

# **CS 380 - GPU and GPGPU Programming**

## **Lecture 15: GPU Compute APIs, Pt. 4; GPU Texturing, Pt. 1**

Markus Hadwiger, KAUST

# Reading Assignment #9 (until Oct 30)



## Read (required):

- Interpolation for Polygon Texture Mapping and Shading,  
Paul Heckbert and Henry Moreton

<https://www.rh.cmu.edu/publications/interpolation-for-polygon-texture-mapping-and-shading/>

- Homogeneous Coordinates

[https://en.wikipedia.org/wiki/Homogeneous\\_coordinates](https://en.wikipedia.org/wiki/Homogeneous_coordinates)

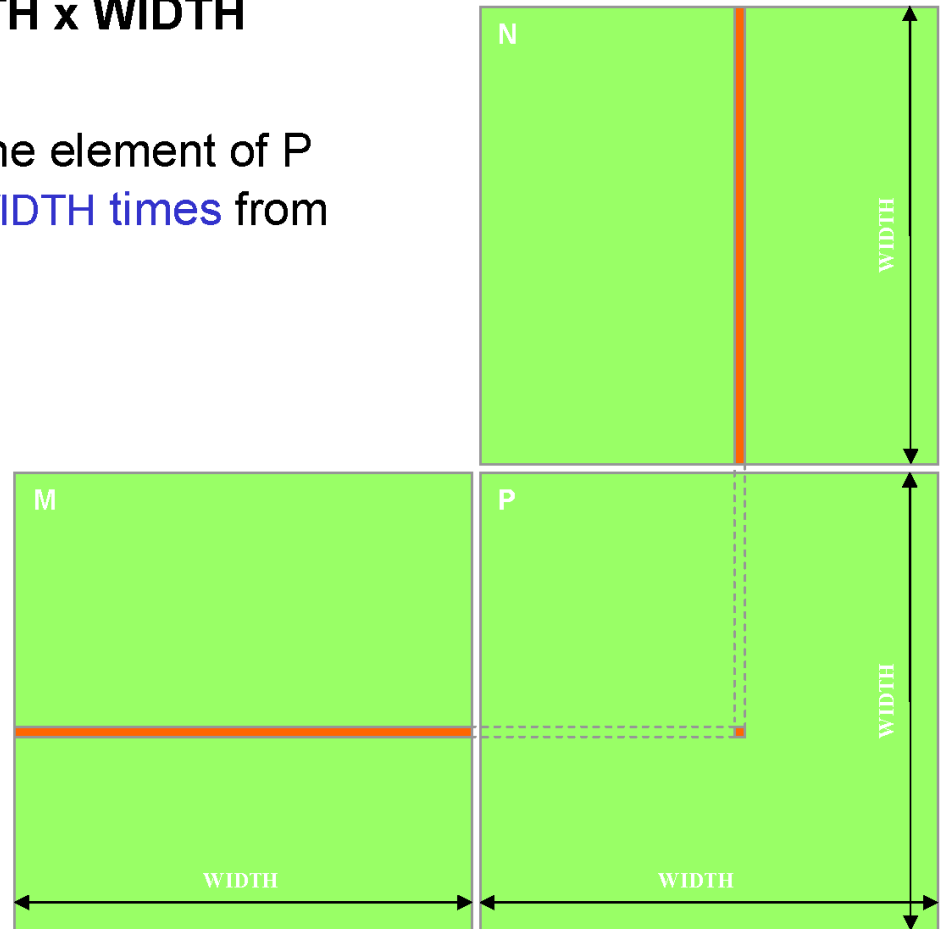
# Code Examples

## Example #2: Matrix Multiply

# Programming Model: Square Matrix Multiplication

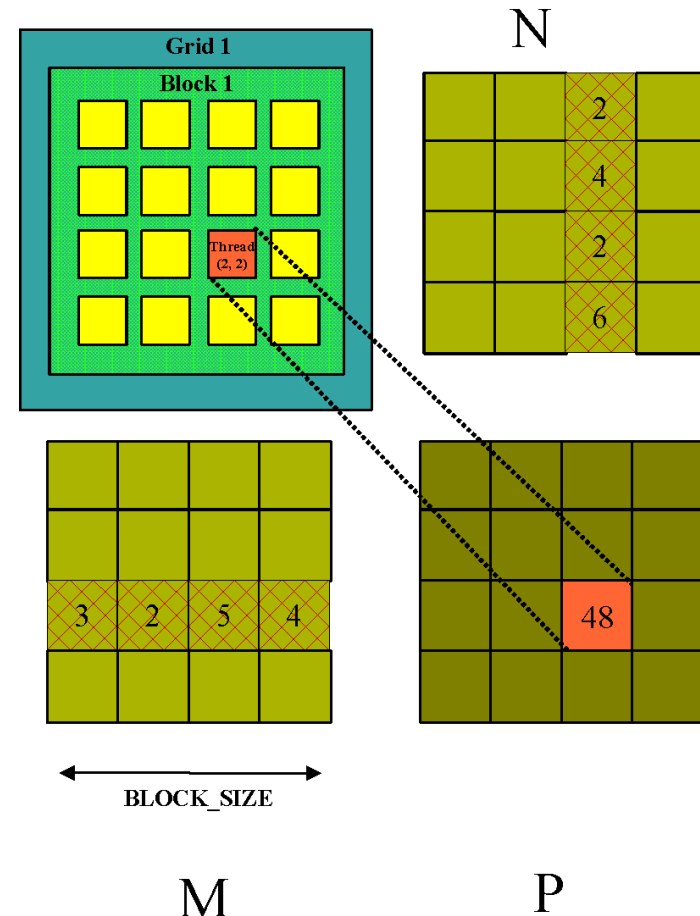
---

- $P = M * N$  of size  $WIDTH \times WIDTH$
- **Without tiling:**
  - One **thread** handles one element of  $P$
  - $M$  and  $N$  are loaded  $WIDTH$  times from global memory



# Multiply Using One Thread Block

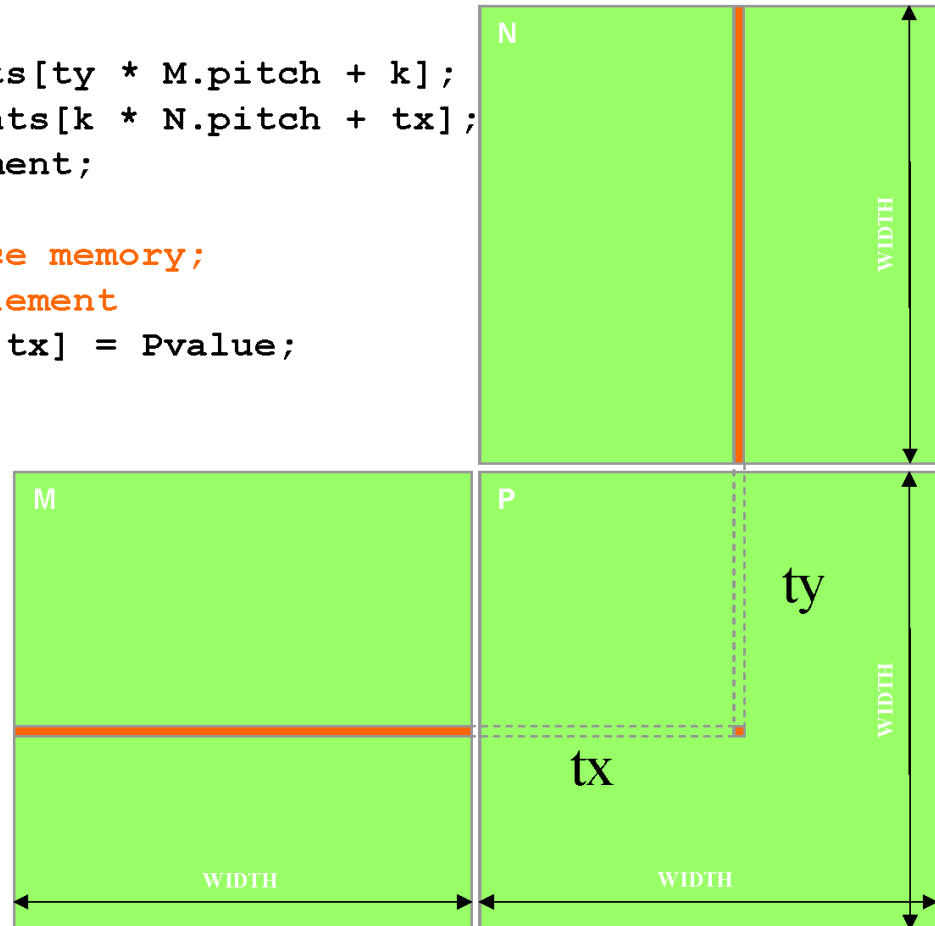
- **One block of threads computes matrix P**
  - Each thread computes one element of P
- **Each thread**
  - Loads a row of matrix M
  - Loads a column of matrix N
  - Perform one multiply and addition for each pair of M and N elements
  - Compute to off-chip memory access ratio close to 1:1 (not very high)
- **Size of matrix limited by the number of threads allowed in a thread block**



# Matrix Multiplication Device-Side Kernel Function (cont.)

...

