

- Many computational problems have a 2D domain (e.g., CV)
 - Many others have a 3D domain (e.g., fluids simulation)
- Solution: layout threads in 2D
 - Simplifies index calculations a lot

Example: Mandelbrot Set Computation

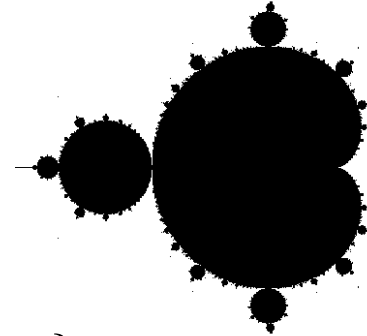
- Definition:

- For each $c \in \mathbb{C}$ consider the (infinity) sequence

$$z_{i+1} = z_i^2 + c, \quad z_0 = 0$$

- Define the Mandelbrot set

$$\mathbb{M} = \{c \in \mathbb{C} \mid \text{sequence } (z_i) \text{ remains bounded} \}$$

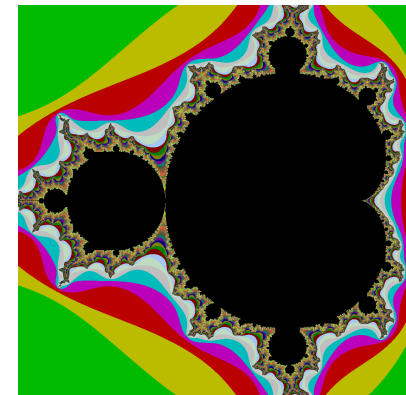


- Theorem (w/o proof):

$$\exists t : |z_t| > 2 \Rightarrow c \notin \mathbb{M}$$

- Visualizing \mathbb{M} nicely:

- Color pixel $c = (x,y)$ black, if $|z|$ remains < 2 after "many" iterations
 - Color c depending on the number of iterations necessary to make $|z_t| > 2$



- A few interesting facts about \mathbb{M}
(with which you can entertain people at a party 😊):
 - The length of the border of \mathbb{M} is *infinite*
 - \mathbb{M} is *connected*
(i.e., all "black" regions are connected with each other)
 - Mandelbrot himself believed \mathbb{M} was disconnected
 - For each color, there is exactly one "ribbon" around \mathbb{M} , i.e., there is exactly one ribbon of values c , such that $|z_1| > 2$, there is exactly one ribbon of values c , such that $|z_2| > 2$, etc. ...
 - Each such "iteration ribbon" reaches goes completely around \mathbb{M} and it is connected (i.e., there are no "self intersections")
 - There is an infinite number of "mini Mandelbrot sets", i.e., smaller copies of \mathbb{M} (self similarity)

Computing the Mandelbrot Set on the GPU

- Embarrassingly parallel: each pixel computes their own z-sequence, then sets the color
- Usual code for allocating memory, here a bitmap:

```
const unsigned int bitmap_size = img_size * img_size * 4;  
h_bitmap = new unsigned char[bitmap_size];  
cudaMalloc( (void**) &d_bitmap, bitmap_size );
```

- Setup threads layout, here a 2D arrangement of blocks

```
dim3 threads( 16, 16 );  
dim3 blocks( img_size/threads.x, img_size/threads.y );
```

- Here, we assume image size = multiple of 32
 - Simplifies calculation of number of blocks
 - Also simplifies kernel: we don't need to check whether thread out of range
 - See example code on web page how to ensure that

