

# CS 179: GPU Programming

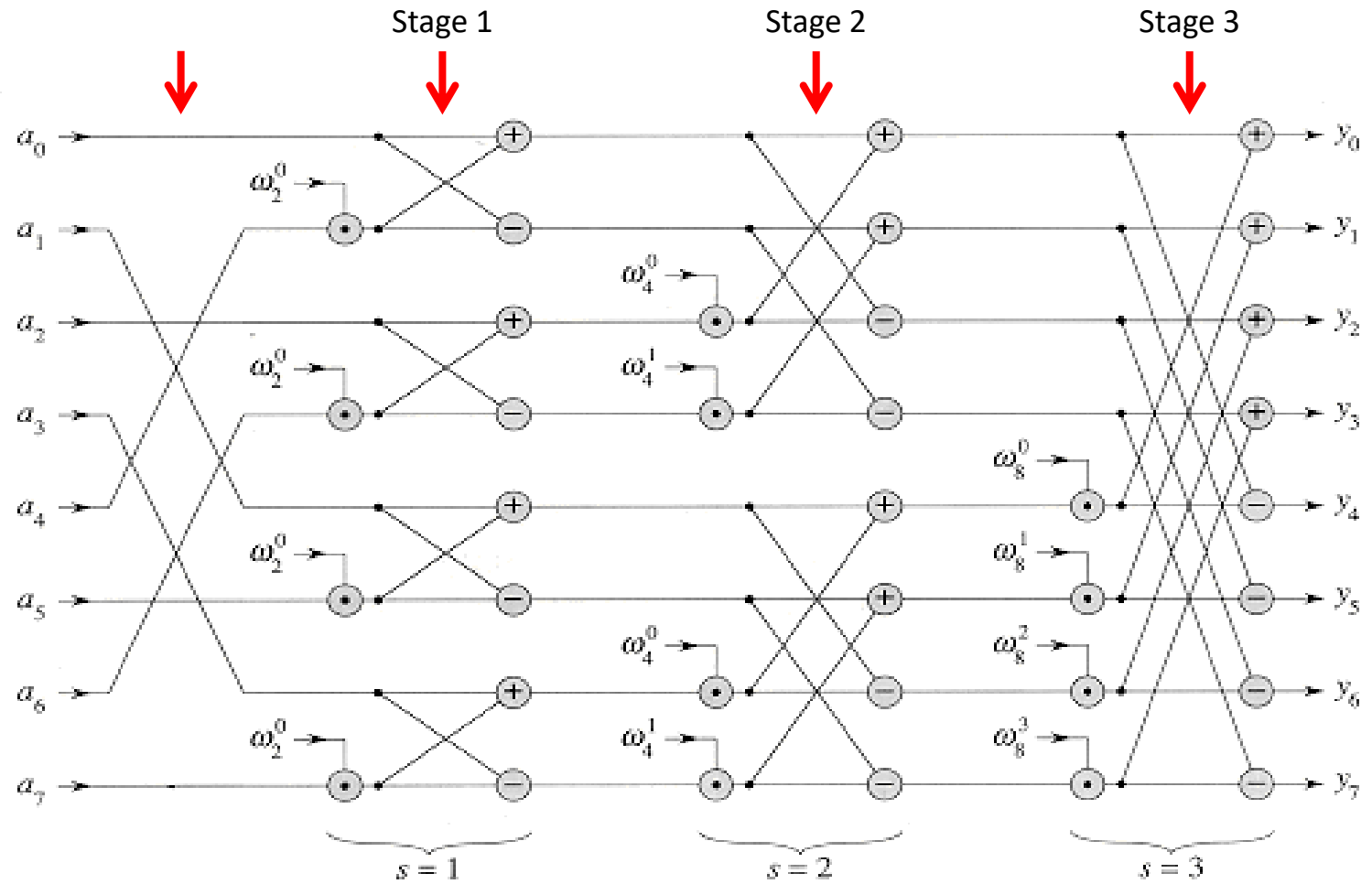
Lecture 9 / Homework 3

# Recap

- Some algorithms are “less obviously parallelizable”:
  - Reduction
  - Sorts
  - FFT (and certain recursive algorithms)

