ECE408/CS483/CSE408  Spring 2020

Applied Parallel Programming

## Lecture 8: Convolution, Constant Memory and Constant Caching

1

---

# Objective

- To learn convolution, an important parallel computation pattern
  – Widely used in signal, image and video processing
  – Foundational to stencil computation used in many science and engineering applications
  – Critical component of Neural Networks and Deep Learning

- Important techniques
  – Taking advantage of cached memories

2

---

# Convolution Applications

- A popular array operation that is used in various forms in signal processing, digital recording, image processing, video processing, computer vision, and machine learning.
- Convolution is often performed as a filter that transforms signals and pixels into more desirable values.
  – Some filters smooth out the signal values so that one can see the big-picture trend
  – Others like Gaussian filters can be used to sharpen boundaries and edges of objects in images..

3

---

# Convolution Computation

- An array operation where each output data element is a weighted sum of a collection of neighboring input elements
- The weights used in the weighted sum calculation are defined by an input mask array, commonly referred to as the *convolution kernel*
  – we will refer to these mask arrays as convolution masks or convolution filters to avoid confusion.
  – The same convolution mask is typically used for all elements of the array.
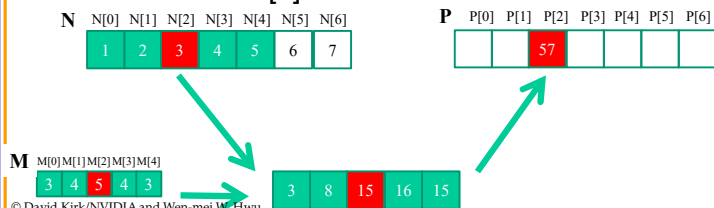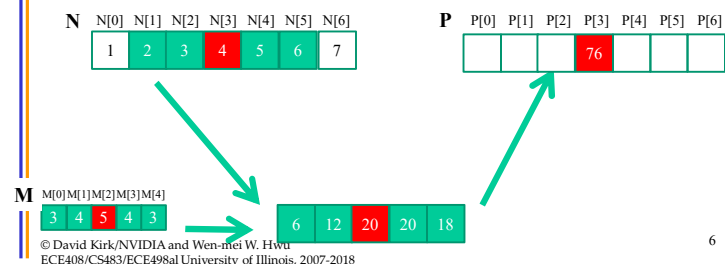
4

## 1D Convolution Example

- Commonly used for audio processing
  - Mask_Width is usually an odd number of elements for symmetry (5 in this example)
  - Mask_Radius is the number of elements used in convolution on each side of the pixel being calculated (2 in this example).
- Calculation of P[2]:

**N**  N[0] N[1] N[2] N[3] N[4] N[5] N[6]

| 1 | 2 | 3 | 4 | 5 | 6 | 7 |

**P**  P[0] P[1] P[2] P[3] P[4] P[5] P[6]

| | | 57 | | | | |

**M**  M[0] M[1] M[2] M[3] M[4]

| 3 | 4 | 5 | 4 | 3 |

| 3 | 8 | 15 | 16 | 15 |

5

---

## 1D Convolution Example - more on inside elements

- Calculation of P[3]

**N**  N[0] N[1] N[2] N[3] N[4] N[5] N[6]

| 1 | 2 | 3 | 4 | 5 | 6 | 7 |

**P**  P[0] P[1] P[2] P[3] P[4] P[5] P[6]

| | | | 76 | | | |

**M**  M[0] M[1] M[2] M[3] M[4]

| 3 | 4 | 5 | 4 | 3 |

| 6 | 12 | 20 | 20 | 18 |

6

---
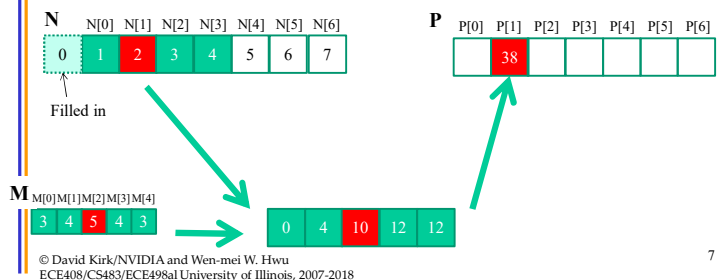
## 1D Convolution Boundary Condition

- Calculation of output elements near the boundaries (beginning and end) of the input array need to deal with "ghost" elements
  - Different policies (0, replicates of boundary values, etc.)