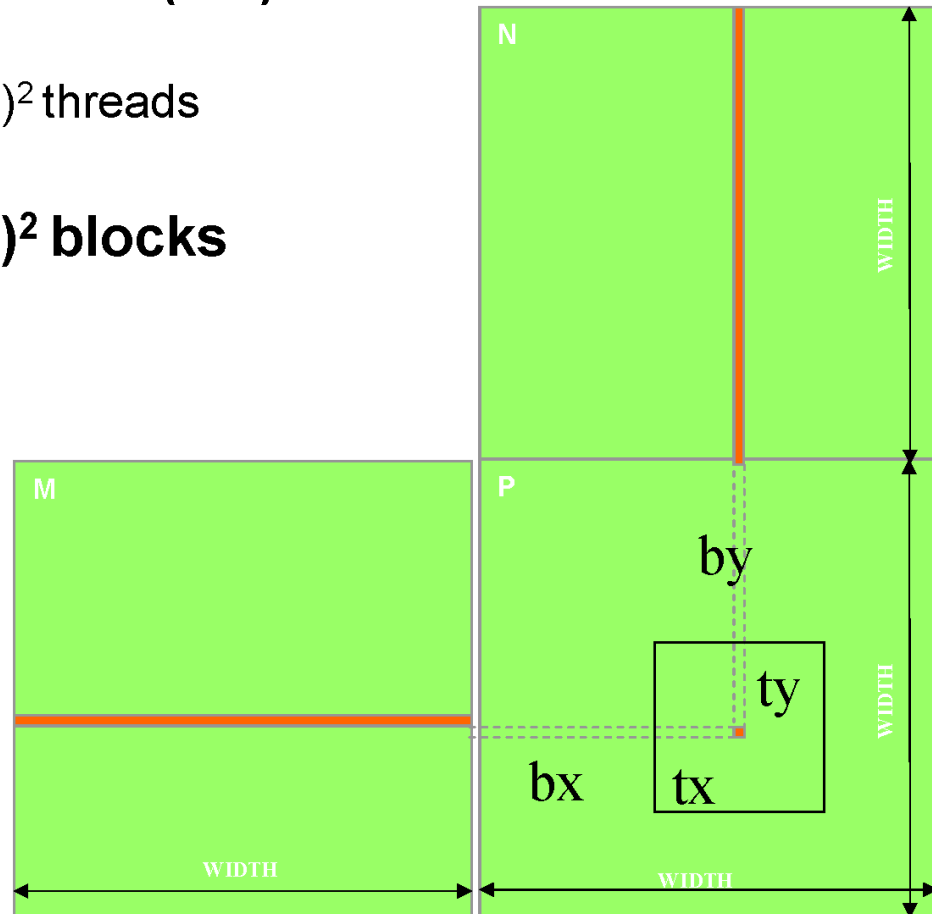
```
for (int k = 0; k < M.width; ++k)
{
    float Melement = M.elements[ty * M.pitch + k];
    float Nelement = Nd.elements[k * N.pitch + tx];
    Pvalue += Melement * Nelement;
}
// Write the matrix to device memory;
// each thread writes one element
P.elements[ty * blockDim.x + tx] = Pvalue;
}
```



# Handling Arbitrary Sized Square Matrices

- Have each 2D thread block to compute a  $(\text{BLOCK\_WIDTH})^2$  sub-matrix (tile) of the result matrix
  - Each has  $(\text{BLOCK\_WIDTH})^2$  threads
- Generate a 2D Grid of  $(\text{WIDTH}/\text{BLOCK\_WIDTH})^2$  blocks

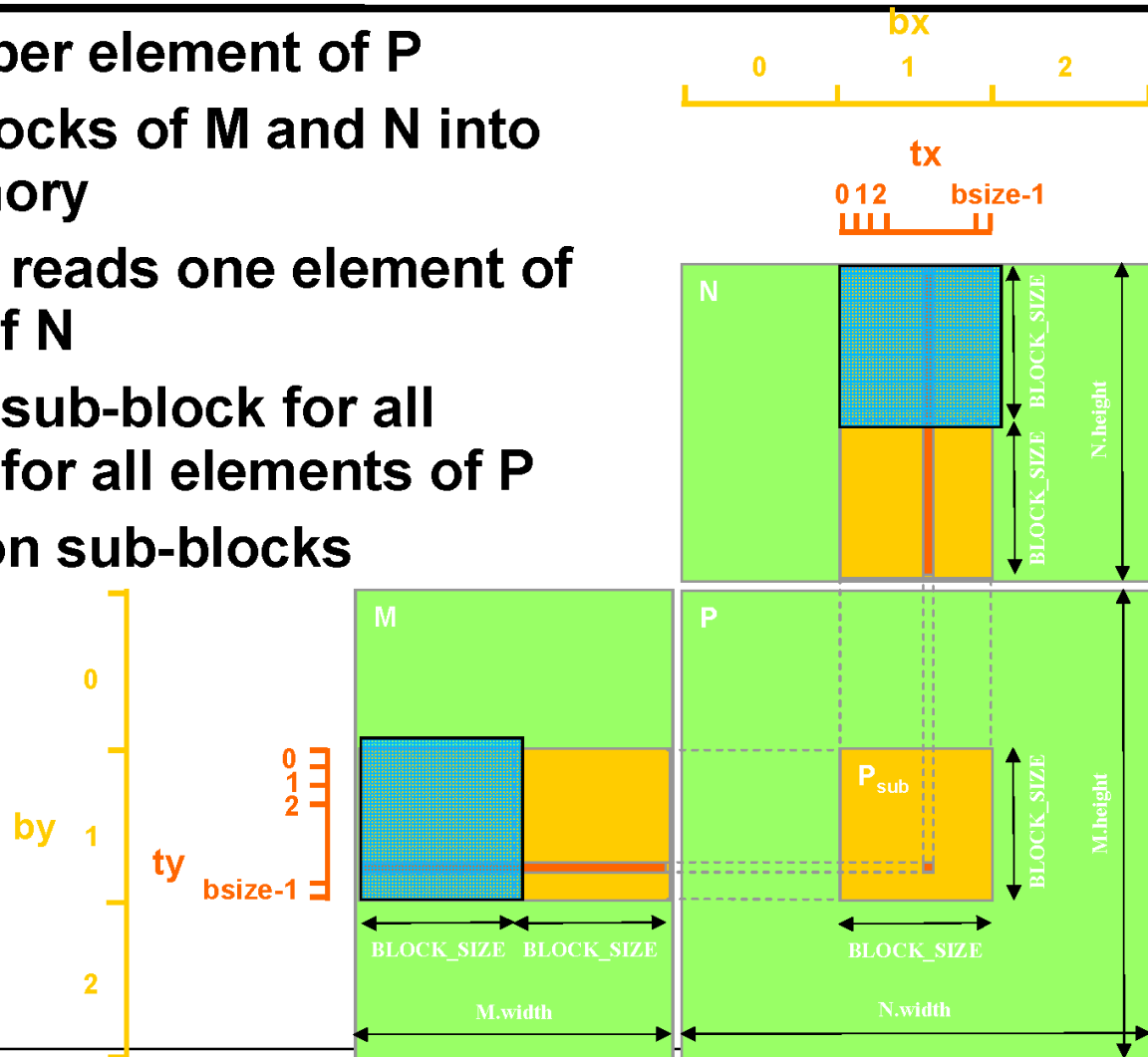
You still need to put a loop around the kernel call for cases where  $\text{WIDTH}$  is greater than Max grid size!