# Parallel FFT structure (radix-2)

Bit-reversed access



# cuFFT 1D example

```
#define NX 262144

cufftComplex *data_host
    = (cufftComplex*)malloc(sizeof(cufftComplex)*NX);
cufftComplex *data_back
    = (cufftComplex*)malloc(sizeof(cufftComplex)*NX);

// Get data...

cufftHandle plan;
cufftComplex *data1;
cudaMalloc((void**)&data1, sizeof(cufftComplex)*NX);
cudaMemcpy(data1, data_host, NX*sizeof(cufftComplex), cudaMemcpyHostToDevice);

/* Create a 1D FFT plan. */
int batch = 1; // Number of transforms to run
cufftPlan1d(&plan, NX, CUFFT_C2C, batch);

/* Transform the first signal in place. */
cufftExecC2C(plan, data1, data1, CUFFT_FORWARD);

/* Inverse transform in place. */
cufftExecC2C(plan, data1, data1, CUFFT_INVERSE);

cudaMemcpy(data_back, data1, NX*sizeof(cufftComplex), cudaMemcpyDeviceToHost);
```

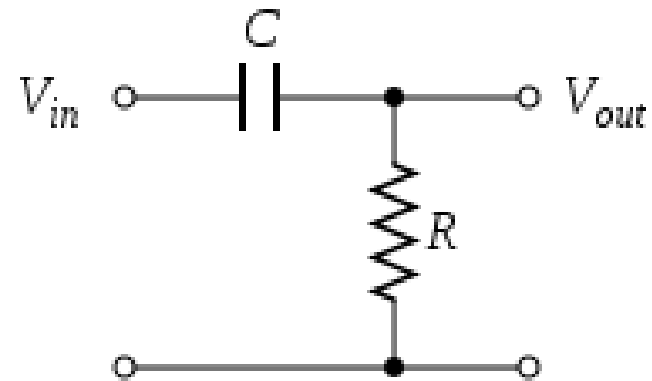
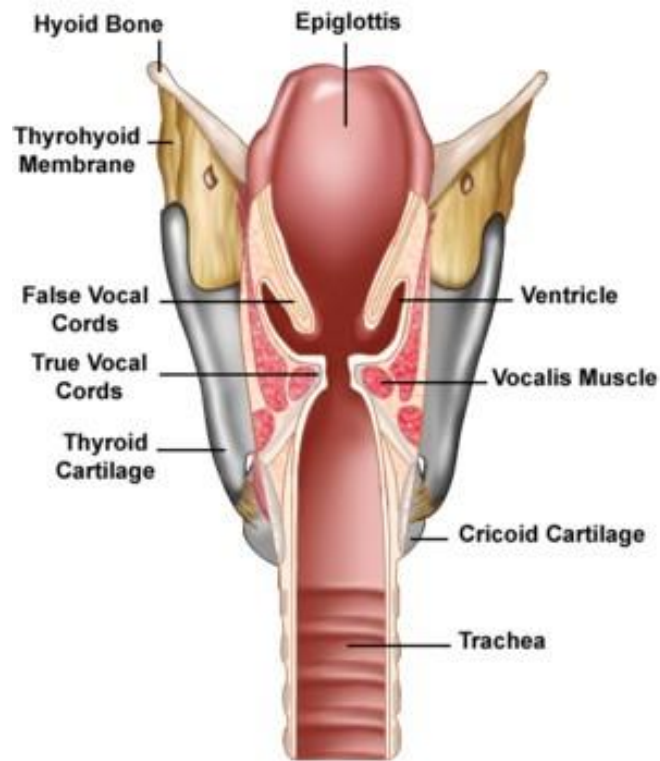
Correction:  
Remember to use  
cufftDestroy(plan)  
when finished with  
transforms

# Today

- Homework 3
  - Large-kernel convolution

# Systems

- Given input signal(s), produce output signal(s)



# LTI system review (Week 1)

- “Linear time-invariant” (LTI) systems
  - Lots of them!
- Can be characterized entirely by “impulse response”  $h[n]$
- Output given from input by *convolution*:

$$y[n] = \sum_{k=-\infty}^{\infty} x[k]h[n-k]$$

# Parallelization

$$y[n] = \sum_{k=-\infty}^{\infty} x[k]h[n-k]$$

- Convolution is parallelizable!
  - Sequential pseudocode (ignoring boundary conditions):

```
(set all y[i] to 0)
For (i from 0 through x.length - 1)
    for (j from 0 through h.length - 1)
        y[i] += (appropriate terms from x and h)
```



# A problem...

- This worked for *small* impulse responses
  - E.g.  $h[n]$ ,  $0 \leq n \leq 20$  in HW 1
- Homework 1 was “small-kernel convolution”:
  - (Vocab alert: Impulse responses are often called “kernels”!)