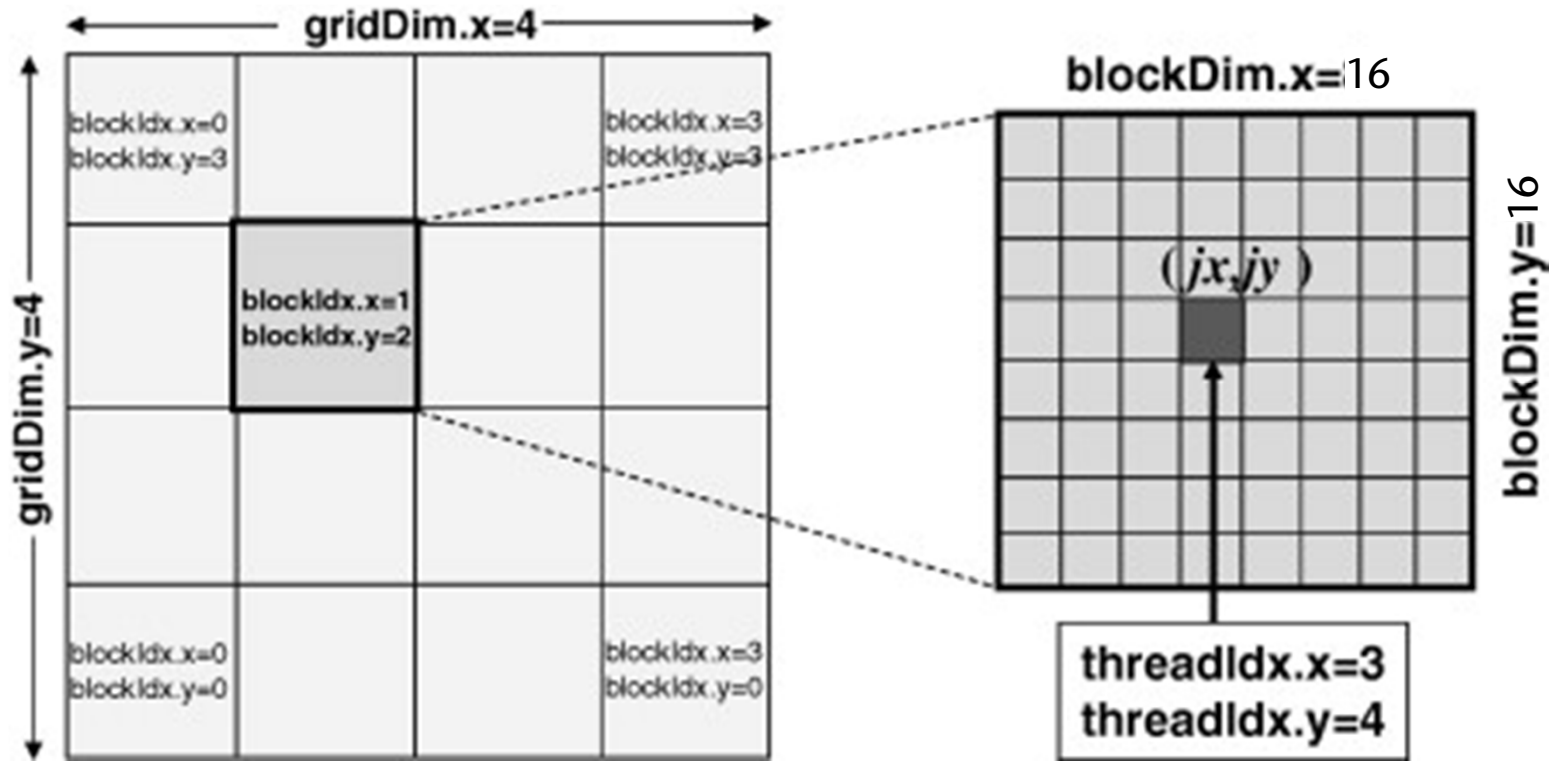
- Launch kernel:

```
mandelImage<<< blocks, threads >>>( d_bitmap, img_size );
```

