

CPSC/ECE 4780/6780

General-Purpose Computation on Graphical Processing Units (GPGPU)

Lecture 6: CUDA Memories by Examples

Recaps of Last Lecture

- Compute to global memory access (CGMA) ratio
- CUDA device memory types
 - Global memory
 - Shared memory
 - Constant memory
 - Registers
 - Texture memory
- Scope and lifetime of CUDA variables
- Reducing global memory access enhances performance
- Memory coalescing
- Shared memory banks and bank conflicts

CUDA Memories by Examples

- Global memory example
 - Julia Set
- Shared memory example
 - Shared memory bitmap
- Constant memory example
 - Ray tracing
- Texture memory example
 - Heat transfer simulation

Julia Set

- Julia Set is the boundary of a certain class of functions over complex numbers
- Julia set evaluates a simple iterative equation for points in the complex plane
 - A point is not in the set if the process of iterating the equation diverges (grows toward infinity) for that point
 - If the values taken by the equation remain bounded, the point is in the set
- Iterative equation

$$Z_{n+1} = Z_n^2 + C$$

involves squaring the current value and adding a constant to get the next value of the equation

CPU Julia Set – Main Function

- Create the appropriate size bitmap image using a utility library provided by NVIDIA
- Pass a pointer to the bitmap data to the kernel function
- Display the bitmap and exit

```
int main( void ) {  
    CPUBitmap bitmap( DIM, DIM );  
    unsigned char *ptr = bitmap.get_ptr();  
  
    kernel( ptr );  
  
    bitmap.display_and_exit();  
}
```

CPU Julia Set – Computation Kernel

- Iterate through all points we care to render
- Call `julia()` on each point to determine its membership in the Julia Set
 - Return 1 if the point is in the set
 - Return 0 otherwise

```
void kernel( unsigned char *ptr ){
    for (int y=0; y<DIM; y++) {
        for (int x=0; x<DIM; x++) {
            int offset = x + y * DIM;

            int juliaValue = julia( x, y );
            ptr[offset*4 + 0] = 255 * juliaValue;
            ptr[offset*4 + 1] = 0;
            ptr[offset*4 + 2] = 0;
            ptr[offset*4 + 3] = 255;
        }
    }
}
```

CPU Julia Set – Julia()

- Translate pixel coordinate to a coordinate in complex space
- Set the arbitrary complex-valued constant C
- Iterate 200 times and determine whether the point is in or out of the Julia Set

```
int julia( int x, int y ) {  
    const float scale = 1.5;  
    float jx = scale * (float)(DIM/2 - x)/(DIM/2);  
    float jy = scale * (float)(DIM/2 - y)/(DIM/2);  
  
    cuComplex c(-0.8, 0.156);  
    cuComplex a(jx, jy);  
  
    int i = 0;  
    for (i=0; i<200; i++) {  
        a = a * a + c;  
        if (a.magnitude2() > 1000)  
            return 0;  
    }  
  
    return 1;  
}
```

Center the complex plane at the image center

Zoom in or out

Ensure the image spans the range of (-1.0, 1.0)

CPU Julia Set – cuComplex Struct

- Define a generic structure to store complex numbers

```
struct cuComplex {  
    float    r;  
    float    i;  
    cuComplex( float a, float b ) : r(a), i(b) {}  
    float magnitude2( void ) { return r * r + i * i; }  
    cuComplex operator*(const cuComplex& a) {  
        return cuComplex(r*a.r - i*a.i, i*a.r + r*a.i);  
    }  
    cuComplex operator+(const cuComplex& a) {  
        return cuComplex(r+a.r, i+a.i);  
    }  
};
```


GPU Julia Set – Main Function

- Declare a pointer `dev_bitmap` to hold a copy of the data on the device
- Allocate memory using `cudaMalloc`
- Specify a two-dimensional grid of blocks for kernel computation
- Copy the results back to the host
- Display the bitmap and exit

```
int main( void ) {  
    CPUBitmap bitmap( DIM, DIM );  
    unsigned char *dev_bitmap;  
  
    HANDLE_ERROR( cudaMalloc( (void**)&dev_bitmap, bitmap.  
        image_size() ) );  
  
    dim3    grid(DIM,DIM);  
    kernel<<<grid,1>>>( dev_bitmap );  
  
    HANDLE_ERROR( cudaMemcpy( bitmap.get_ptr(), dev_bitmap,  
        bitmap.image_size(),  
        cudaMemcpyDeviceToHost ) );  
    HANDLE_ERROR( cudaFree( dev_bitmap ) );  
    bitmap.display_and_exit();  
}
```

CPU version counterpart

```
int main( void ) {  
    CPUBitmap bitmap( DIM, DIM );  
    unsigned char *ptr = bitmap.get_ptr();  
  
    kernel( ptr );  
  
    bitmap.display_and_exit();  
}
```

GPU Julia Set – Computation Kernel

- Declare kernel() as a `__global__` function so it runs on the device but can be called from the host
- Calculate pixel index using `blockIdx` instead of nested `for()` loops
- Calculate a linear offset using `gridDim`
- Call `Julia()` on each to determine membership in the Julia Set

```
__global__ void kernel( unsigned char *ptr ) {  
    // map from blockIdx to pixel position  
    int x = blockIdx.x;  
    int y = blockIdx.y;  
    int offset = x + y * gridDim.x;  
  
    // now calculate the value at that position  
    int juliaValue = julia( x, y );  
    ptr[offset*4 + 0] = 255 * juliaValue;  
    ptr[offset*4 + 1] = 0;  
    ptr[offset*4 + 2] = 0;  
    ptr[offset*4 + 3] = 255;  
}
```

```
void kernel( unsigned char *ptr ){  
    for (int y=0; y<DIM; y++) {  
        for (int x=0; x<DIM; x++) {  
            int offset = x + y * DIM;  
  
            int juliaValue = julia( x, y );  
            ptr[offset*4 + 0] = 255 * juliaValue;  
            ptr[offset*4 + 1] = 0;  
            ptr[offset*4 + 2] = 0;  
            ptr[offset*4 + 3] = 255;  
        }  
    }  
}
```