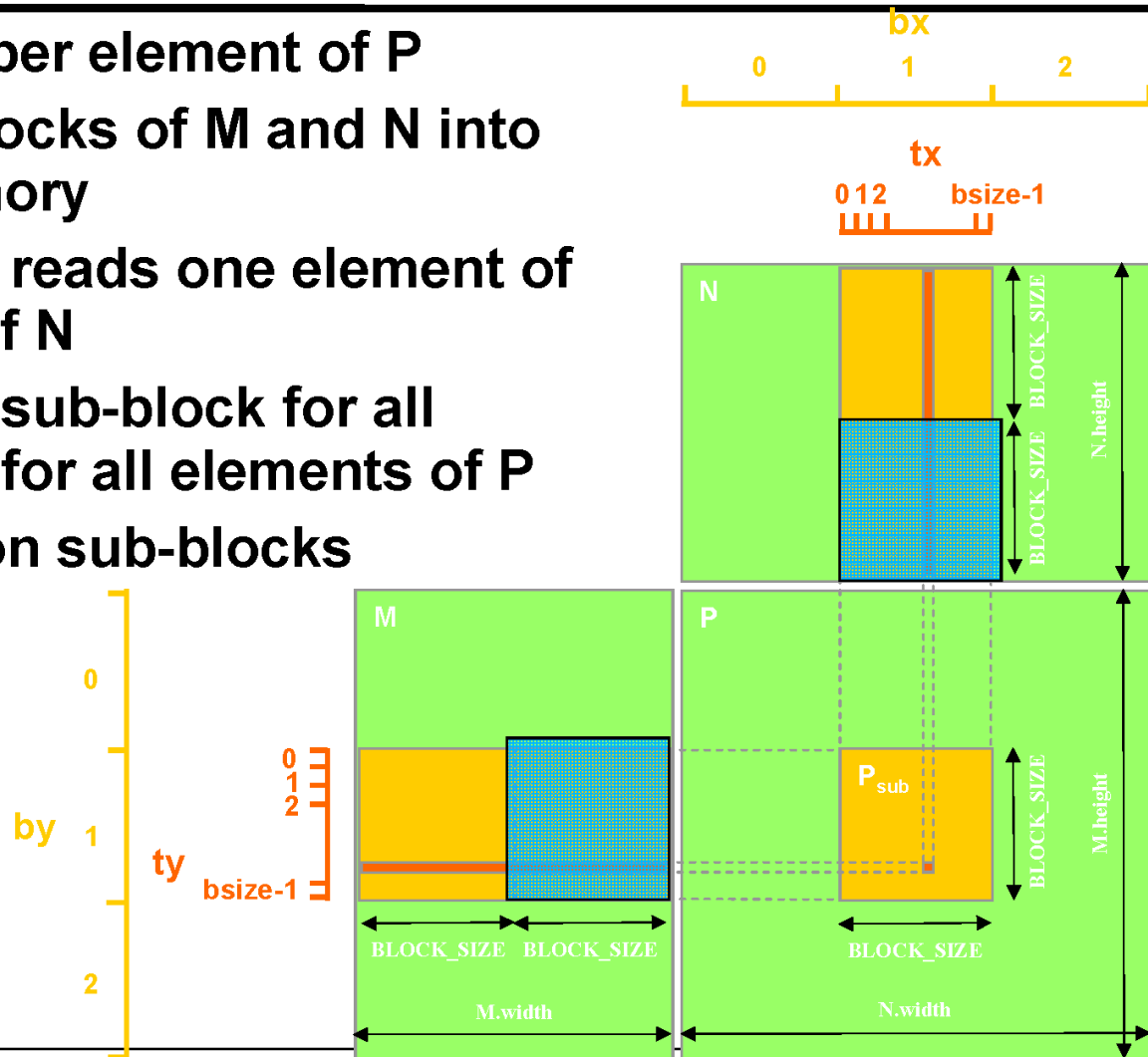
# Multiply Using Several Blocks - Idea

- One thread per element of P
- Load sub-blocks of M and N into shared memory
- Each thread reads one element of M and one of N
- Reuse each sub-block for all threads, i.e. for all elements of P
- Outer loop on sub-blocks



# Multiply Using Several Blocks - Idea

- One thread per element of P
- Load sub-blocks of M and N into shared memory
- Each thread reads one element of M and one of N
- Reuse each sub-block for all threads, i.e. for all elements of P
- Outer loop on sub-blocks



# Example: Matrix Multiplication (1)



- Copy matrices to device; invoke kernel; copy result matrix back to host

```
// Matrix multiplication - Host code
// Matrix dimensions are assumed to be multiples of BLOCK_SIZE
void MatMul(const Matrix A, const Matrix B, Matrix C)
{
    // Load A and B to device memory
    Matrix d_A;
    d_A.width = d_A.stride = A.width; d_A.height = A.height;
    size_t size = A.width * A.height * sizeof(float);
    cudaMalloc((void**)&d_A.elements, size);
    cudaMemcpy(d_A.elements, A.elements, size,
               cudaMemcpyHostToDevice);

    Matrix d_B;
    d_B.width = d_B.stride = B.width; d_B.height = B.height;
    size = B.width * B.height * sizeof(float);
    cudaMalloc((void**)&d_B.elements, size);
    cudaMemcpy(d_B.elements, B.elements, size,
               cudaMemcpyHostToDevice);
}
```

## Example: Matrix Multiplication (2)



```
// Allocate C in device memory
Matrix d_C;
d_C.width = d_C.stride = C.width; d_C.height = C.height;
size = C.width * C.height * sizeof(float);
cudaMalloc((void**)&d_C.elements, size);

// Invoke kernel
dim3 dimBlock(BLOCK_SIZE, BLOCK_SIZE);
dim3 dimGrid(B.width / dimBlock.x, A.height / dimBlock.y);
MatMulKernel<<<dimGrid, dimBlock>>>(d_A, d_B, d_C);

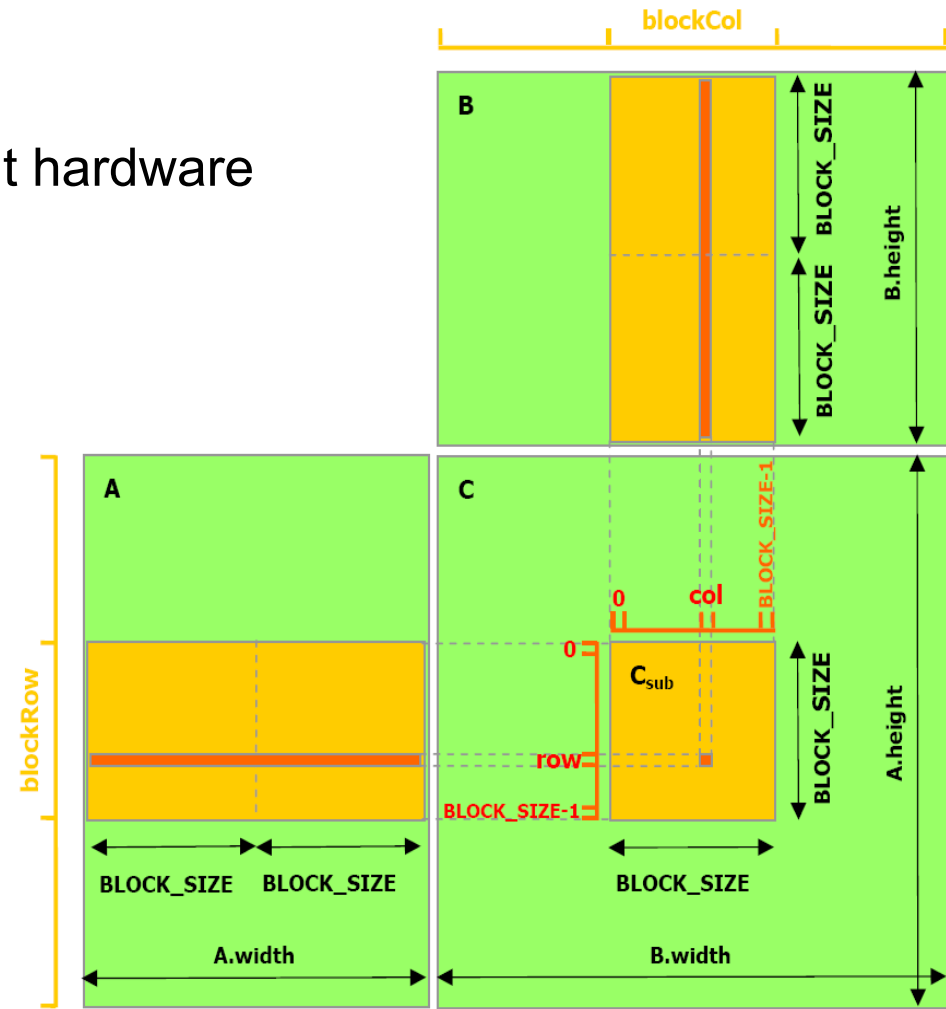
// Read C from device memory
cudaMemcpy(C.elements, d_C.elements, size,
           cudaMemcpyDeviceToHost);

// Free device memory
cudaFree(d_A.elements);
cudaFree(d_B.elements);
cudaFree(d_C.elements);
}
```