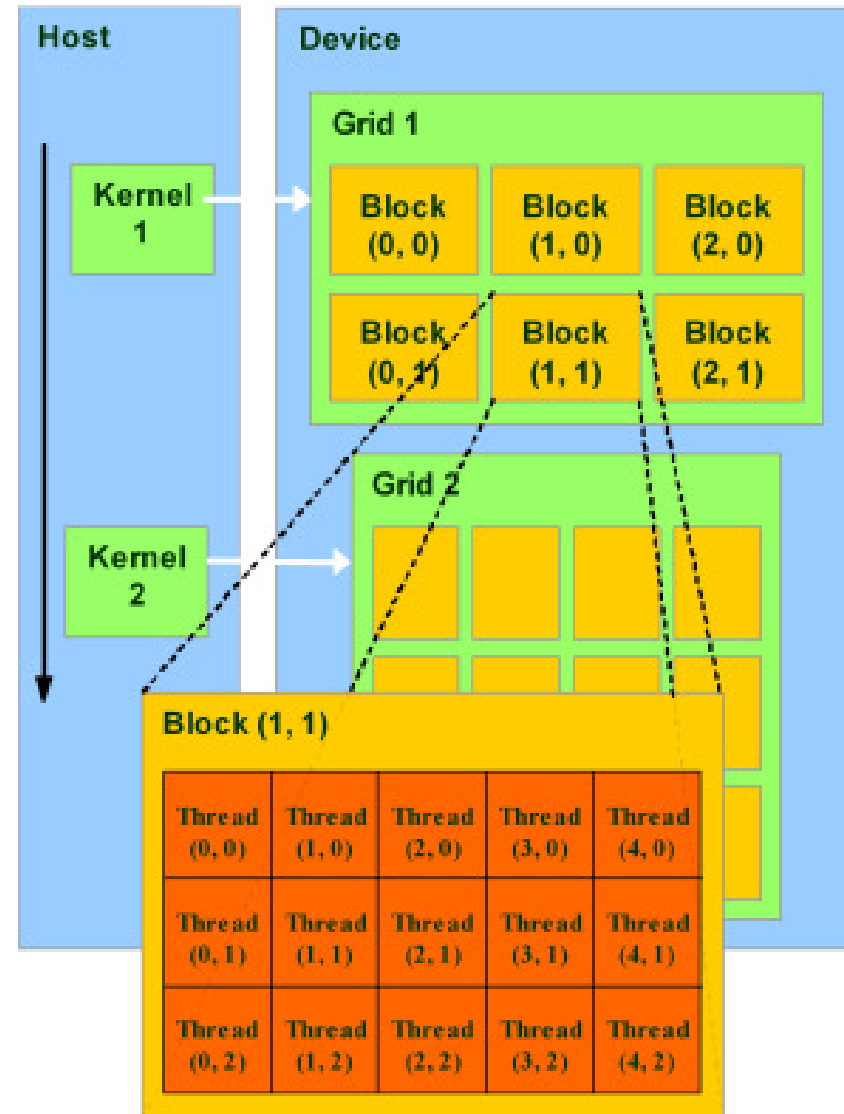
- Implementation of the kernel (simplified):

```
__global__  
void mandelImage( char4 * bitmap, const int img_size )  
{  
    int x = blockIdx.x * blockDim.x + threadIdx.x;  
    int y = blockIdx.y * blockDim.y + threadIdx.y;  
    int offset = x + y * (gridDim.x * blockDim.x); // x + y * width  
  
    int isOutsideM = isPointInMandelbrot( x, y, img_size );  
  
    bitmap[offset].x = 255 * isOutsideM; // red = outside  
    bitmap[offset].y = bitmap[offset].z = 0;  
    bitmap[offset].w = 255;  
}
```

- Visualization of our layout:



- In general, the layout of threads can change from kernel to kernel:



- Definition (done by CUDA):

```
struct dim3    // is actually a C++ class
{
    unsigned int x, y, z;
};
```

- Usage:

```
dim3 layout (nx);    = dim3 layout (nx, 1);    = dim3 layout (nx, 1, 1);
```

```
dim3 layout (nx, ny);    = dim3 layout (nx, ny, 1);
```

- Launching a kernel like this: `kernel<<<N,M>>> (...);`

is equivalent:

```
dim3 threads (M, 1);
dim3 blocks (N, 1);
kernel<<<blocks, threads>>> (...);
```


Implementation of the Kernel

```
__device__
int isPointInMandelbrot( int x, int y,
                        const int img_size, float scale )
{
    cuComplex c( (float)(x - img_size/2) / (img_size/2),
                (float)(y - img_size/2) / (img_size/2) );
    c *= scale;
    cuComplex z( 0.0, 0.0 );           // z_i of the sequence

    for ( int i = 0; i < 200; i ++ )
    {
        z = z*z + c;
        if ( z.magnitude2() > 4.0f ) // |z|^2 > 2^2 -> outside
            return i;
    }

    return 0;
}
```

```
struct cuComplex    // define a class for complex numbers
{
    float r, i;      // real and imaginary part

    __device__      // constructor
    cuComplex( float a, float b ) : r(a), i(b) {}

    __device__      // |z|^2
    float magnitude2( void )
    {
        return r * r + i * i;
    }

    __device__      // z1 * z2
    cuComplex operator * (const cuComplex & a)
    {
        return cuComplex(r*a.r - i*a.i, i*a.r + r*a.i);
    }
    // for more: see example code on web page
};
```

	Executed on:	Only callable from:
<code>__device__ float DeviceFunc();</code>	device	device
<code>__global__ void KernelFunc();</code>	device	host
<code>__host__ float HostFunc();</code>	host	host

- Remarks:
 - `__global__` defines a kernel function
 - Each `'__'` consists of two underscore characters
 - A kernel function must return `void`
 - `__device__` and `__host__` can be used together