CPU version counterpart

GPU Julia Set – Julia()

```
__device__ int julia( int x, int y ) {  
    const float scale = 1.5;  
    float jx = scale * (float)(DIM/2 - x)/(DIM/2);  
    float jy = scale * (float)(DIM/2 - y)/(DIM/2);  
  
    cuComplex c(-0.8, 0.156);  
    cuComplex a(jx, jy);  
  
    int i = 0;  
    for (i=0; i<200; i++) {  
        a = a * a + c;  
        if (a.magnitude2() > 1000)  
            return 0;  
    }  
  
    return 1;  
}
```

```
int julia( int x, int y ) {  
    const float scale = 1.5;  
    float jx = scale * (float)(DIM/2 - x)/(DIM/2);  
    float jy = scale * (float)(DIM/2 - y)/(DIM/2);  
  
    cuComplex c(-0.8, 0.156);  
    cuComplex a(jx, jy);  
  
    int i = 0;  
    for (i=0; i<200; i++) {  
        a = a * a + c;  
        if (a.magnitude2() > 1000)  
            return 0;  
    }  
  
    return 1;  
}
```

CPU version counterpart

GPU Julia Set – cuComplex Struct

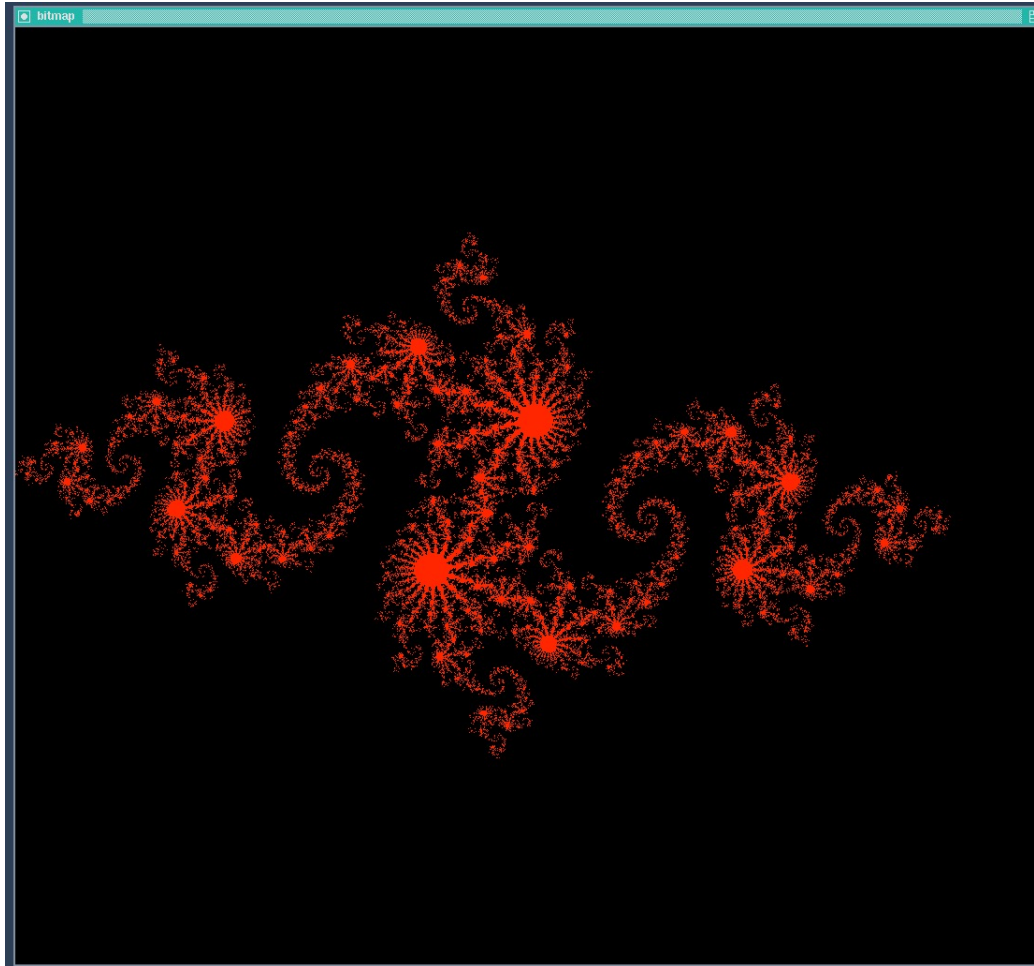
```
struct cuComplex {
    float    r;
    float    i;
    __host__ __device__ cuComplex( float a, float b ) : r(a), i(b)
        {}
    __device__ float magnitude2( void ) {
        return r * r + i * i;
    }
    __device__ cuComplex operator*(const cuComplex& a) {
        return cuComplex(r*a.r - i*a.i, i*a.r + r*a.i);
    }
    __device__ cuComplex operator+(const cuComplex& a) {
        return cuComplex(r+a.r, i+a.i);
    }
};
```

```
struct cuComplex {
    float    r;
    float    i;
    cuComplex( float a, float b ) : r(a), i(b) {}
    float magnitude2( void ) { return r * r + i * i; }
    cuComplex operator*(const cuComplex& a) {
        return cuComplex(r*a.r - i*a.i, i*a.r + r*a.i);
    }
    cuComplex operator+(const cuComplex& a) {
        return cuComplex(r+a.r, i+a.i);
    }
};
```

CPU version counterpart

GPU Julia Set – Compilation

- `nvcc -lglut -lGL -lGLU juliaSet.cu`



Shared Memory Bitmap – Main Function

```
struct DataBlock {
    unsigned char    *dev_bitmap;
};

int main( void ) {
    DataBlock    data;
    CPUBitmap bitmap( DIM, DIM, &data );
    unsigned char    *dev_bitmap;

    HANDLE_ERROR( cudaMalloc( (void**)&dev_bitmap,
                               bitmap.image_size() ) );
    data.dev_bitmap = dev_bitmap;

    dim3    grids( DIM/16, DIM/16 );
    dim3    threads( 16, 16 );
    kernel<<<grids, threads>>>( dev_bitmap );

    HANDLE_ERROR( cudaMemcpy( bitmap.get_ptr(), dev_bitmap,
                               bitmap.image_size(),
                               cudaMemcpyDeviceToHost ) );

    HANDLE_ERROR( cudaFree( dev_bitmap ) );

    bitmap.display_and_exit();
}
```

→ Launch multiple threads per block

Shared Memory Bitmap – Kernel Function

```
__global__ void kernel( unsigned char *ptr ) {  
    // map from threadIdx/BlockIdx to pixel position  
    int x = threadIdx.x + blockIdx.x * blockDim.x;  
    int y = threadIdx.y + blockIdx.y * blockDim.y;  
    int offset = x + y * blockDim.x * gridDim.x;
```

→ Compute x and y location in the output image

```
    __shared__ float shared[16][16];
```