# A problem...

$$y[n] = \sum_{k=-\infty}^{\infty} x[k]h[n-k]$$

- Sequential runtime:  $O(n*m)$ 
  - (n: size of x)
  - (m: size of h)
- Troublesome for large m! (i.e. large impulse responses)

(set all  $y[i]$  to 0)

For (i from 0 through  $x.length - 1$ )

    for (j from 0 through  $h.length - 1$ )

$y[i] +=$  (appropriate terms from x and h)

# DFT/FFT

- Same problem with Discrete Fourier Transform!

$$X_k \stackrel{\text{def}}{=} \sum_{n=0}^{N-1} x_n \cdot e^{-2\pi i k n / N}, \quad k \in \mathbb{Z}$$

- Successfully optimized *and* GPU-accelerated!
  - $O(n^2)$  to  $O(n \log n)$
  - Can we do the same here?

# “Circular” convolution

- “equivalent” of convolution for periodic signals

# “Circular” convolution

- Linear convolution:

$$y[n] = \sum_{k=-\infty}^{\infty} x[k] h[n - k]$$

- Circular convolution:

$$y[n] = \sum_{k=0}^{N-1} x[k] h[(n - k) \bmod N]$$

x0	x1	x2	x3
----	----	----	----

h0	h1	0	0
----	----	---	---

# Example:

$$y[n] = \sum_{k=-\infty}^{\infty} x[k]h[n-k]$$

- $x[0..3], h[0..1]$
- Linear convolution:

$$y[0] = x[0]h[0]$$

0	0	<b>x0</b>	<b>x1</b>	<b>x2</b>	<b>x3</b>	0	0
0	<b>h1</b>	<b>h0</b>	0	0	0	0	0

x0	x1	x2	x3
----	----	----	----

h0	h1	0	0
----	----	---	---

# Example:

$$y[n] = \sum_{k=-\infty}^{\infty} x[k]h[n-k]$$

- $x[0..3], h[0..1]$
- Linear convolution:

$$y[0] = x[0]h[0]$$

$$y[1] = x[0]h[1] + x[1]h[0]$$

0	0	<b>x0</b>	<b>x1</b>	<b>x2</b>	<b>x3</b>	0	0
0	0	<b>h1</b>	<b>h0</b>	0	0	0	0

x0	x1	x2	x3
----	----	----	----

h0	h1	0	0
----	----	---	---

# Example:

$$y[n] = \sum_{k=-\infty}^{\infty} x[k]h[n-k]$$

- $x[0..3], h[0..1]$
- Linear convolution:

$$y[0] = x[0]h[0]$$

$$y[1] = x[0]h[1] + x[1]h[0]$$

$$y[2] = x[1]h[1] + x[2]h[0]$$

$$y[3] = x[2]h[1] + x[3]h[0]$$

0	0	<b>x0</b>	<b>x1</b>	<b>x2</b>	<b>x3</b>	0	0
0	0	0	<b>h1</b>	<b>h0</b>	0	0	0



x0	x1	x2	x3
----	----	----	----

h0	h1	0	0
----	----	---	---

# Example:

$$y[n] = \sum_{k=-\infty}^{\infty} x[k]h[n-k]$$

- $x[0..3], h[0..1]$
- Linear convolution:

$$y[0] = x[0]h[0]$$

$$y[1] = x[0]h[1] + x[1]h[0]$$

$$y[2] = x[1]h[1] + x[2]h[0]$$

$$y[3] = x[2]h[1] + x[3]h[0]$$

$$y[4] = x[3]h[1] + x[4]h[0]$$

0	0	<b>x0</b>	<b>x1</b>	<b>x2</b>	<b>x3</b>	0	0
0	0	0	0	<b>h1</b>	<b>h0</b>	0	0

x0	x1	x2	x3
----	----	----	----

h0	h1	0	0
----	----	---	---

# Example:

$$y[n] = \sum_{k=-\infty}^{\infty} x[k]h[n-k]$$

- $x[0..3]$ ,  $h[0..1]$
- Linear convolution:

$$y[0] = x[0]h[0]$$

$$y[1] = x[0]h[1] + x[1]h[0]$$

$$y[2] = x[1]h[1] + x[2]h[0]$$

$$y[3] = x[2]h[1] + x[3]h[0]$$

$$y[4] = x[3]h[1] + x[4]h[0]$$

0	0	<b>x0</b>	<b>x1</b>	<b>x2</b>	<b>x3</b>	0	0
0	<b>h1</b>	<b>h0</b>	0	0	0	0	0

- Circular convolution:

$$y[n] = \sum_{k=0}^{N-1} x[k] h[(n-k) \bmod N]$$

$$y[0] = x[3]h[1] + x[0]h[0] + x[1]h[3] + x[2]h[2]$$

$$y[1] = x[0]h[1] + x[1]h[0] + x[2]h[3] + x[3]h[2]$$

$$y[2] = x[1]h[1] + x[2]h[0] + x[3]h[3] + x[0]h[2]$$

$$y[3] = x[2]h[1] + x[3]h[0] + x[0]h[3] + x[1]h[2]$$

<b>x0</b>	<b>x1</b>	<b>x2</b>	<b>x3</b>
<b>h0</b>	0	0	<b>h1</b>

# Padding

- $x[n]$  and  $h[n]$  are different lengths?
  - $h[2]$  and  $h[3]$  are undefined in the previous slide
  - We will treat them as zeros
  - To do this, store  $h$  and  $x$  in bigger arrays
  - Pad them with enough zeros so they both have length at least  $N = x.length + h.length - 1$

# Padding

- $x[n]$  and  $h[n]$  are different lengths?
  - $h[2]$  and  $h[3]$  are undefined in the previous slide
  - We will treat them as zeros
  - To do this, store  $h$  and  $x$  in bigger arrays
  - Pad them with enough zeros so they both have length at least  $N = x.length + h.length - 1$
- This will let us linearly convolve using circular convolution

# Example: (Padding)

- $x[0..3], h[0..1]$
- Linear convolution:

$$y[n] = \sum_{k=-\infty}^{\infty} x[k]h[n-k]$$

$$y[0] = x[0]h[0]$$

$$y[1] = x[0]h[1] + x[1]h[0]$$

$$y[2] = x[1]h[1] + x[2]h[0]$$

$$y[3] = x[2]h[1] + x[3]h[0]$$

$$y[4] = x[3]h[1] + x[4]h[0]$$

- Circular convolution:  $y[n] = \sum_{k=0}^{N-1} x[k] h[(n-k) \bmod N]$

$$y[0] = x[3]h[1] + x[0]h[0] + x[1]h[3] + x[2]h[2]$$

$$y[1] = x[0]h[1] + x[1]h[0] + x[2]h[3] + x[3]h[2]$$

$$y[2] = x[1]h[1] + x[2]h[0] + x[3]h[3] + x[0]h[2]$$

$$y[3] = x[2]h[1] + x[3]h[0] + x[0]h[3] + x[1]h[2]$$

# Example: (Padding)

- $x[0..3], h[0..1]$
- Linear convolution:

$$y[n] = \sum_{k=-\infty}^{\infty} x[k]h[n-k]$$

$$\begin{aligned}y[0] &= x[0]h[0] \\y[1] &= x[0]h[1] + x[1]h[0] \\y[2] &= x[1]h[1] + x[2]h[0] \\y[3] &= x[2]h[1] + x[3]h[0] \\y[4] &= x[3]h[1] + x[4]h[0]\end{aligned}$$

$N$  is now  $(4 + 2 - 1) = 5$



- Circular convolution:  $y[n] = \sum_{k=0}^{N-1} x[k] h[(n-k) \bmod N]$

$$\begin{aligned}y[0] &= x[4]h[1] + x[0]h[0] + x[1]h[4] + x[2]h[3] + x[3]h[2] \\y[1] &= x[0]h[1] + x[1]h[0] + x[2]h[4] + x[3]h[3] + x[4]h[2] \\y[2] &= x[1]h[1] + x[2]h[0] + x[3]h[4] + x[4]h[3] + x[0]h[2] \\y[3] &= x[2]h[1] + x[3]h[0] + x[4]h[4] + x[0]h[3] + x[1]h[2] \\y[4] &= x[3]h[1] + x[4]h[0] + x[0]h[4] + x[1]h[3] + x[2]h[2]\end{aligned}$$