**N**  N[0] N[1] N[2] N[3] N[4] N[5] N[6]

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

Filled in

**P**  P[0] P[1] P[2] P[3] P[4] P[5] P[6]

| | 38 | | | | | |

**M**  M[0] M[1] M[2] M[3] M[4]

| 3 | 4 | 5 | 4 | 3 |

| 0 | 4 | 10 | 12 | 12 |

7

---

## A 1D Convolution Kernel with Boundary Condition Handling

- This kernel forces all elements outside the valid data index range to 0
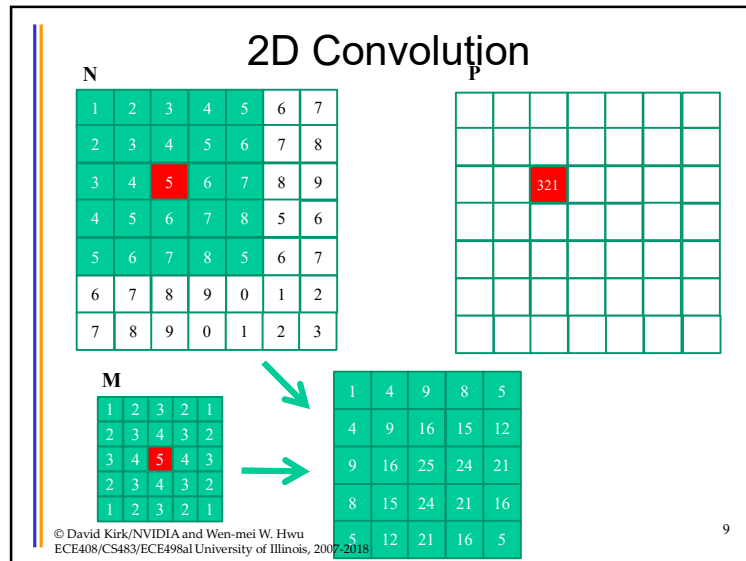
```
__global__ void convolution_1D_basic_kernel(float *N, float *M, float *P,
int Mask_Width, int Width) {

    int i = blockIdx.x*blockDim.x + threadIdx.x;

    float Pvalue = 0;
    int N_start_point = i - (Mask_Width/2);
    for (int j = 0; j < Mask_Width; j++) {
        if (N_start_point + j >= 0 && N_start_point + j < Width) {
            Pvalue += N[N_start_point + j]*M[j];
        }
    }
    P[i] = Pvalue;

}
```
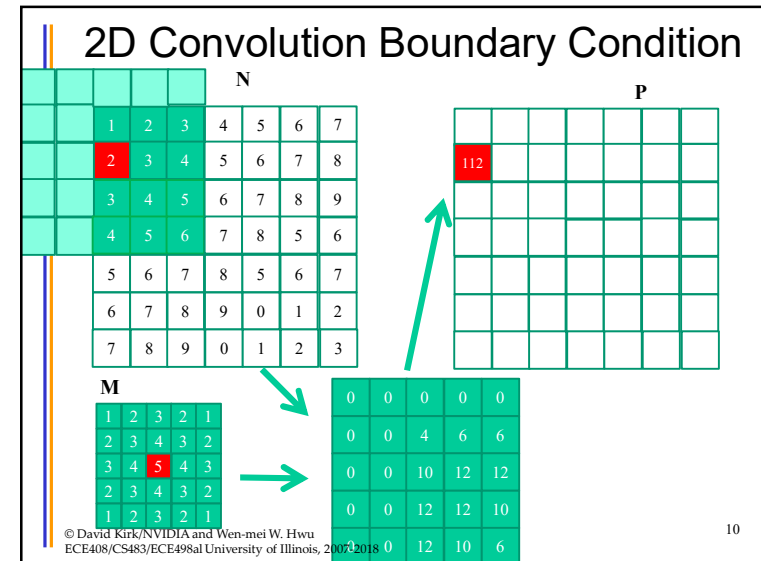
8