→ Use a shared memory buffer to cache the computation

```
    // now calculate the value at that position
```

```
    const float period = 128.0f;
```

```
    shared[threadIdx.x][threadIdx.y] =  
        255 * (sinf(x*2.0f*PI/ period) + 1.0f) *  
        (sinf(y*2.0f*PI/ period) + 1.0f) / 4.0f;
```

→ Each thread computes a value to be stored into the buffer

```
    __syncthreads();
```

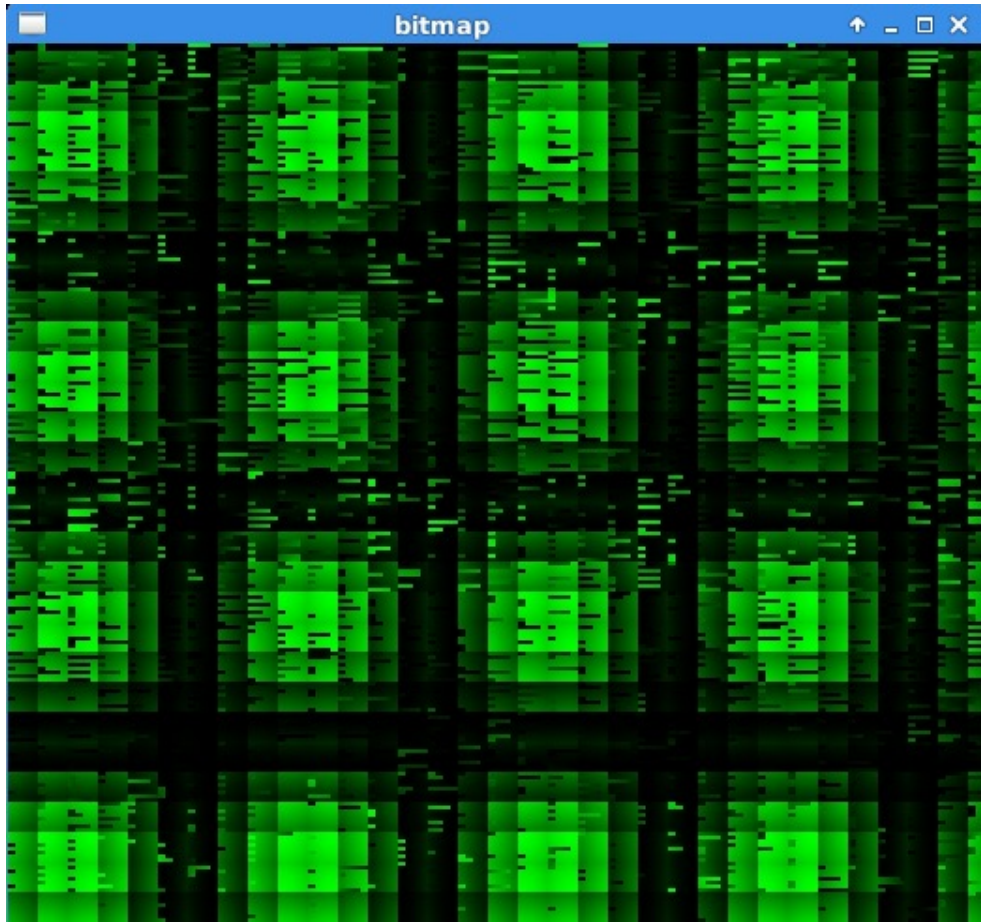
→ To make sure all the values are written before the subsequent read from it

```
    ptr[offset*4 + 0] = 0;  
    ptr[offset*4 + 1] = shared[15-threadIdx.x][15-threadIdx.y];  
    ptr[offset*4 + 2] = 0;  
    ptr[offset*4 + 3] = 255;
```

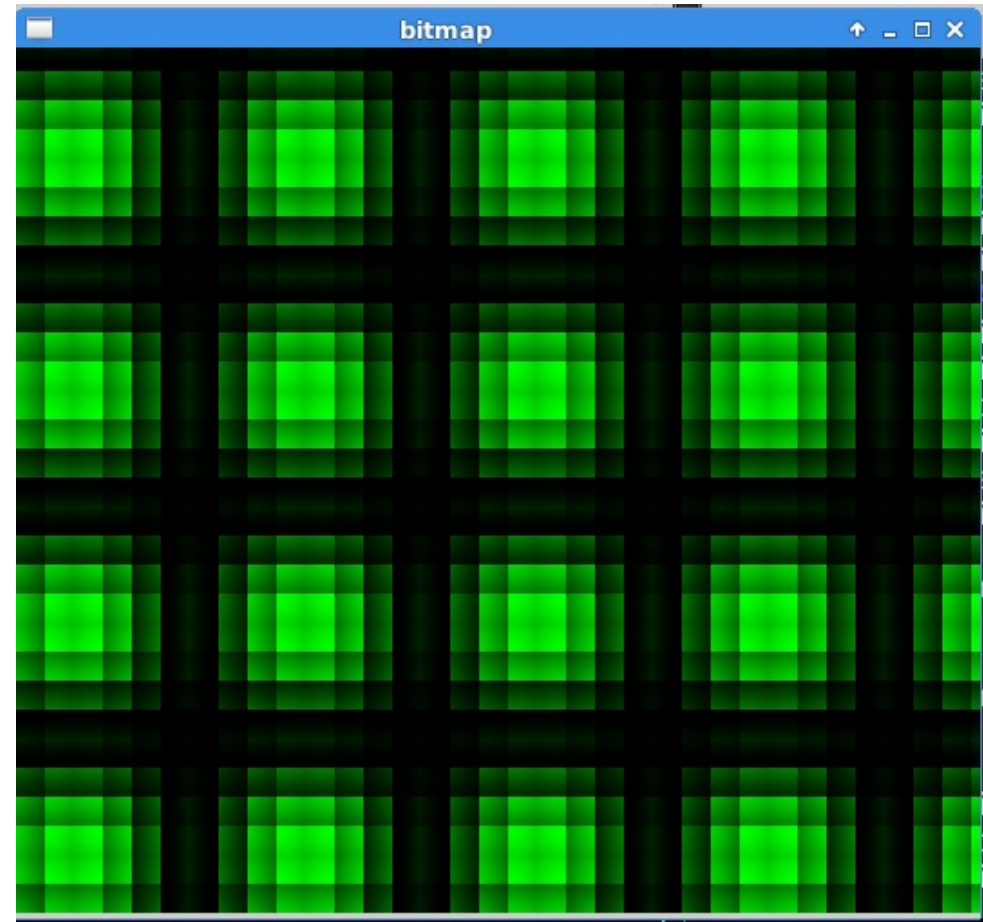
→ Store the values back to the pixel, reversing the order of x and y

```
}
```