# Example: Matrix Multiplication (3)



- Multiply matrix block-wise
- Set BLOCK\_SIZE for efficient hardware use, e.g., to 16 on cc. 1.x or 16 or 32 on cc. 2.x +
- Maximize parallelism
  - Launch as many threads per block as block elements
  - Each thread fetches one element of block
  - Perform row \* column dot products in parallel



# Example: Matrix Multiplication (4)



```
__global__ void MatrixMul( float *matA, float *matB, float *matC, int w )
{
    __shared__ float blockA[ BLOCK_SIZE ][ BLOCK_SIZE ];
    __shared__ float blockB[ BLOCK_SIZE ][ BLOCK_SIZE ];

    int bx = blockIdx.x; int tx = threadIdx.x;
    int by = blockIdx.y; int ty = threadIdx.y;

    int col = bx * BLOCK_SIZE + tx;
    int row = by * BLOCK_SIZE + ty;

    float out = 0.0f;
    for ( int m = 0; m < w / BLOCK_SIZE; m++ ) {

        blockA[ ty ][ tx ] = matA[ row * w + m * BLOCK_SIZE + tx ];
        blockB[ ty ][ tx ] = matB[ col + ( m * BLOCK_SIZE + ty ) * w ];
        __syncthreads();

        for ( int k = 0; k < BLOCK_SIZE; k++ ) {
            out += blockA[ ty ][ k ] * blockB[ k ][ tx ];
        }
        __syncthreads();
    }

    matC[ row * w + col ] = out;
}
```

Caveat: for brevity, this code assumes matrix sizes are a multiple of the block size (either because they really are, or because padding is used; otherwise guard code would need to be added)

**What About Memory Performance?  
(more to come later...)**

# Memory Layout of a Matrix in C

$M_{0,0}$	$M_{1,0}$	$M_{2,0}$	$M_{3,0}$
$M_{0,1}$	$M_{1,1}$	$M_{2,1}$	$M_{3,1}$
$M_{0,2}$	$M_{1,2}$	$M_{2,2}$	$M_{3,2}$
$M_{0,3}$	$M_{1,3}$	$M_{2,3}$	$M_{3,3}$

M

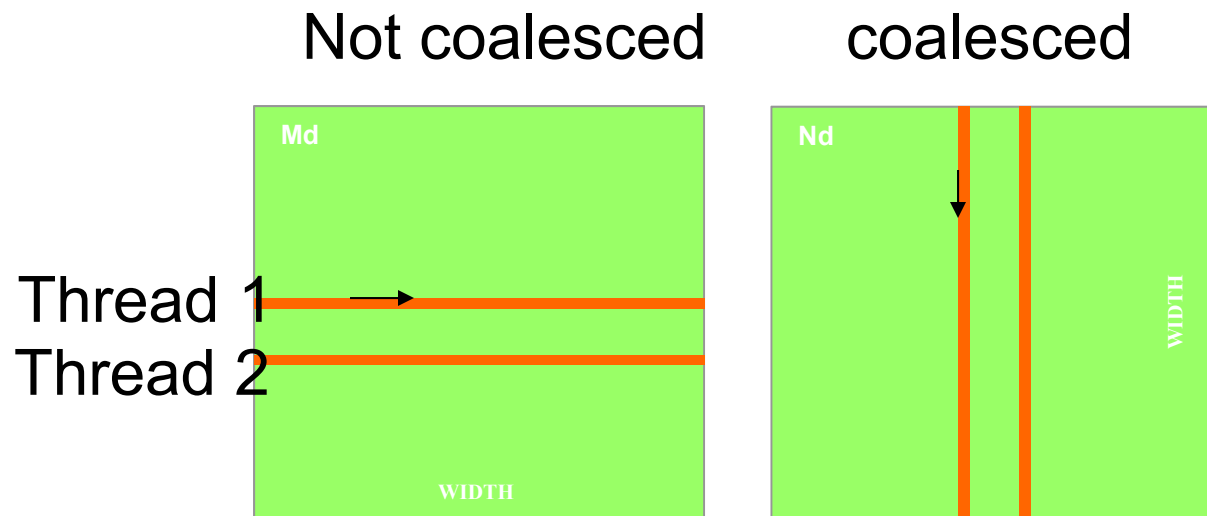


$M_{0,0}$	$M_{1,0}$	$M_{2,0}$	$M_{3,0}$	$M_{0,1}$	$M_{1,1}$	$M_{2,1}$	$M_{3,1}$	$M_{0,2}$	$M_{1,2}$	$M_{2,2}$	$M_{3,2}$	$M_{0,3}$	$M_{1,3}$	$M_{2,3}$	$M_{3,3}$
-----------	-----------	-----------	-----------	-----------	-----------	-----------	-----------	-----------	-----------	-----------	-----------	-----------	-----------	-----------	-----------



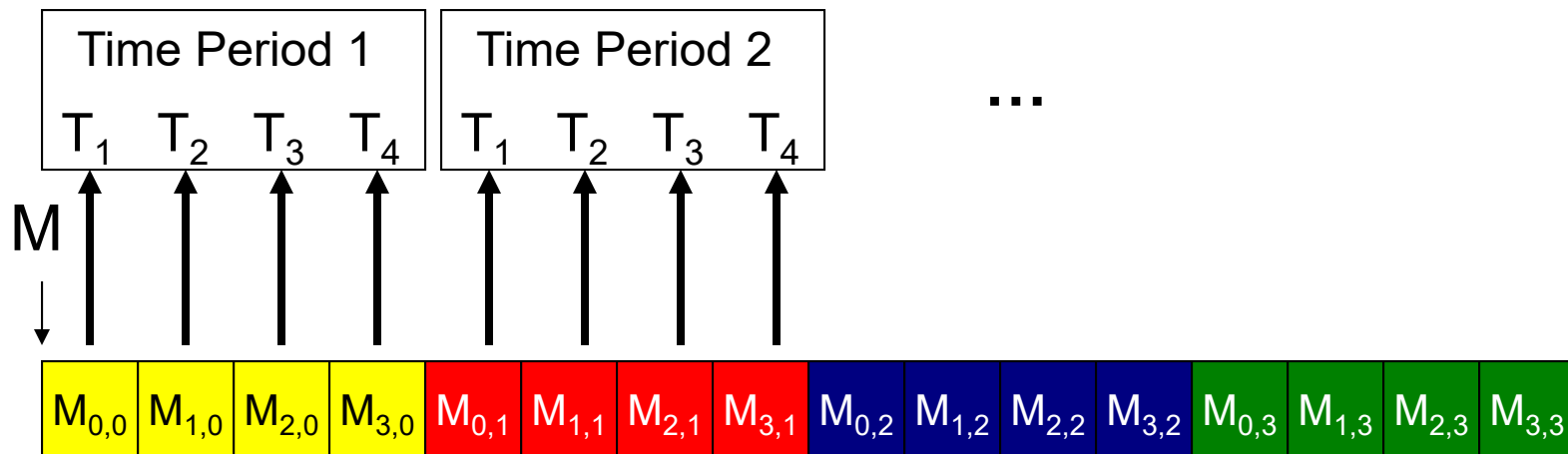
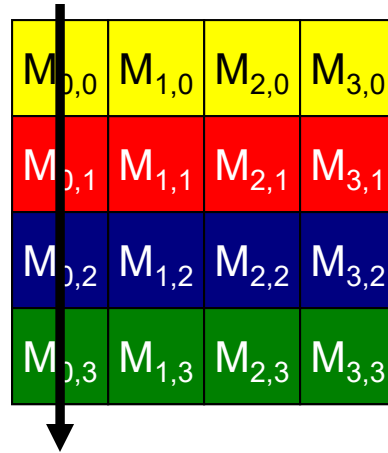
# Memory Coalescing