- Example for the latter: make **cuComplex** usable on both device and host

```
struct cuComplex    // define a class for complex numbers
{
    float r, i;      // real, imaginary part

    __device__ __host__
    cuComplex( float a, float b ) : r(a), i(b) {}

    __device__ __host__
    float magnitude2( void )
    {
        return r * r + i * i;
    }
    // etc. ...
};
```

- An "Optimization":

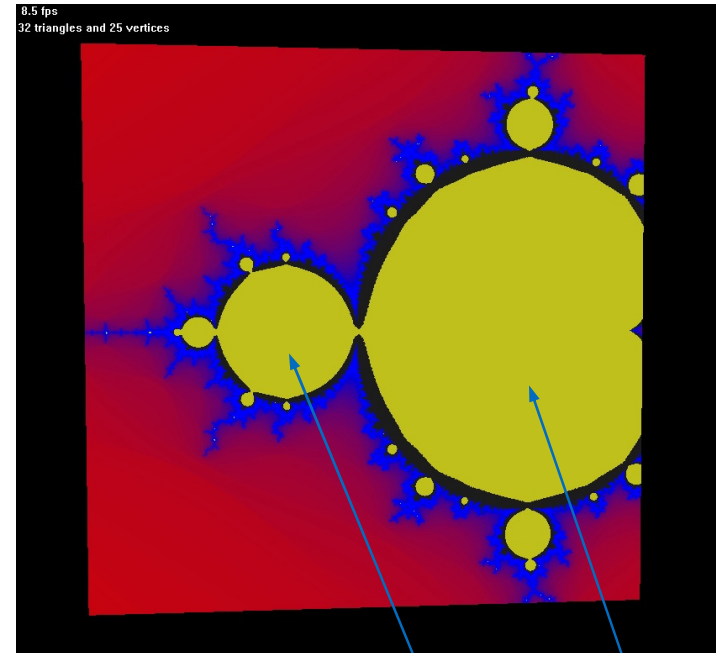
- The sequence of z_i can either converge towards single (complex) value,
 - or it can end up in a cycle of values,
 - or it can be chaotic.

- Idea:

- Try to recognize such cycles; if you realize that a thread's is caught in a cycle, exit immediately (should happen much earlier in most cases)

- Maintain an array of the k most recent elements of the sequence

- Last time I checked: 4x slower than the brute-force version!



All points here converge towards cycle of length 2

All points here converge towards fixed point

Querying the Device for its Capabilities

- How do you know how many threads can be in a block, etc.?
- Query your GPU, like so:

```
int devID;
cudaGetDevice( &devID );           // GPU currently in use
cudaDeviceProp props;
cudaGetDeviceProperties( &props, devID );

unsigned int threads_per_block = props.maxThreadsPerBlock;
```

For Your Reference: the Complete Table of the cudaDeviceProp

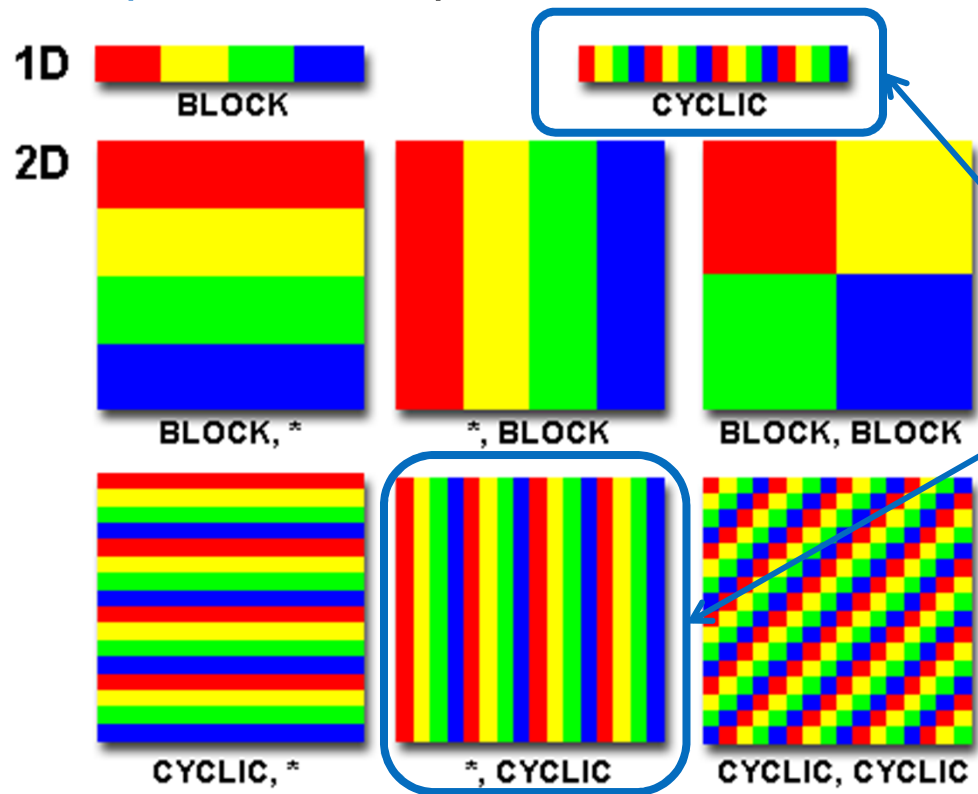
DEVICE PROPERTY	DESCRIPTION
<code>char name[256];</code>	An ASCII string identifying the device (e.g., "GeForce GTX 280")
<code>size_t totalGlobalMem</code>	The amount of global memory on the device in bytes
<code>size_t sharedMemPerBlock</code>	The maximum amount of shared memory a single block may use in bytes
<code>int regsPerBlock</code>	The number of 32-bit registers available per block
<code>int warpSize</code>	The number of threads in a warp
<code>size_t memPitch</code>	The maximum pitch allowed for memory copies in bytes
<code>int maxThreadsPerBlock</code>	The maximum number of threads that a block may contain
<code>int maxThreadsDim[3]</code>	The maximum number of threads allowed along each dimension of a block
<code>int maxGridSize[3]</code>	The number of blocks allowed along each dimension of a grid
<code>size_t totalConstMem</code>	The amount of available constant memory