Shared Memory Bitmap – Comparison on Correctness



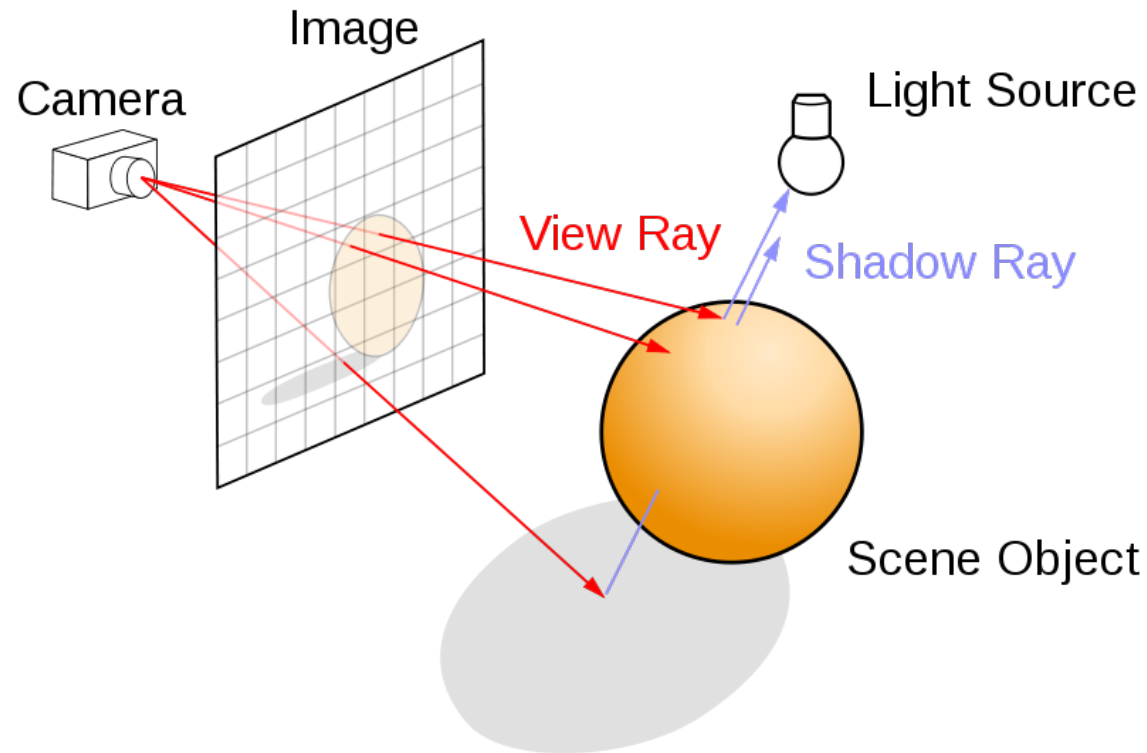
No `__syncthreads()`



with `__syncthreads()`

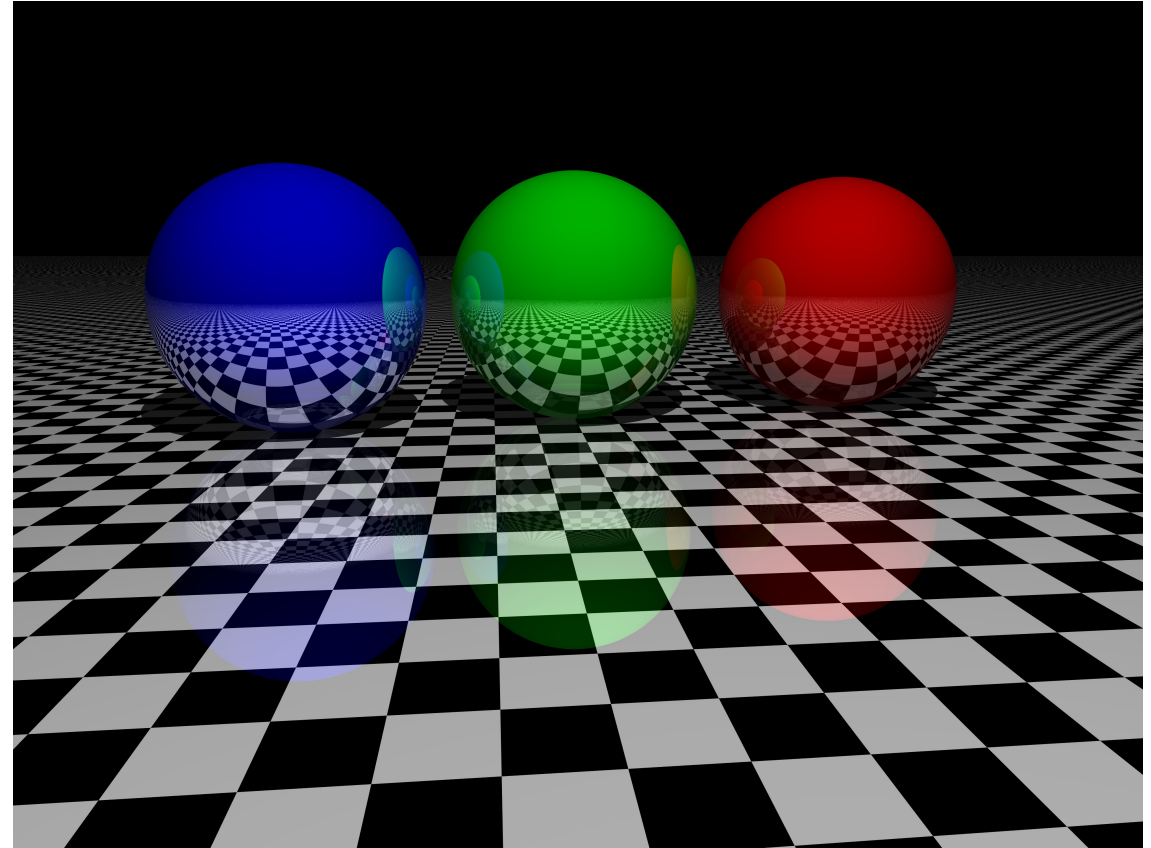
Ray Tracing

- Ray tracing is a technique for generating an image by tracing the path of light as pixels in an image plane and simulating the effects of its encounters with virtual objects