- When accessing global memory, peak performance utilization occurs when all threads in a half warp (**full warp on Fermi**) access continuous memory locations.
- Requirements relaxed on  $\geq 1.2$  devices; L1 cache on Fermi!

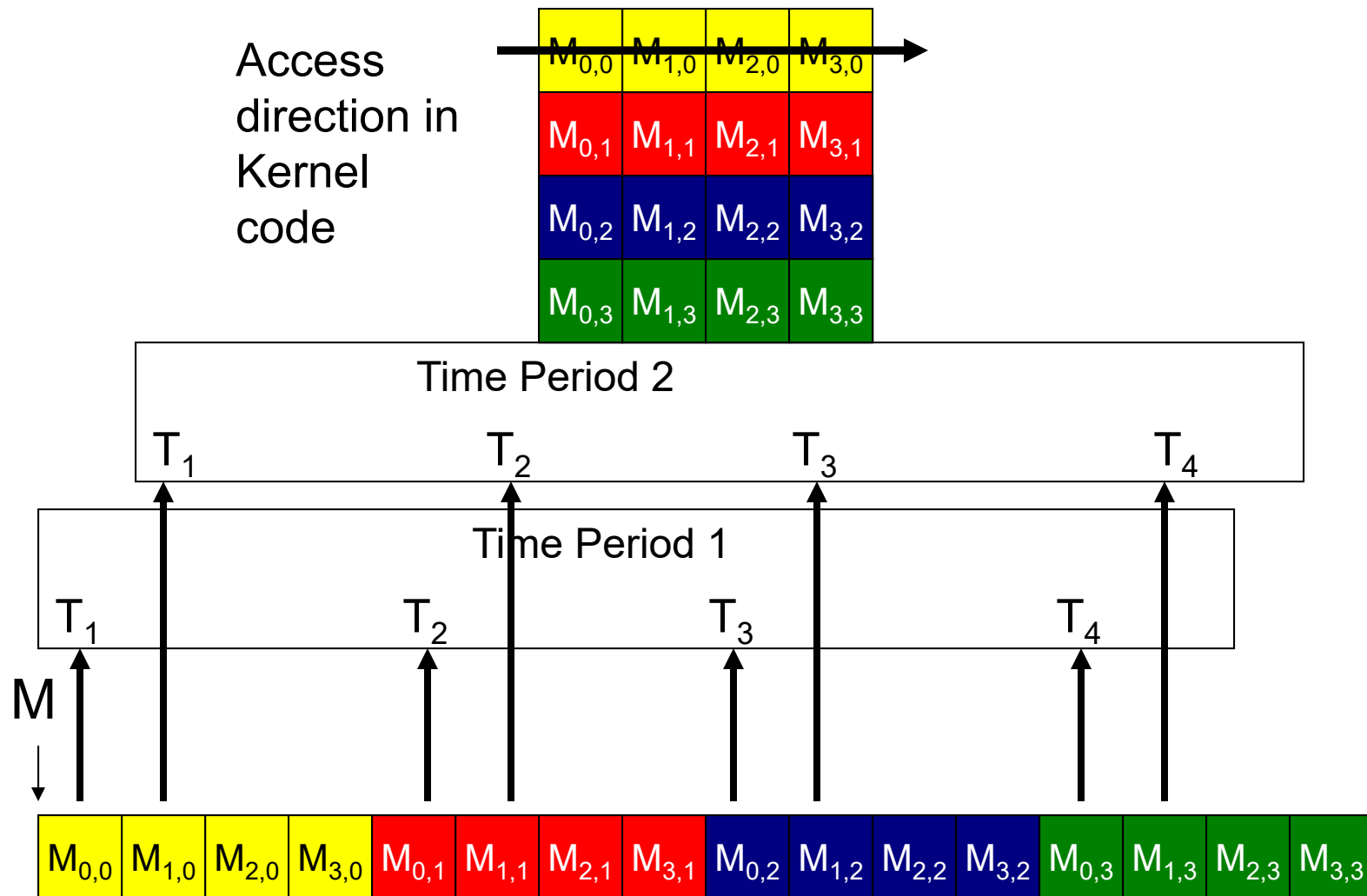


# Memory Layout of a Matrix in C

Access  
direction in  
Kernel  
code



# Memory Layout of a Matrix in C



# GPU Texturing

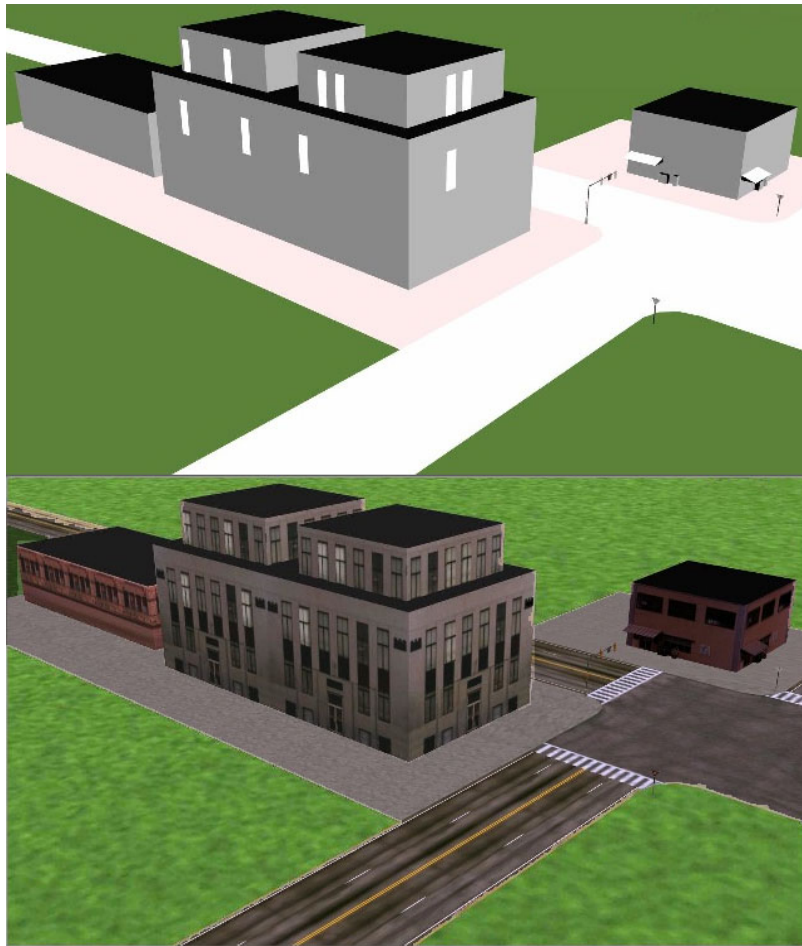
# GPU Texturing



Rage / id Tech 5 (id Software)

# Why Texturing?

- Idea: enhance visual appearance of surfaces by applying fine / high-resolution details



- Basis for most real-time rendering effects
- Look and feel of a surface
- Definition:
  - A *regularly sampled function* that is *mapped* onto every *fragment* of a surface
  - Traditionally an image, but...
- Can hold arbitrary information
  - Textures become general data structures
  - Sampled and interpreted by fragment programs
  - Can render into textures → important!



