



# Parallel Computing

Prof. Marco Bertini



### Data parallelism: GPU computing



## 2D convolution: tile boundaries



### **2D Image Matrix with Automated Padding**

- It is sometimes desirable to pad each row of a 2D matrix to multiples of DRAM bursts
  - So each row starts at the DRAM burst boundary
  - Effectively adding columns
  - This is usually done automatically by matrix allocation function

width

M<sub>0,0</sub>

 $M_{1,0}$ 

 $M_{2,0}$ 

| M<sub>0,1</sub>

 $M_{1,1}$ 

M<sub>2,1</sub>

 $M_{0.2}$ 

 $M_{1.2}$ 

 $M_{22}$ 

pitch

Padded

elements

- Pitch can be different for different hardware
- Example: a 3X3 matrix padded into a 3X4 matrix





### **Row-Major Layout with Pitch**





## Sample image struct

// Image Matrix Structure declaration

//

typedef struct {

int width;

int height;

int pitch;

int channels;

float\* data;





# Setting Block Size

#### #define O\_TILE\_WIDTH 12

#### #define BLOCK\_WIDTH (0\_TILE\_WIDTH + 4)

dim3 dimBlock(BLOCK\_WIDTH,BLOCK\_WIDTH);

dim3 dimGrid((Image\_Width-1)/0\_TILE\_WIDTH+1,
(Image\_Height-1)/0\_TILE\_WIDTH+1, 1)

In general, BLOCK\_WIDTH should be

• O\_TILE\_WIDTH + (MASK\_WIDTH-1)

### Using constant memory and caching for Mask

- Mask is used by all threads but not modified in the convolution kernel
  - All threads in a warp access the same locations at each point in time
- CUDA devices provide constant memory whose contents are aggressively cached
  - Cached values are broadcast to all threads in a warp
  - Effectively magnifies memory bandwidth without consuming shared memory
- Use of const \_\_\_restrict\_\_ qualifiers for the mask parameter informs the compiler that it is eligible for constant caching, for example:

\_\_global\_\_ void convolution\_2D\_kernel(float \*P, float \*N, int height, int width, int channels, const float \_\_restrict\_\_ \*M);



## Shifting from output coordinates to input coordinate

int tx = threadIdx.x;

int ty = threadIdx.y;

int row\_o = blockIdx.y\*0\_TILE\_WIDTH + ty;

int col\_o = blockIdx.x\*0\_TILE\_WIDTH + tx;





}

### Taking Care of Boundaries (1 channel example)

if((row\_i >= 0) && (row\_i < height) &&

(col\_i >= 0) && (col\_i < width)) {</pre>

Ns[ty][tx] = data[row\_i \* width + col\_i];

} else{

Ns[ty][tx] = 0.0f;

Use of width here is OK if pitch is set to width (no padding)



## Calculating output

Some threads do not participate in calculating output

float output = 0.0f;

if(ty <  $0_TILE_WIDTH \&\& tx < 0_TILE_WIDTH$ ){

for(i = 0; i < MASK\_WIDTH; i++) {</pre>

for(j = 0; j < MASK\_WIDTH; j++) {</pre>

output += M[i][j] \* Ns[i+ty][j+tx];





# Writing output

- Some threads do not write output (1 channel example)
  - if(row\_o < height && col\_o < width)</pre>
    - data[row\_o\*width + col\_o] =
      output;







- These slides report material from:
  - NVIDIA GPU Teaching Kit







 Programming Massively Parallel Processors: A Hands-on Approach, D. B. Kirk and W-M. W. Hwu, Morgan Kaufman - Chapt. 8