How Ray Tracing Works

- 1. For each pixel create a ray
- 2. Follow the ray into the scene and find the closest object that is intersected and assign the color of the object to the pixel
- 3. If the object is made of glass, mirror or something similar, find the reflected and/or refracted ray and trace on
- 4. Else, compute the color of the pixel based on light sources and material information



A Basic Ray Tracer

- Restrictions
 - Support scenes of spheres only
 - The camera is restricted to the z-axis, facing the origin
 - No support of any lighting of the scene
 - Each sphere will be assigned a color and shaded with precomputed function if they are visible
- What will it do?
 - Fire a ray from each pixel and keep track of which rays hit which spheres
 - Track the depth of each of these hits to find the closest sphere

Ray Tracing on GPU – Sphere Struct

```
#define INF      2e10f

struct Sphere {
    float    r,b,g;
    float    radius;
    float    x,y,z;
    __device__ float hit( float ox, float oy, float *n ) {
        float dx = ox - x;
        float dy = oy - y;
        if (dx*dx + dy*dy < radius*radius) {
            float dz = sqrtf( radius*radius - dx*dx - dy*dy );
            *n = dz / sqrtf( radius * radius );
            return dz + z;
        }
        return -INF;
    }
};
```

- The “hit” method computes whether the ray intersects the sphere

Ray Tracing on GPU – Main Function

```
74 int main( void ) {
75     CPUBitmap bitmap( DIM, DIM );
76     unsigned char *dev_bitmap;
77     Sphere *s;
78
79     // allocate memory on the GPU for the output bitmap
80     HANDLE_ERROR( cudaMalloc( (void*)&dev_bitmap,
81                               bitmap.image_size() ) );
82     // allocate memory for the Sphere dataset
83     HANDLE_ERROR( cudaMalloc( (void*)&s,
84                               sizeof(Sphere) * SPHERES ) );
85
86     // allocate temp memory, initialize it, copy to
87     // memory on the GPU, then free our temp memory
88     Sphere *temp_s = (Sphere*)malloc( sizeof(Sphere) * SPHERES );
89     for (int i=0; i<SPHERES; i++) {
90         temp_s[i].r = rnd( 1.0f );
91         temp_s[i].g = rnd( 1.0f );
92         temp_s[i].b = rnd( 1.0f );
93         temp_s[i].x = rnd( 1000.0f ) - 500;
94         temp_s[i].y = rnd( 1000.0f ) - 500;
95         temp_s[i].z = rnd( 1000.0f ) - 500;
96         temp_s[i].radius = rnd( 100.0f ) + 20;
97     }
98
99     HANDLE_ERROR( cudaMemcpy( s, temp_s,
100                              sizeof(Sphere) * SPHERES,
101                              cudaMemcpyHostToDevice ) );
102     free( temp_s );
103
104     // generate a bitmap from our sphere data
105     dim3 grids(DIM/16,DIM/16);
106     dim3 threads(16,16);
107     kernel<<<grids,threads>>>( s, dev_bitmap );
108
109     // copy our bitmap back from the GPU for display
110     HANDLE_ERROR( cudaMemcpy( bitmap.get_ptr(), dev_bitmap,
111                               bitmap.image_size(),
112                               cudaMemcpyDeviceToHost ) );
113
114     HANDLE_ERROR( cudaFree( dev_bitmap ) );
115     HANDLE_ERROR( cudaFree( s ) );
116
117     // display
118     bitmap.display_and_exit();
119 }
```

Ray Tracing on GPU – Kernel Function

```
44 __global__ void kernel( Sphere *s, unsigned char *ptr ) {  
45     // map from threadIdx/BlockIdx to pixel position  
46     int x = threadIdx.x + blockIdx.x * blockDim.x;  
47     int y = threadIdx.y + blockIdx.y * blockDim.y;  
48     int offset = x + y * blockDim.x * gridDim.x;  
49     float ox = (x - DIM/2);  
50     float oy = (y - DIM/2);  
51  
52     float r=0, g=0, b=0;  
53     float maxz = -INF;  
54     for(int i=0; i<SPHERES; i++) {  
55         float n;  
56         float t = s[i].hit( ox, oy, &n );  
57         if (t > maxz) {  
58             float fscale = n;  
59             r = s[i].r * fscale;  
60             g = s[i].g * fscale;  
61             b = s[i].b * fscale;  
62             maxz = t;  
63         }  
64     }  
65  
66     ptr[offset*4 + 0] = (int) (r * 255);  
67     ptr[offset*4 + 1] = (int) (g * 255);  
68     ptr[offset*4 + 2] = (int) (b * 255);  
69     ptr[offset*4 + 3] = 255;  
70 }
```