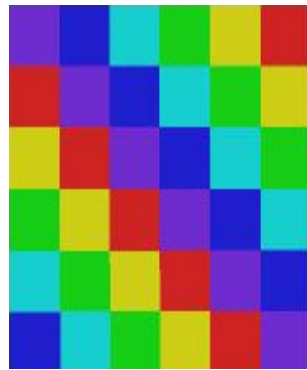


- Spatial layout
  - Cartesian grids: 1D, 2D, 3D, 2D\_ARRAY, ...
  - Cube maps, ...
- Formats (too many), e.g. OpenGL
  - GL\_LUMINANCE16\_ALPHA16
  - GL\_RGB8, GL\_RGBA8, ...: integer texture formats
  - GL\_RGB16F, GL\_RGBA32F, ...: float texture formats
  - compressed formats, high dynamic range formats, ...
- External (CPU) format vs. internal (GPU) format
  - OpenGL driver converts from external to internal

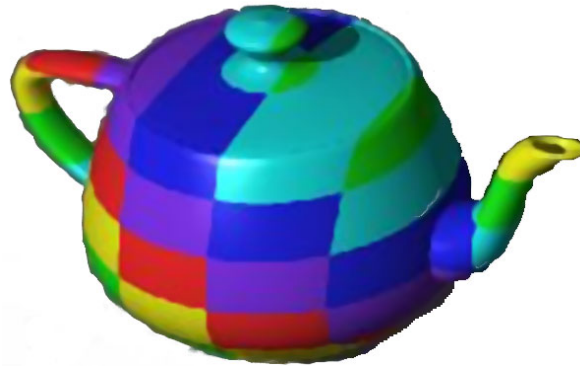




# Texturing: General Approach



Texels



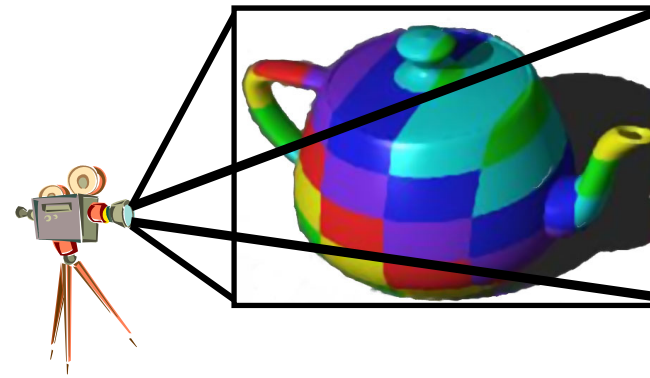
Texture space  $(u, v)$

Object space  $(x_O, y_O, z_O)$

Image Space  $(x_I, y_I)$

Parametrization

Rendering  
(Projection etc.)