DEVICE PROPERTY	DESCRIPTION
<code>int</code> major	The major revision of the device's compute capability
<code>int</code> minor	The minor revision of the device's compute capability
<code>size_t</code> textureAlignment	The device's requirement for texture alignment
<code>int</code> deviceOverlap	A boolean value representing whether the device can simultaneously perform a <code>cudaMemcpy()</code> and kernel execution
<code>int</code> multiProcessorCount	The number of multiprocessors on the device
<code>int</code> kernelExecTimeoutEnabled	A boolean value representing whether there is a runtime limit for kernels executed on this device
<code>int</code> integrated	A boolean value representing whether the device is an integrated GPU (i.e., part of the chipset and not a discrete GPU)
<code>int</code> canMapHostMemory	A boolean value representing whether the device can map host memory into the CUDA device address space
<code>int</code> computeMode	A value representing the device's computing mode: default, exclusive, or prohibited
<code>int</code> maxTexture1D	The maximum size supported for 1D textures

DEVICE PROPERTY	DESCRIPTION
<code>int</code> <code>maxTexture2D</code> [2]	The maximum dimensions supported for 2D textures
<code>int</code> <code>maxTexture3D</code> [3]	The maximum dimensions supported for 3D textures
<code>int</code> <code>maxTexture2DArray</code> [3]	The maximum dimensions supported for 2D texture arrays
<code>int</code> <code>concurrentKernels</code>	A boolean value representing whether the device supports executing multiple kernels within the same context simultaneously

Problem Partitioning

- Problem: your input, e.g. the vectors, is larger than the maximally allowed size along one dimension?
 - I.e., what if `vec_len > maxThreadsDim[0] * maxGridSize[0]`?
- Solution: **partition** the problem (color = thread ID)



For most kind of applications these two partitionings are best for GPUs!

Example: Adding Huge Vectors

- Vectors of size 100,000,000 are not uncommon in high-performance computing (HPC) ...
- The thread layout:

```
dim3 threads(16,16);    // = 256 threads per block
int  n_threads_pb = threads.x * threads.y;
int  n_blocks = (vec_len + n_threads_pb - 1) / n_threads_pb;
int  nb_sqrt = (int)( ceilf( sqrtf( n_blocks ) ) );
dim3 blocks( nb_sqrt, nb_sqrt );
```