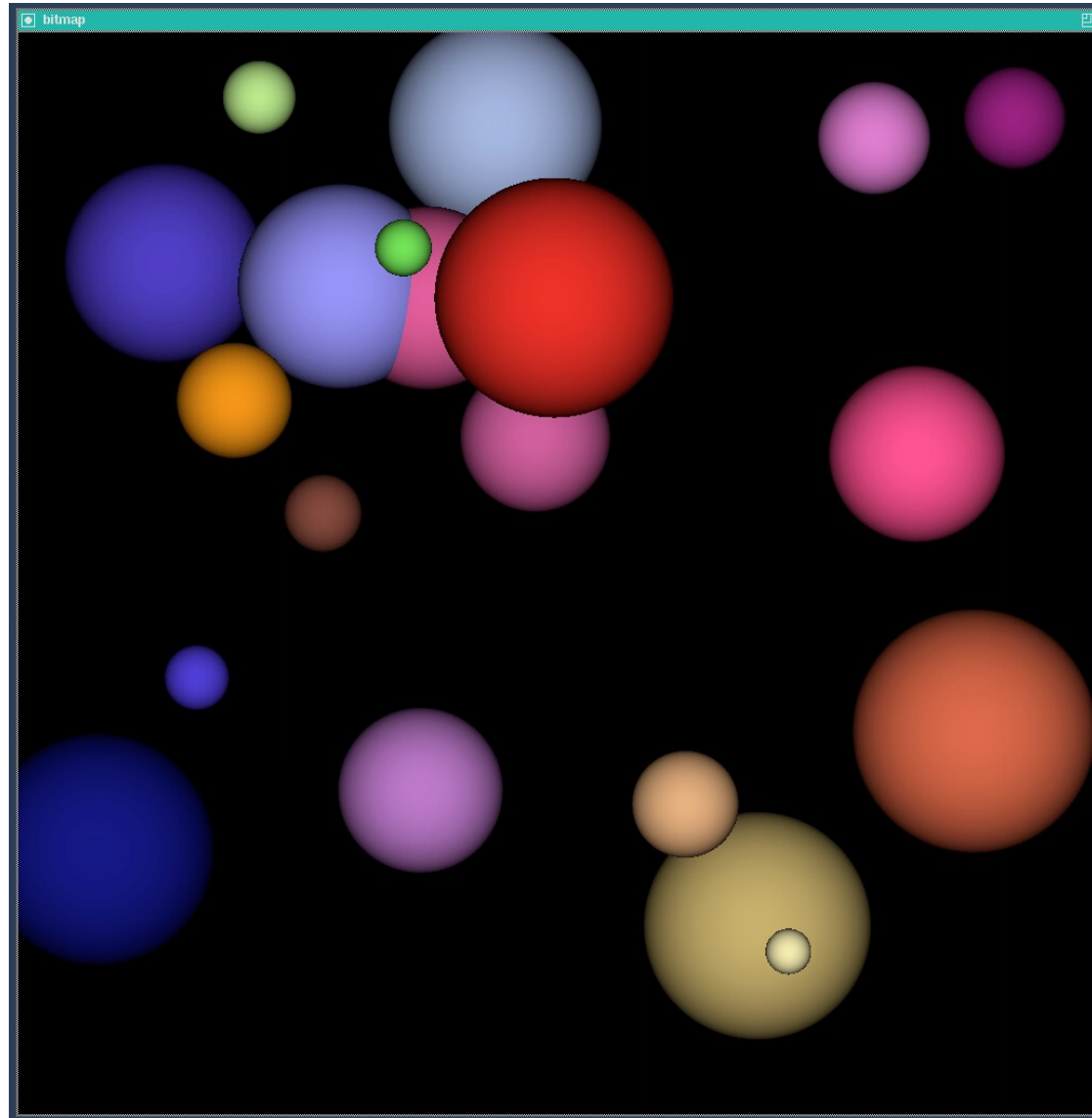
Each thread generates one pixel for the output image

Initialize the background to be black if no spheres have been hit

Iterate through the array of spheres to check each sphere for intersection

Store the current color into the output image

Ray Tracing on GPU – Screenshot



Ray Tracing on GPU with Constant Memory

- Store the input array of spheres in constant memory

- Instead of declaring the array like

`Sphere *s;`

→ Declare a pointer and use `cudaMalloc()` to allocate GPU memory for it

We add the modifier `__constant__` before it:

`__constant__ Sphere s[SPHERES];`

→ Statically allocate the space in constant memory

- Need to commit to a size for this array at compile-time
- Changes in `main()`:

```
HANDLE_ERROR( cudaMemcpy( s, temp_s,
                          sizeof(Sphere) * SPHERES,
                          cudaMemcpyHostToDevice ) );
```



```
HANDLE_ERROR( cudaMemcpyToSymbol( s, temp_s,
                                   sizeof(Sphere) * SPHERES) );
```


Benefits of Using Constant Memory for Ray Tracing on GPU

- Receive the data in a half-warp broadcast is efficient
- Retrieve the data from the constant memory cache is fast

Measuring Performance with Events

- CUDA event API: An event in CUDA is essentially a GPU time stamp that is recorded at a user-specified point in time.
- Two steps:
 - Creating an event
 - Subsequently recording an event

```
cudaEvent_t start, stop;  
cudaEventCreate(&start);  
cudaEventCreate(&stop);  
cudaEventRecord(start, 0);  
  
// do some work on the GPU  
  
cudaEventRecord(stop, 0);  
cudaEventSynchronize(stop);
```

Make sure all GPU work before the stop event has completed so it is safe to read the time stamp recorded

- Do not use CUDA events to time mixtures of device and host code!

Measuring Ray Tracing Performance

```
80 // capture the start time
81 cudaEvent_t start, stop;
82 HANDLE_ERROR( cudaEventCreate( &start ) );
83 HANDLE_ERROR( cudaEventCreate( &stop ) );
84 HANDLE_ERROR( cudaEventRecord( start, 0 ) );
```

```
119 // get stop time, and display the timing results
120 HANDLE_ERROR( cudaEventRecord( stop, 0 ) );
121 HANDLE_ERROR( cudaEventSynchronize( stop ) );
122 float elapsedTime;
123 HANDLE_ERROR( cudaEventElapsedTime( &elapsedTime,
124                                     start, stop ) );
125 printf( "Time to generate:  %3.1f ms\n", elapsedTime );
126
127 HANDLE_ERROR( cudaEventDestroy( start ) );
128 HANDLE_ERROR( cudaEventDestroy( stop ) );
```

Computes the elapsed time between two previously recorded events

Free the created events

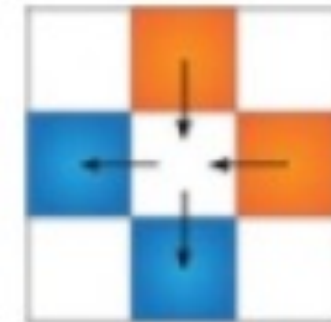
Heat Transfer Simulation

- Simple heating model
 - A rectangular room divide into a grid with “heaters” of various temperature
 - Cells with heaters in them always remain a constant temperature
 - Heat “flows” between a cell and its neighbors
 - New temperature in a cell is the sum of differences between its temperature and the temperatures of its neighbor

$$T_{NEW} = T_{OLD} + \sum_{NEIGHBORS} k \cdot (T_{NEIGHBOR} - T_{OLD})$$



$$T_{NEW} = T_{OLD} + k \cdot (T_{TOP} + T_{BOTTOM} + T_{LEFT} + T_{RIGHT} - 4 \cdot T_{OLD})$$



Heat Transfer Simulation – Computing Temperature Updates

- Step 1. Given some grid of input temperatures, copy the temperature of cells with heaters to this grid, overwriting any previously computed temperatures in these cells
 - `copy_const_kernel()`
- Step 2. Given the input temperature grid, compute the output temperatures based on the update equation
 - `Blend_kernel()`
- Step 3. Swap the input and output buffers in preparation of the next time step

Heat Transfer Simulation – Step 1

```
__global__ void copy_const_kernel( float *iptr,
                                   const float *cptr ) {
    // map from threadIdx/BlockIdx to pixel position
    int x = threadIdx.x + blockIdx.x * blockDim.x;
    int y = threadIdx.y + blockIdx.y * blockDim.y;
    int offset = x + y * blockDim.x * gridDim.x;

    if (cptr[offset] != 0) iptr[offset] = cptr[offset];
}
```