# Texture Mapping

---

2D (3D) Texture Space

| Texture Transformation

2D Object Parameters

| Parameterization

3D Object Space

| Model Transformation

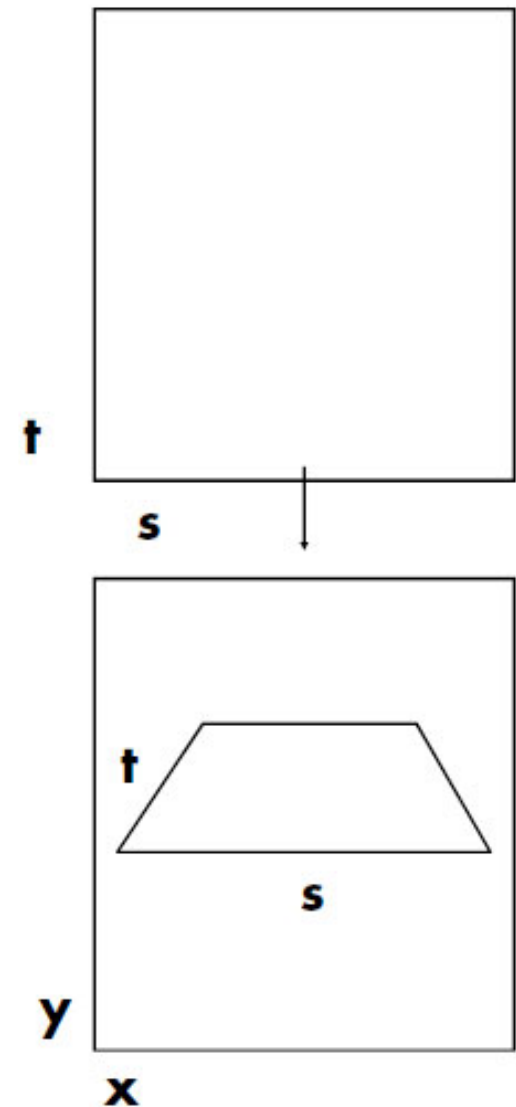
3D World Space

| Viewing Transformation

3D Camera Space

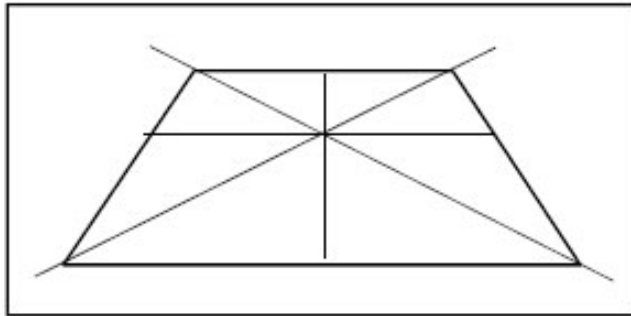
| Projection

2D Image Space

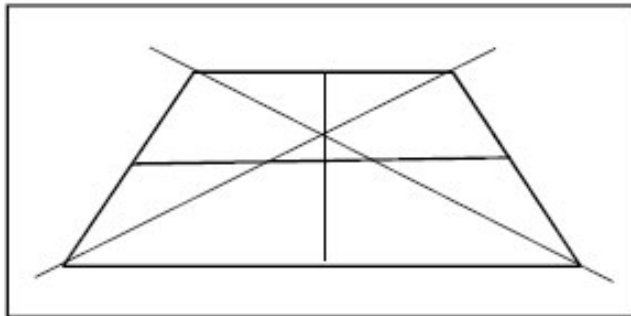


# Linear Perspective

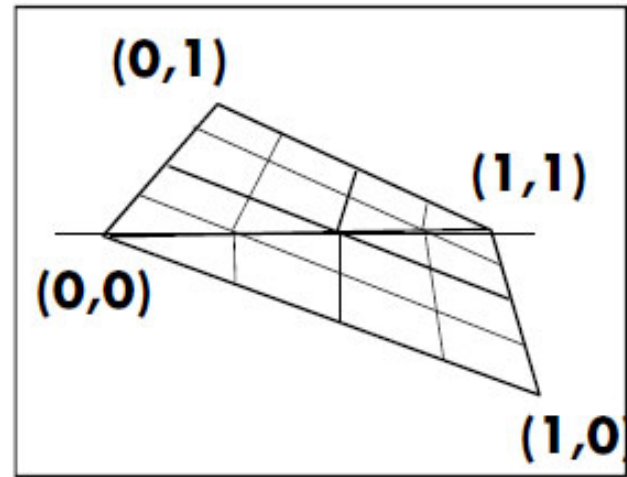
---



**Correct Linear Perspective**



**Incorrect Perspective**



**Linear Interpolation, *Bad***

**Perspective Interpolation, *Good***

# Texture Mapping Polygons

---

Forward transformation: linear projective map

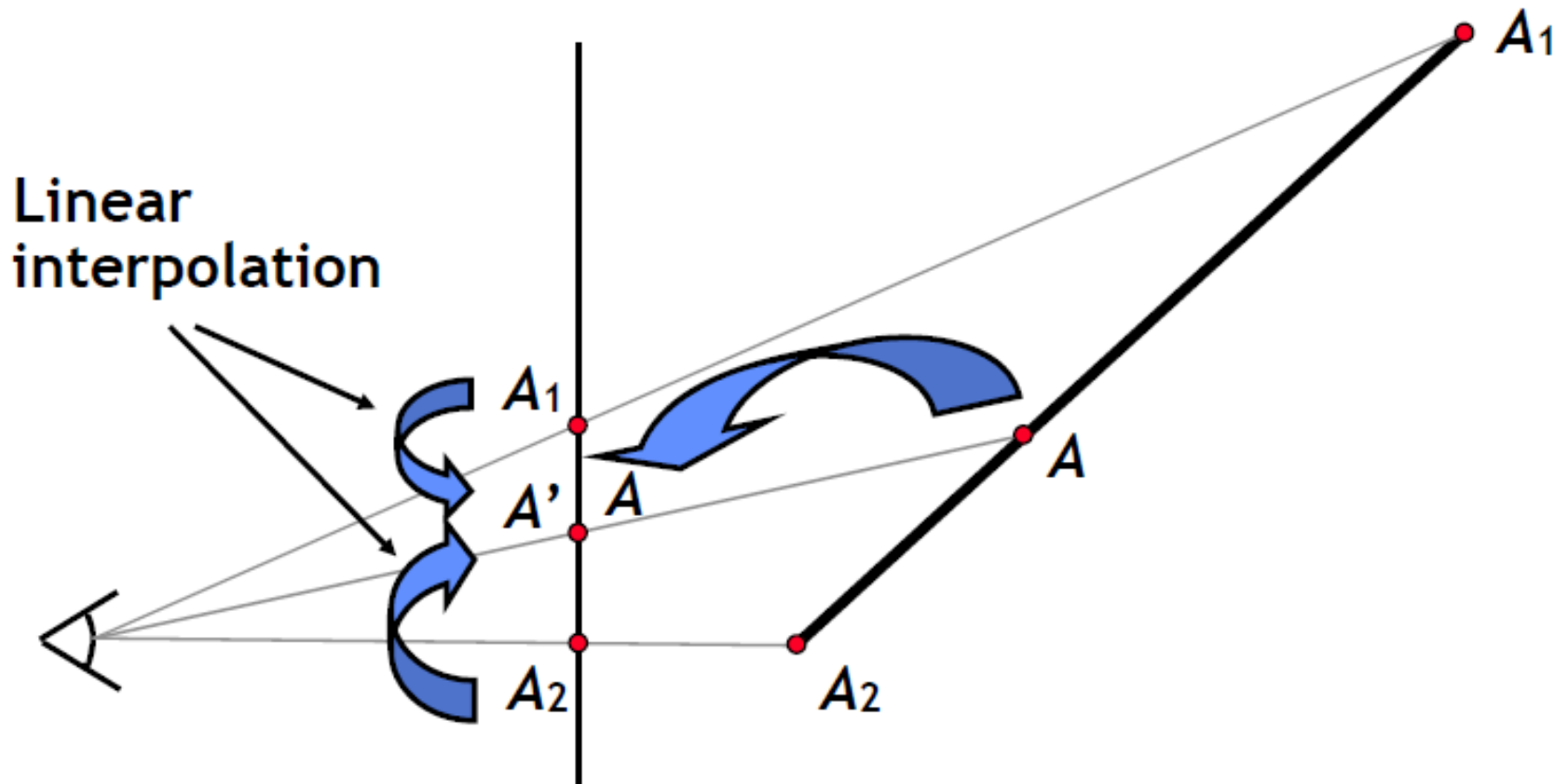
$$\begin{bmatrix} x \\ y \\ w \end{bmatrix} = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \begin{bmatrix} s \\ t \\ r \end{bmatrix}$$

Backward transformation: linear projective map

$$\begin{bmatrix} s \\ t \\ r \end{bmatrix} = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}^{-1} \begin{bmatrix} x \\ y \\ w \end{bmatrix}$$

# Incorrect attribute interpolation

---



# Linear interpolation

---

Compute intermediate attribute value

- Along a line:  $A = aA_1 + bA_2$ ,  $a+b=1$
- On a plane:  $A = aA_1 + bA_2 + cA_3$ ,  $a+b+c=1$

Only projected values interpolate linearly in screen space (straight lines project to straight lines)

- $x$  and  $y$  are projected (divided by  $w$ )
- Attribute values are not naturally projected

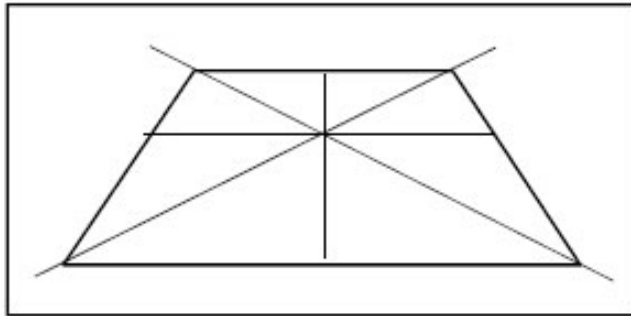
Choice for attribute interpolation in screen space

- Interpolate unprojected values
  - Cheap and easy to do, but gives wrong values
  - Sometimes OK for color, but
  - Never acceptable for texture coordinates
- Do it right

