WIDTH + col];
74          sweepData[myId*WIDTH + col] = mySum;
75      }
76  }
```

1. Perform sweeps in CUDA, write results in linear memory

**5. Copy results into the result texture**
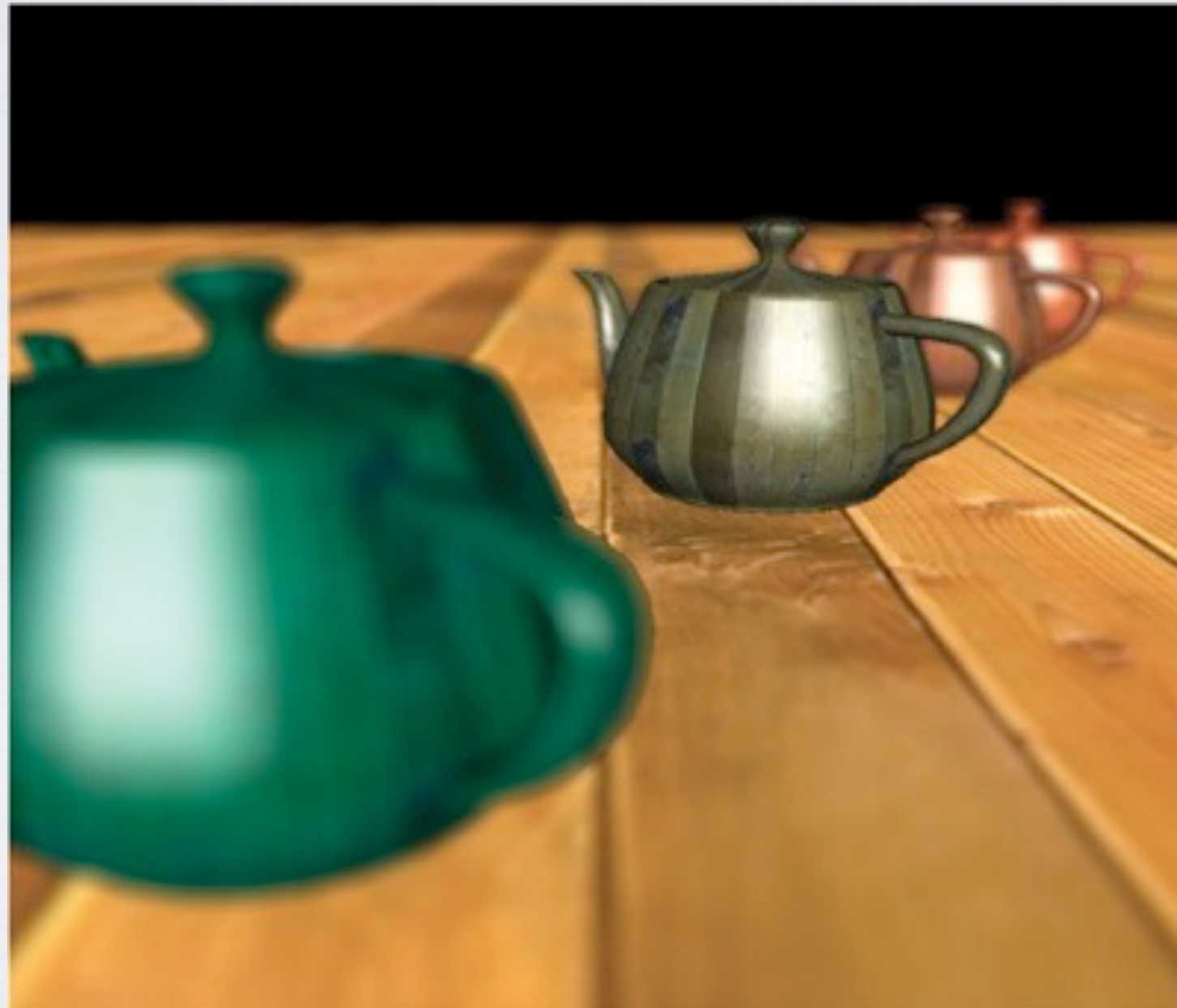
6. Use the result texture in OpenGL in the usual fashion

```
80  cudaMemcpyToArray(satTexArray, 0, 0, // wOffset, hOffset
81          sweepData, WIDTH*HEIGHT*sizeof(float4),
82          cudaMemcpyDeviceToDevice);
83
84  // Unmapping
85  cudaGraphicsUnmapResources(1, &satTex_cuda);
86  cudaGraphicsUnmapResources(1, &renBuf_cuda);
87
88  ...
89
90  // At exit: unregister
91  cudaGraphicsUnregisterResource(satTex_cuda);
92  cudaGraphicsUnregisterResource(renBuf_cuda);
```

5. Copy results into the result texture

**6. Use the result texture in OpenGL in the usual fash..**

# QUESTIONS?