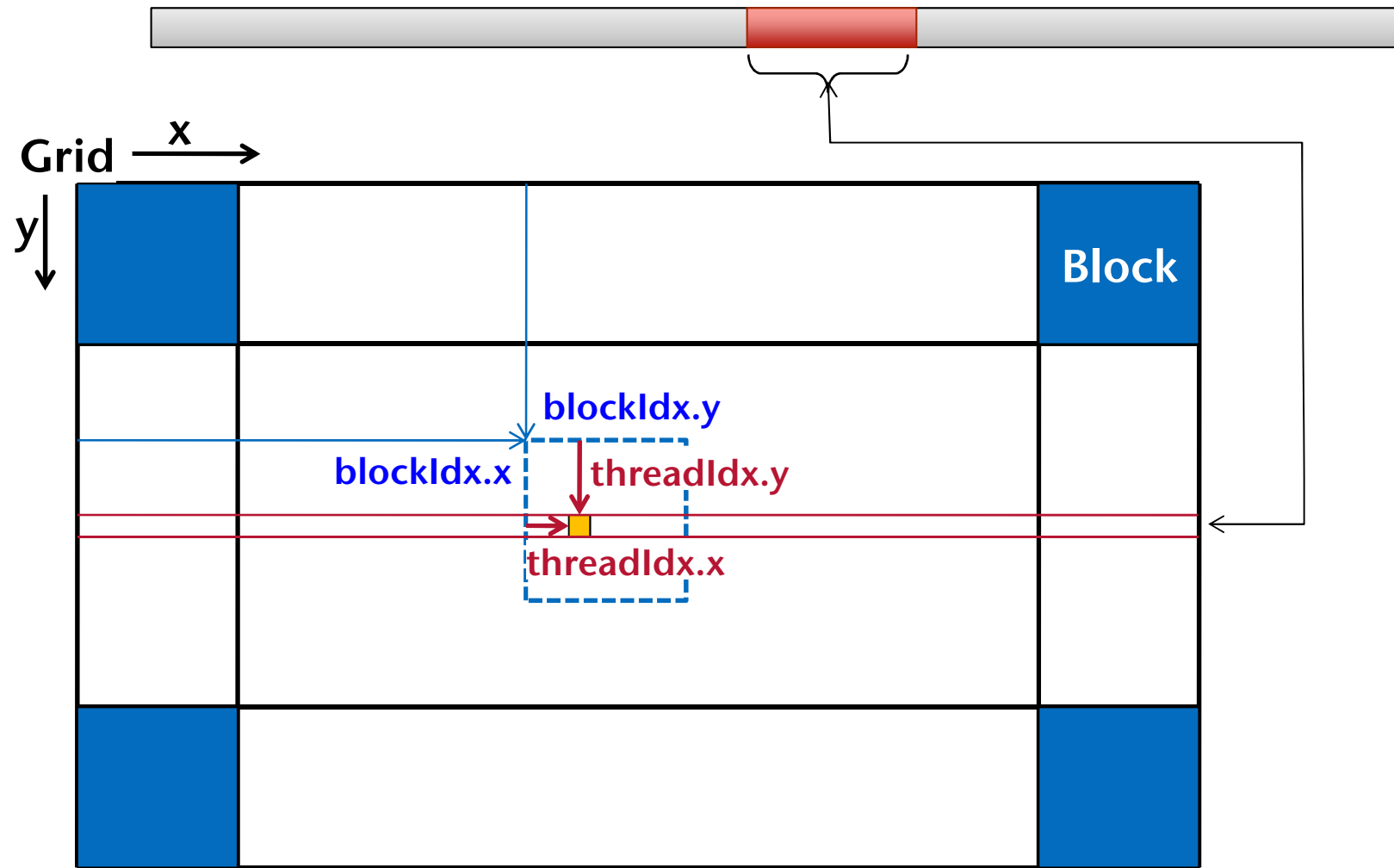
- Kernel launch:

```
addVectors<<< threads, blocks >>>( d_a, d_b, d_c, n );
```

- Index computation in the kernel:

```
unsigned int tid_x = blockDim.x * blockIdx.x + threadIdx.x;
unsigned int tid_y = blockDim.y * blockIdx.y + threadIdx.y;
unsigned int i     = tid_y * (blockDim.x * gridDim.x) + tid_x;
```

- Visualization of this index computation:

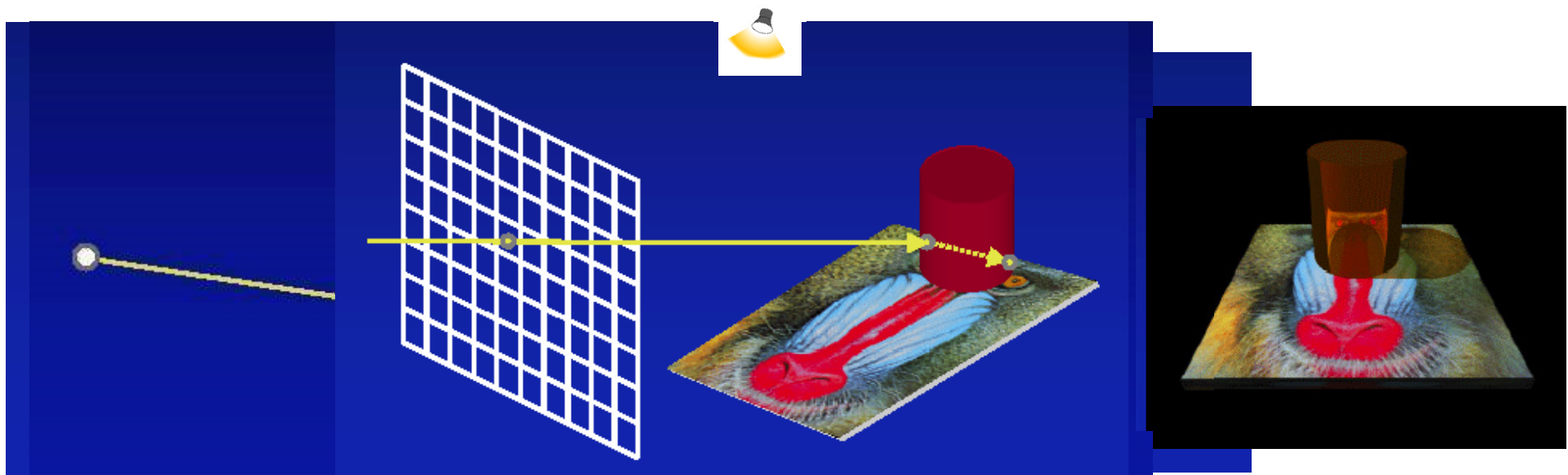


Constant Memory

- Why is it so important to declare constant variables/instances in C/C++ as **const** ?
- It allows the compiler to ...
 - optimize your program a lot
 - do more type-checking
- Something similar exists in CUDA → constant memory

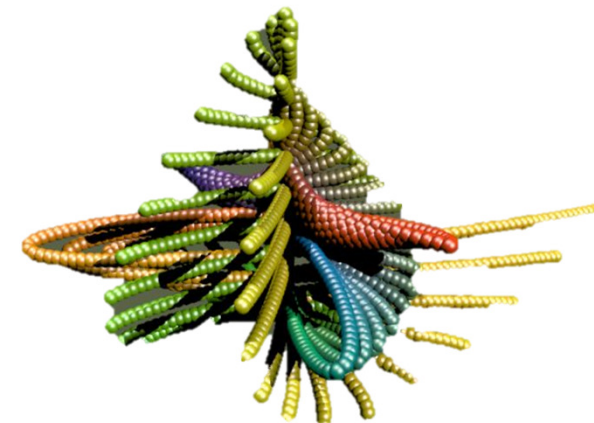
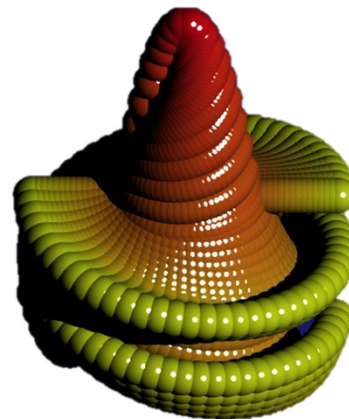
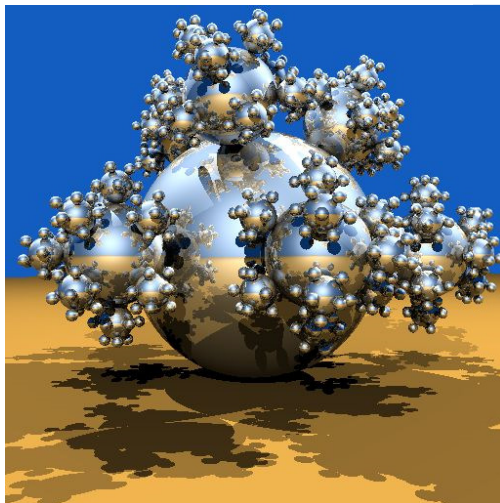
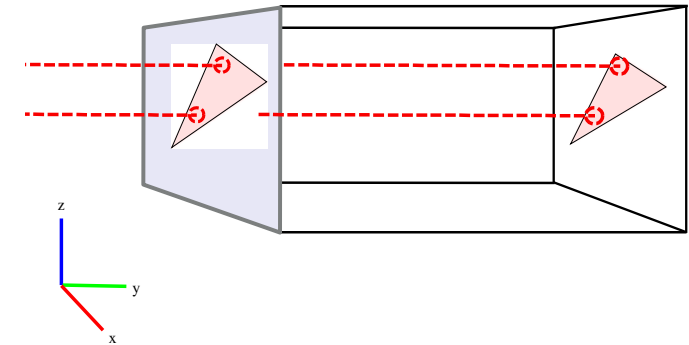
Example: a Simple Raytracer