Heat Transfer Simulation – Step 2

```
__global__ void blend_kernel( float *outSrc,  
                             const float *inSrc ) {  
    // map from threadIdx/BlockIdx to pixel position  
    int x = threadIdx.x + blockIdx.x * blockDim.x;  
    int y = threadIdx.y + blockIdx.y * blockDim.y;  
    int offset = x + y * blockDim.x * gridDim.x;
```

```
    int left = offset - 1;  
    int right = offset + 1;  
    if (x == 0)    left++;  
    if (x == DIM-1) right--;  
  
    int top = offset - DIM;  
    int bottom = offset + DIM;  
    if (y == 0)    top += DIM;  
    if (y == DIM-1) bottom -= DIM;
```

→ Determine the offsets of the left, right, top, and bottom neighbors

$$T_{NEW} = T_{OLD} + k \cdot (T_{TOP} + T_{BOTTOM} + T_{LEFT} + T_{RIGHT} - 4 \cdot T_{OLD})$$

```
    outSrc[offset] = inSrc[offset] + SPEED * ( inSrc[top] +  
                                              inSrc[bottom] + inSrc[left] + inSrc[right] -  
                                              inSrc[offset]*4);
```

→ Adding the old temperature and the scaled differences of that temperature and the cell's neighbors' temperature

```
}
```


Animating the Simulation – anim_gpu()

```
void anim_gpu( DataBlock *d, int ticks ) {
```

→ Called by the animation framework on every frame

.....

```
for (int i=0; i<90; i++) {  
    copy_const_kernel<<<blocks,threads>>>( d->dev_inSrc,  
                                              d->dev_constSrc );  
    blend_kernel<<<blocks,threads>>>( d->dev_outSrc,  
                                       d->dev_inSrc );  
    swap( d->dev_inSrc, d->dev_outSrc );  
}
```

→ Compute a single time step of the simulation as described by Step 1 to 3

```
float_to_color<<<blocks,threads>>>( d->output_bitmap,  
                                     d->dev_inSrc );
```

→ Convert the temperatures to colors

```
HANDLE_ERROR( cudaMemcpy( bitmap->get_ptr(),  
                          d->output_bitmap,  
                          bitmap->image_size(),  
                          cudaMemcpyDeviceToHost ) );
```

→ Copy the resultant image back to the CPU

.....

```
}
```


Using Texture Memory for Heat Transfer Simulation – Set Up Textures

- First, declare inputs as texture references
 - `texture<float> texConstSrc;`
 - `texture<float> texIn;`
 - `texture<float> texOut;`
- Next, bind the references to the memory buffer using `cudaBindTexture()`

```
HANDLE_ERROR( cudaMalloc( (void*)&data.dev_inSrc,  
                           imageSize ) );  
HANDLE_ERROR( cudaMalloc( (void*)&data.dev_outSrc,  
                           imageSize ) );  
HANDLE_ERROR( cudaMalloc( (void*)&data.dev_constSrc,  
                           imageSize ) );  
  
HANDLE_ERROR( cudaBindTexture( NULL, texConstSrc,  
                               data.dev_constSrc,  
                               imageSize ) );  
  
HANDLE_ERROR( cudaBindTexture( NULL, texIn,  
                               data.dev_inSrc,  
                               imageSize ) );  
  
HANDLE_ERROR( cudaBindTexture( NULL, texOut,  
                               data.dev_outSrc,  
                               imageSize ) );
```

Using Texture Memory for Heat Transfer Simulation – Reading from Texture Memory

- Modify `blend_kernel()` to use `tex1Dfetch()` when reading from memory
- Texture references must be declared globally at file scope
 - No longer pass the input and output buffers as parameters to `blend_kernel()`
 - Pass to `blend_kernel()` a boolean flag `dstOut` that indicates which buffer to use as input and which to use as output

```
float    t, l, c, r, b;
if (dstOut) {
    t = tex1Dfetch(texIn,top);
    l = tex1Dfetch(texIn,left);
    c = tex1Dfetch(texIn,offset);
    r = tex1Dfetch(texIn,right);
    b = tex1Dfetch(texIn,bottom);

} else {
    t = tex1Dfetch(texOut,top);
    l = tex1Dfetch(texOut,left);
    c = tex1Dfetch(texOut,offset);
    r = tex1Dfetch(texOut,right);
    b = tex1Dfetch(texOut,bottom);
}
dst[offset] = c + SPEED * (t + b + r + l - 4 * c);
```

```
__global__ void blend_kernel( float *dst,
                             bool dstOut ) {
```

Using Texture Memory for Heat Transfer Simulation – Changes in copy_const_kernel()

```
__global__ void copy_const_kernel( float *iptr ) {  
    // map from threadIdx/BlockIdx to pixel position  
    int x = threadIdx.x + blockIdx.x * blockDim.x;  
    int y = threadIdx.y + blockIdx.y * blockDim.y;  
    int offset = x + y * blockDim.x * gridDim.x;  
  
    float c = tex1Dfetch(texConstSrc,offset);  
    if (c != 0)  
        iptr[offset] = c;  
}
```

```
__global__ void copy_const_kernel( float *iptr,  
                                   const float *cptr ) {  
    // map from threadIdx/BlockIdx to pixel position  
    int x = threadIdx.x + blockIdx.x * blockDim.x;  
    int y = threadIdx.y + blockIdx.y * blockDim.y;  
    int offset = x + y * blockDim.x * gridDim.x;  
  
    if (cptr[offset] != 0) iptr[offset] = cptr[offset];  
}
```

Non-texture version counterpart

Using Texture Memory for Heat Transfer Simulation – Changes in anim_gpu()

```
volatile bool dstOut = true;
for (int i=0; i<90; i++) {
    float *in, *out;
    if (dstOut) {
        in = d->dev_inSrc;
        out = d->dev_outSrc;
    } else {
        out = d->dev_inSrc;
        in = d->dev_outSrc;
    }
    copy_const_kernel<<<blocks,threads>>>( in );
    blend_kernel<<<blocks,threads>>>( out, dstOut );
    dstOut = !dstOut;
}
float_to_color<<<blocks,threads>>>( d->output_bitmap,
                                     d->dev_inSrc );

HANDLE_ERROR( cudaMemcpy( bitmap->get_ptr(),
                          d->output_bitmap,
                          bitmap->image_size(),
                          cudaMemcpyDeviceToHost ) );
```