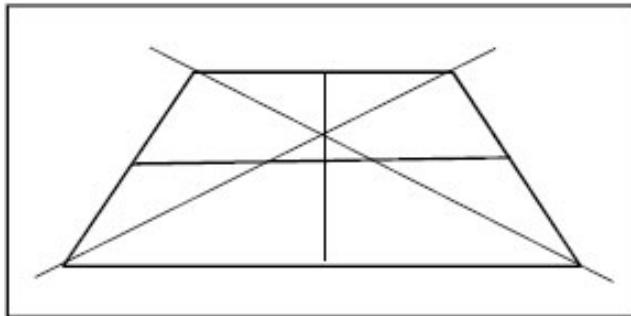


# Linear Perspective

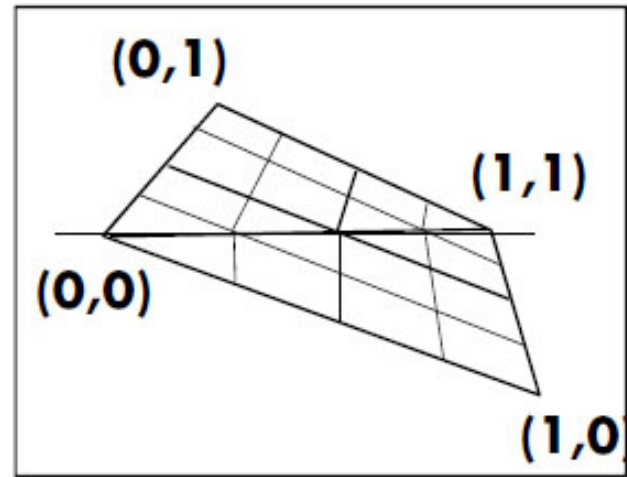
---



**Correct Linear Perspective**



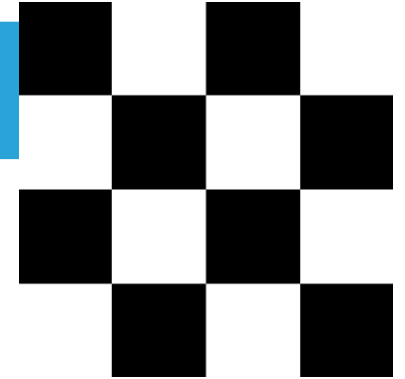
**Incorrect Perspective**



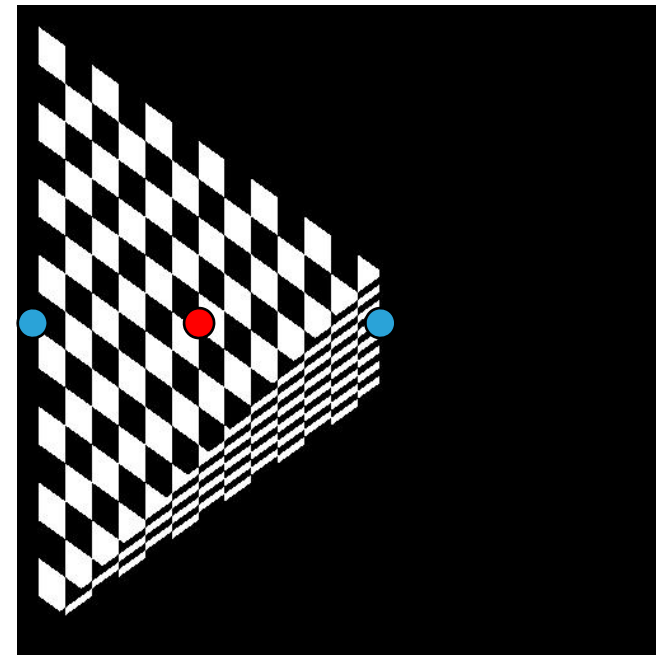
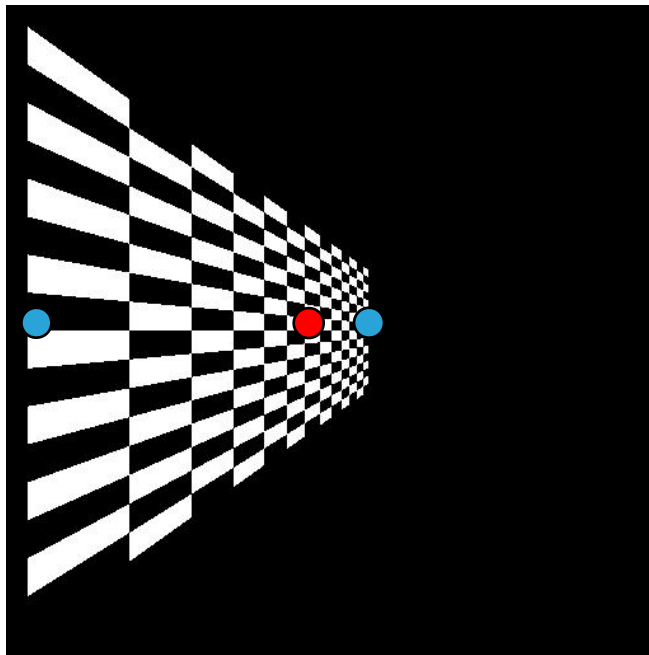
**Linear Interpolation, *Bad***

**Perspective Interpolation, *Good***

# Perspective Texture Mapping



linear interpolation in object space  $\frac{ax_1 + bx_2}{aw_1 + bw_2} \neq a \frac{x_1}{w_1} + b \frac{x_2}{w_2}$  linear interpolation in screen space



$$a = b = 0.5$$





# Homogeneous Coordinates (1)



## Projective geometry

- (Real) projective spaces  $\mathbb{RP}^n$ :  
Real projective line  $\mathbb{RP}^1$ , real projective plane  $\mathbb{RP}^2$ , ...
- A point in  $\mathbb{RP}^n$  is a line through the origin (i.e., all the scalar multiples of the same vector) in an  $(n+1)$ -dimensional (real) vector space



## Homogeneous coordinates of 2D projective point in $\mathbb{RP}^2$

- Coordinates differing only by a non-zero factor  $\lambda$  map to the same point  
 $(\lambda x, \lambda y, \lambda)$       dividing out the  $\lambda$  gives  $(x, y, 1)$ , corresponding to  $(x, y)$  in  $\mathbb{R}^2$
- Coordinates with last component = 0 map to “points at infinity”  
 $(\lambda x, \lambda y, 0)$       division by last component not allowed; but again this is the same point if it only differs by a scalar factor, e.g., this is the same point as  $(x, y, 0)$

# Homogeneous Coordinates (2)



## Examples of usage

- Translation (with translation vector  $\vec{b}$ )
- Affine transformations (linear transformation + translation)

$$\vec{y} = A\vec{x} + \vec{b}.$$

- With homogeneous coordinates:

$$\begin{bmatrix} \vec{y} \\ 1 \end{bmatrix} = \left[ \begin{array}{ccc|c} & A & & \vec{b} \\ 0 & \dots & 0 & 1 \end{array} \right] \begin{bmatrix} \vec{x} \\ 1 \end{bmatrix}$$

- Setting the last coordinate = 1 and the last row of the matrix to  $[0, \dots, 0, 1]$  results in translation of the point  $\vec{x}$  (via addition of translation vector  $\vec{b}$ )
- The matrix above is a linear map, but because it is one dimension higher, it does not have to move the origin in the  $(n+1)$ -dimensional space for translation

# Homogeneous Coordinates (3)

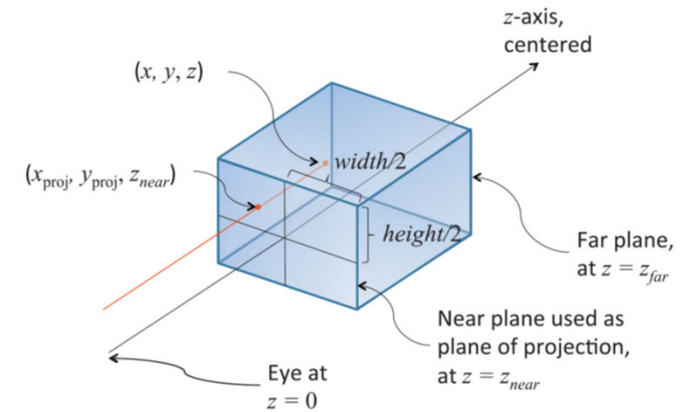


## Examples of usage

- Projection (e.g., OpenGL projection matrices)

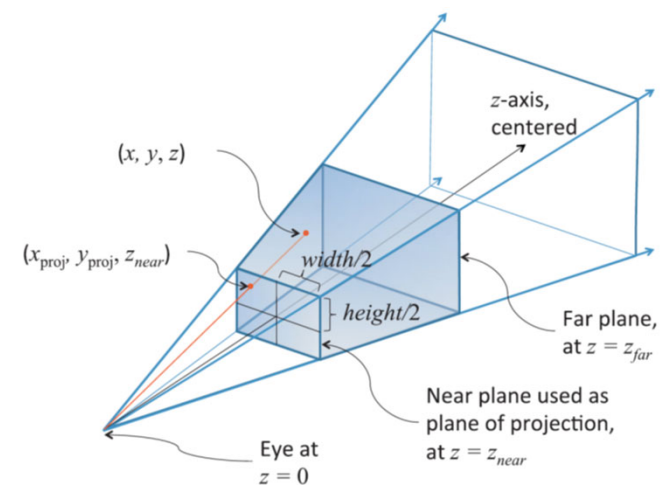
$$\begin{bmatrix} \frac{2}{\text{right} - \text{left}} & 0 & 0 & \frac{\text{right} + \text{left}}{\text{right} - \text{left}} \\ 0 & \frac{2}{\text{top} - \text{bottom}} & 0 & \frac{\text{top} + \text{bottom}}{\text{top} - \text{bottom}} \\ 0 & 0 & \frac{-2}{\text{far} - \text{near}} & \frac{\text{far} + \text{near}}{\text{far} - \text{near}} \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

orthographic



$$\begin{bmatrix} \frac{z_{\text{near}}}{\text{width}/2} & 0.0 & \frac{\text{left} + \text{right}}{\text{width}/2} & 0.0 \\ 0.0 & \frac{z_{\text{near}}}{\text{height}/2} & \frac{\text{top} + \text{bottom}}{\text{height}/2} & 0.0 \\ 0.0 & 0.0 & \frac{z_{\text{far}} + z_{\text{near}}}{z_{\text{far}} - z_{\text{near}}} & \frac{2z_{\text{far}}z_{\text{near}}}{z_{\text{far}} - z_{\text{near}}} \\ 0.0 & 0.0 & -1.0 & 0.0 \end{bmatrix}$$

perspective



Thank you.