# Circular Convolution Theorem\*

$$y[n] = \sum_{k=0}^{N-1} x[k] h[(n - k) \bmod N]$$

- Can be calculated by:  $\text{IFFT}(\text{FFT}(x) .* \text{FFT}(h))$
- i.e.

$$\vec{X} = \text{FFT}(\vec{x})$$

$$\vec{H} = \text{FFT}(\vec{h})$$

– For all  $i$ :

$$Y_i = X_i H_i$$

– Then:

$$\vec{y} = \text{IFFT}(\vec{Y})$$

# Circular Convolution Theorem\*

$$y[n] = \sum_{k=0}^{N-1} x[k] h[(n - k) \bmod N]$$

- Can be calculated by:  $IFFT( FFT(x) .* FFT(h) )$
- i.e.

$$\vec{X} = FFT(\vec{x}) \quad O(n \log n) \quad \text{Assume } n > m$$

$$\vec{H} = FFT(\vec{h}) \quad O(m \log m)$$

– For all  $i$ :

$$Y_i = X_i H_i \quad O(n) \quad \text{Total: } O(n \log n)$$

– Then:

$$\vec{y} = IFFT(\vec{Y}) \quad O(n \log n)$$

\* DFT case



# Summary

- Alternate algorithm for large impulse response convolution!
  - Serial:  $O(n \log n)$  vs.  $O(mn)$ 
    - Small vs. large  $m$  determines algorithm choice
    - Runtime does “carry over” to parallel situations (to some extent)

# Homework 3, Part 1

- Implement FFT (“large-kernel”) convolution
  - Use cuFFT for FFT/IFFT (if brave, try your own)
    - ~~Use “batch” variable to save FFT calculations~~  
Correction: Good practice in general, but results in poor performance on Homework 3
    - Don’t forget padding
  - Complex multiplication kernel:
    - Multiply the FFT values pointwise!

# Complex numbers

- `cufftComplex`: cuFFT complex number type
  - Example usage:

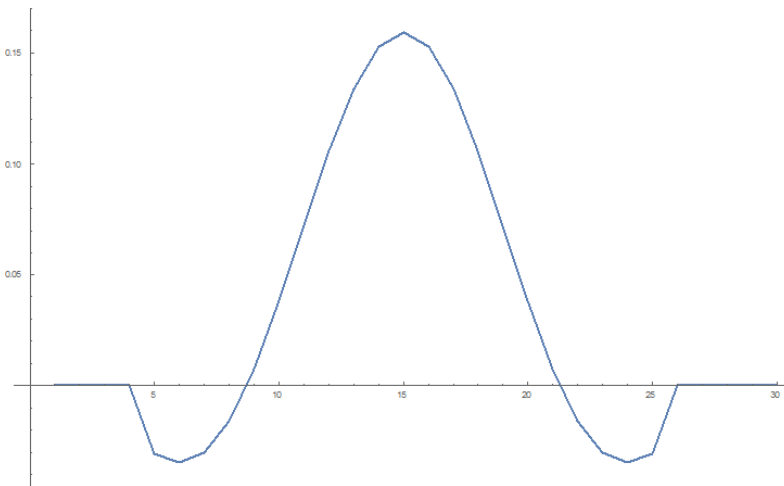
```
cufftComplex a;  
a.x = 3;          // Real part  
a.y = 4;          // Imaginary part
```

- Complex Multiplication:

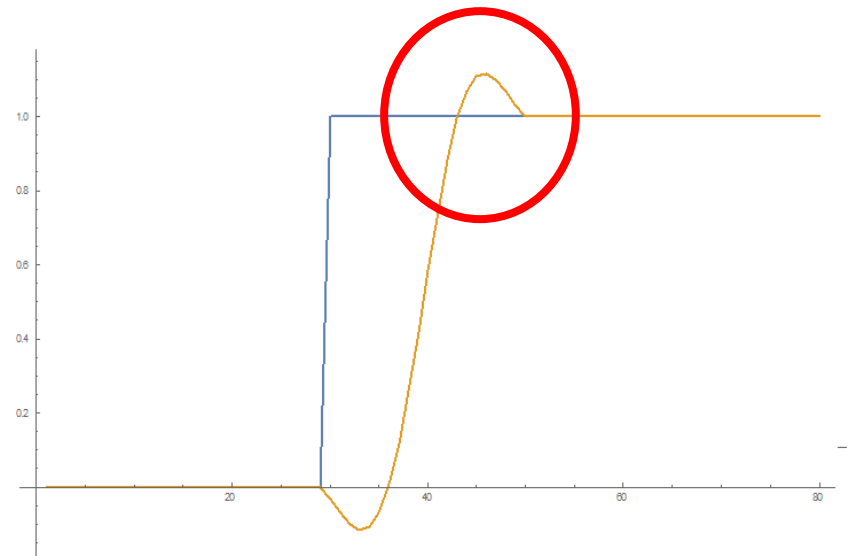
$$(a + bi)(c + di) = (ac - bd) + (ad + bc)i$$