Using Texture Memory for Heat Transfer Simulation – Unbind Texture

```
// clean up memory allocated on the GPU
void anim_exit( DataBlock *d ) {
    cudaUnbindTexture( texIn );
    cudaUnbindTexture( texOut );
    cudaUnbindTexture( texConstSrc );
    HANDLE_ERROR( cudaFree( d->dev_inSrc ) );
    HANDLE_ERROR( cudaFree( d->dev_outSrc ) );
    HANDLE_ERROR( cudaFree( d->dev_constSrc ) );

    HANDLE_ERROR( cudaEventDestroy( d->start ) );
    HANDLE_ERROR( cudaEventDestroy( d->stop ) );
}
```


Using Two-Dimensional Texture Memory

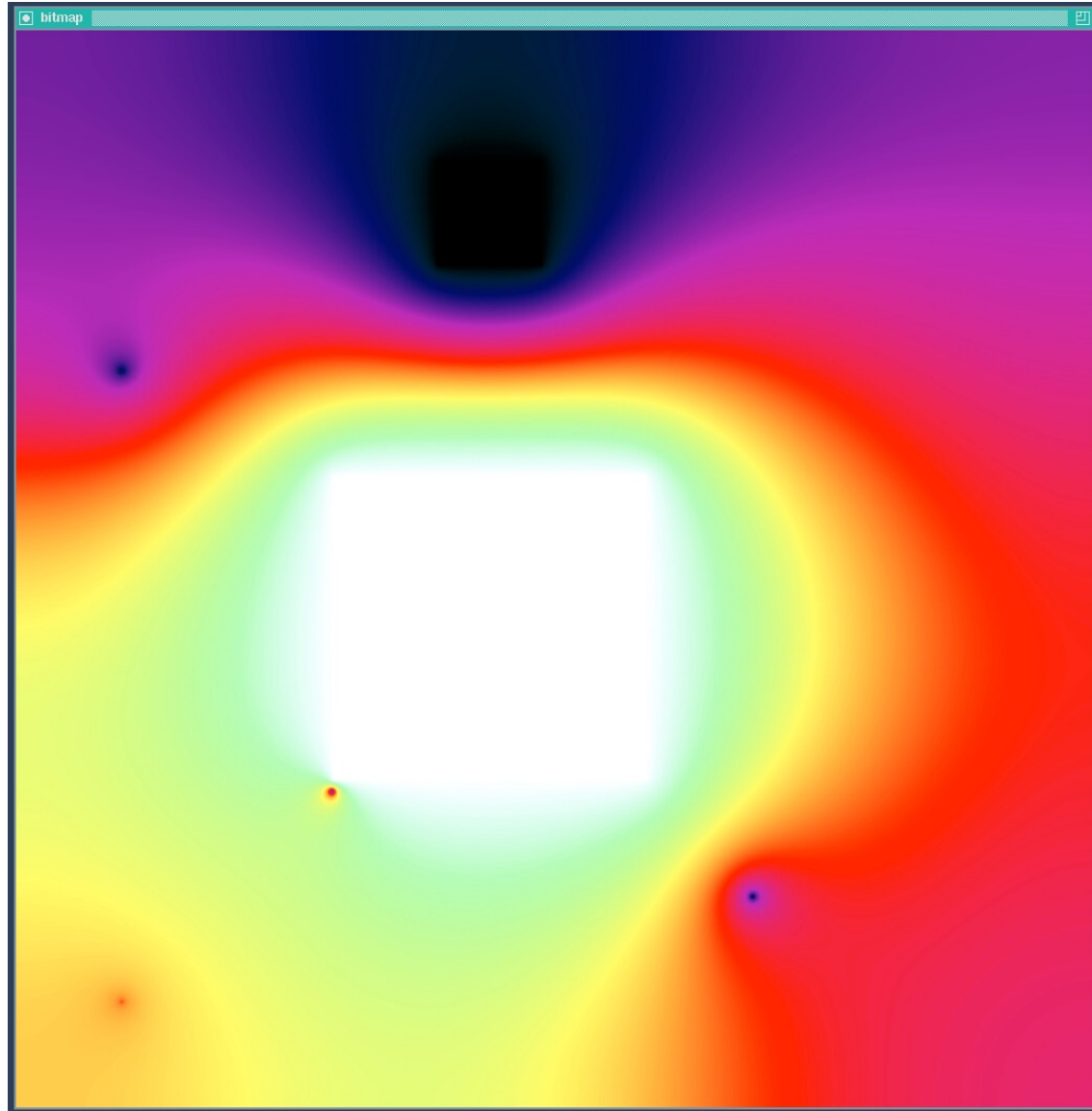
- Add a dimensionality argument of 2 to declare two-dimensional textures
 - `texture<float, 2> texConstSrc;`
 - `texture<float, 2> texIn;`
 - `texture<float, 2> texOut;`
- Change `tex1Dfetch()` calls to `tex2D()` calls
- Provide a `cudaChannelFormatDesc` when binding two-dimensional textures

```
float t, l, c, r, b;
if (dstOut) {
    t = tex2D(texIn, x, y-1);
    l = tex2D(texIn, x-1, y);
    c = tex2D(texIn, x, y);
    r = tex2D(texIn, x+1, y);
    b = tex2D(texIn, x, y+1);
} else {
    t = tex2D(texOut, x, y-1);
    l = tex2D(texOut, x-1, y);
    c = tex2D(texOut, x, y);
    r = tex2D(texOut, x+1, y);
    b = tex2D(texOut, x, y+1);
}
dst[offset] = c + SPEED * (t + b + r + l - 4 * c);
```

```
float c = tex2D(texConstSrc, x, y);
if (c != 0)
    iptr[offset] = c;
```

```
cudaChannelFormatDesc desc = cudaCreateChannelDesc<float>();
```

Heat Transfer -- Screenshot



1D Texture vs. 2D Texture?

- The version of our heat transfer simulation that uses two-dimensional textures has essentially identical performance characteristics as the version that uses one-dimensional textures
- From a performance standpoint, the decision between one- and two-dimensional textures is likely to be inconsequential
- Make the decision between one- and two-dimensional textures on a case-by-case basis

Takeaways

- Global memory is large but slow
- Shared memory is small but fast, great for applications that exhibit locality of data access
- Use constant memory for data that will not change over the course of a kernel execution to reduce the required memory bandwidth
- Use texture memory for application where memory access patterns exhibit a great deal of “spatial locality” – the use of data elements within relatively close storage locations