2

## 2D Convolution

**N**

| 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|
| 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 4 | 5 | 6 | 7 | 8 | 5 | 6 |
| 5 | 6 | 7 | 8 | 5 | 6 | 7 |
| 6 | 7 | 8 | 9 | 0 | 1 | 2 |
| 7 | 8 | 9 | 0 | 1 | 2 | 3 |

**P**

321

**M**

| 1 | 2 | 3 | 2 | 1 |
|---|---|---|---|---|
| 2 | 3 | 4 | 3 | 2 |
| 3 | 4 | 5 | 4 | 3 |
| 2 | 3 | 4 | 3 | 2 |
| 1 | 2 | 3 | 2 | 1 |

| 1 | 4 | 9 | 8 | 5 |
|---|---|---|---|---|
| 4 | 9 | 16 | 15 | 12 |
| 9 | 16 | 25 | 24 | 21 |
| 8 | 15 | 24 | 21 | 16 |
| 5 | 12 | 21 | 16 | 5 |

9

---

## 2D Convolution Boundary Condition

**N**

| 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|
| 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 4 | 5 | 6 | 7 | 8 | 5 | 6 |
| 5 | 6 | 7 | 8 | 5 | 6 | 7 |
| 6 | 7 | 8 | 9 | 0 | 1 | 2 |
| 7 | 8 | 9 | 0 | 1 | 2 | 3 |

**P**

112

**M**

| 1 | 2 | 3 | 2 | 1 |
|---|---|---|---|---|
| 2 | 3 | 4 | 3 | 2 |
| 3 | 4 | 5 | 4 | 3 |
| 2 | 3 | 4 | 3 | 2 |
| 1 | 2 | 3 | 2 | 1 |

| 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|
| 0 | 0 | 4 | 6 | 6 |
| 0 | 0 | 10 | 12 | 12 |
| 0 | 0 | 12 | 12 | 10 |
| 0 | 0 | 12 | 10 | 6 |

10

---

## 2D Convolution – Ghost Cells

**N**

| 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|
| 0 | 3 | 4 | 5 | 6 |
| 0 | 2 | 3 | 4 | 5 |
| 0 | 3 | 5 | 6 | 7 |
| 0 | 1 | 1 | 3 | 1 |

**P**

79

**M**

| 1 | 2 | 3 | 2 | 1 |
|---|---|---|---|---|
| 2 | 3 | 4 | 3 | 2 |
| 3 | 4 | 5 | 4 | 3 |
| 2 | 3 | 4 | 3 | 2 |
| 1 | 2 | 3 | 2 | 1 |

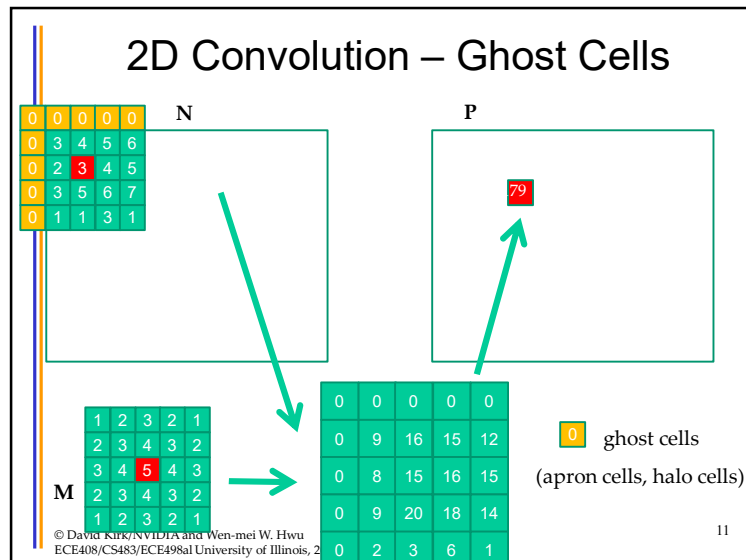| 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|
| 0 | 9 | 16 | 15 | 12 |
| 0 | 8 | 15 | 16 | 15 |
| 0 | 9 | 20 | 18 | 14 |
| 0 | 2 | 3 | 6 | 1 |