- The ray-tracing principle:
 1. Shoot rays from camera through every pixel into scene (**primary rays**)
 2. If the rays hits more than one object, then consider only the first hit
 3. From there, shoot rays to all light sources (**shadow feelers**)
 4. If a shadow feeler hits another obj → point is in shadow w.r.t. that light source.
Otherwise, evaluate a lighting model (e.g., Phong [see "Computer graphics"])
 5. If the hit object is glossy, then shoot reflected rays into scene (**secondary rays**) → recursion
 6. If the hit object is transparent, then shoot refracted ray → more recursion



■ Simplifications (for now):

- Only primary rays
- Camera at infinity → primary rays are orthogonal to image plane
- Only spheres
 - They are so easy, every raytracer has them ☺



- The data structures:

```
struct Sphere
{
    Vec3 center;           // center of sphere
    float radius;
    Color r, g, b;        // color of sphere

    __device__
    bool intersect( const Ray & ray, Hit * hit )
    {
        ...
    }
};
```


- The mechanics on the host side:

```
int main( void )
{
    // create host/device bitmaps (see Mandelbrot ex.)
    ...
    Sphere * h_spheres = new Sphere[n_spheres];
    // generate spheres, or read from file

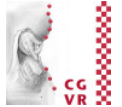
    // transfer spheres to device (later)

    // generate image by launching kernel
    // assumption: img_size = multiple of 16!
    dim3 threads(16,16);
    dim3 blocks( img_size/16, img_size/16 );
    raytrace<<<blocks,threads>>>( d_bitmap );

    // display, clean up, and exit
};
```



The mechanics on the device side



```

__global__
void raytrace( unsigned char * bitmap ) {
    // map from thread id to pixel pos
    int x = blockIdx.x * blockDim.x + threadIdx.x;
    int y = blockIdx.y * blockDim.y + threadIdx.y;
    int offset = x + y * (gridDim.x * blockDim.x);

    Ray primary( x, y, camera );           // generate primary ray

    // check intersection with scene, take closest one
    min_dist = INF;
    int hit_sphere = MAX_INT;
    Hit hit;
    for ( int i = 0; i < n_spheres; i ++ ) {
        if ( intersect(ray, i, & hit) ) {
            if ( hit.dist < min_dist ) {
                min_dist = hit.dist;       // found a closer hit
                hit_sphere = i;           // remember sphere; hit info
            }                               // is already filled
        }
    }
    // compute color at hit point (if any) and set in bitmap[offset]
}

```

Declaration & transfer

- Since it's constant memory, we declare it as such:

```
const int MAX_NUM_SPHERES 1000;
__constant__ Sphere c_spheres[MAX_NUM_SPHERES];
```

- Transfer now works by a different kind of Memcpy:

```
int main( void )
{
    ...
    // transfer spheres to device (later)
    cudaMemcpyToSymbol( c_spheres, h_spheres,
                       n_spheres * sizeof(Sphere) );
    ...
};
```

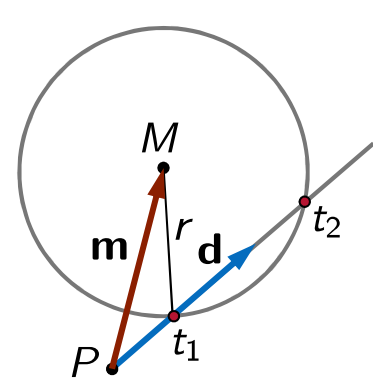
- Access of constant memory on the device (i.e., from a kernel) works just like with any globally declared variable
- Example:

```

__constant__ Sphere c_spheres[MAX_NUM_SPHERES];

__device__
bool intersect( const Ray & ray, int s, Hit * hit )
{
    Vec3 m( c_spheres[s].center - ray.orig );
    float q = m*m - c_spheres[s].radius*c_spheres[s].radius;
    float p = ...
    solve_pq( p, q, *t1, *t2 );
    ...
}

```



$$(t \cdot d - m)^2 = r^2 \Rightarrow t^2 - 2t \cdot md + m^2 - r^2 = 0$$