# Homework 3, Part 2

- For a normalized Gaussian filter, output values cannot be larger than the largest input
- Not true in general



Low Pass Filter



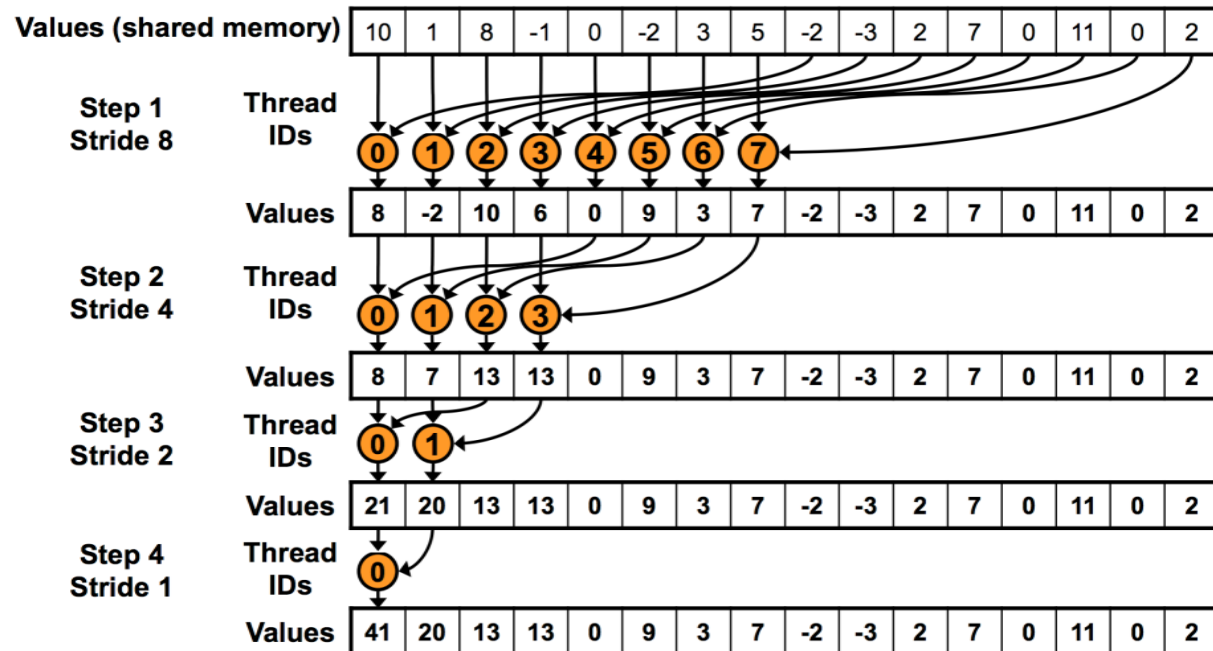
Test Signal and Response

# Normalization

- Amplitudes must lie in range  $[-1, 1]$ 
  - Normalize s.t. maximum magnitude is 1 (or  $1 - \epsilon$ )
- How to find maximum amplitude?

# Reduction

- This time, maximum (instead of sum)
  - Lecture 7 strategies
  - “Optimizing Parallel Reduction in CUDA” (Harris)



# Homework 3, Part 2

- Implement GPU-accelerated normalization
  - Find maximum (reduction)
    - The max amplitude may be a negative sample
  - Divide by maximum to normalize

# (Demonstration)

- Rooms can be modeled as LTI systems!



# Other notes

- Machines:
  - Normal mode: titan, haru, mako, mx, minuteman
  - Audio mode: titan, haru, mako
  
- Due date:
  - Wednesday (4/25), 3 PM