**0** ghost cells
(apron cells, halo cells)
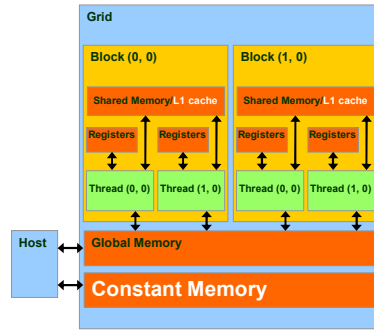
11

---

## Access Pattern for M

- Elements of M are called mask (kernel, filter) coefficients (weights)
  - Calculation of all output P elements needs M
  - M is not changed during grid execution

- Bonus - M elements are accessed in the same order when calculating all P elements

- M is a good candidate for Constant Memory

12

---

## Programmer View of CUDA Memories (Review)

- Each thread can:
  - Read/write per-thread **registers (~1 cycle)**
  - Read/write per-block **shared memory (~5 cycles)**
  - Read/write per-grid **global memory (~500 cycles)**
  - Read/only per-grid **constant memory (~5 cycles with caching)**

13

13

## Memory Hierarchies

- Review: If we had to go to global memory to access data all the time, the execution speed of GPUs would be limited by the global memory bandwidth
  - We saw the use of shared memory (scratchpad) in tiled matrix multiplication.

- Another important solution: Caches

14

14

## Caches Store Lines of Memory

Recall: memory bursts
- contain around **1024 bits** (**128B**) from
- consecutive (linear) addresses.
- Let's call a single burst a **line**.

What's a **cache**?
- An **array of cache lines** (and tags).
- Memory **read produces** a **line**,
- **cache stores** a **copy** of the line, and
- tag records line's memory address.

15

15

## Memory Accesses Show Locality

An executing program
- loads and store data from memory.
- **Consider sequence of addresses** accessed.
**Sequence** usually **shows** two types of **locality**:
  - **spatial**: accessing **X implies** accessing **X+1** (and X+2, and so forth) **soon**
  - **temporal**: accessing **X** implies accessing **X again soon**
  (Caches improve performance for both types.)

16

16

4

## Caches Can't Hold Everything

Caches are smaller than memory.

**When cache is full**,
- must make room for new line,
- usually by **discard**ing **least recently used line**.

17

17

## GPU Has Constant and L1 Caches

**To support writes** (modification of lines),
- **changes** must be **copied back to memory**, and
- cache must **track** modification **status**.
- **L1 cache** in GPU (for global memory accesses) **supports writes**.

**Cache for constant** / texture **memory**
- Special case: **lines are read-only**
- Enables higher-throughput access than L1 for common GPU kernel access patterns.

18

18

## Cache vs. Scratchpad (GPU Shared Mem.)

- Caches vs. shared memory
  - Both on chip*, with similar performance
  - (As of Volta generation, both using the same physical resources, allocated dynamically!)

**What's the difference?**
- **Programmer controls shared memory** contents (called a scratchpad)
- **Microarchitecture** automatically **determines contents of cache**.
  - *Static RAM, not DRAM, by the way—see ECE120/CS233.

19

19

## How to Use Constant Memory

**Host code** is **similar** to previous versions, but…

**Allocate** device memory for **M** (the mask)
- **outside of all functions**
- **using __constant__** (tells GPU that caching is safe).

**For copying** to device memory, **use**
- **cudaMemcpyToSymbol(dst, src, size)**
- with destination defined as above.

20

20

## Example of Host Code

(**MASK_WIDTH** is the size of the mask.)

```
// outside of any kernel/function
static __constant__ float Mc[MASK_WIDTH][MASK_WIDTH];

// allocate N, P, initialize N elements, copy N to Nd

// in host code:
    float* M; // host memory copy of mask
    // initialize M
    cudaMemcpyToSymbol(Mc, M,
        MASK_WIDTH * MASK_WIDTH * sizeof(M[0]));
    ConvolutionKernel<<<dimGrid, dimBlock>>>(Nd, Pd);
    // (note that file-scope Mc is visible to kernel)
```

21

21

## ANY MORE QUESTIONS?
## READ CHAPTER 7

22

22