# 1D Convolution & Shared Memory

By Nicolas Agostini
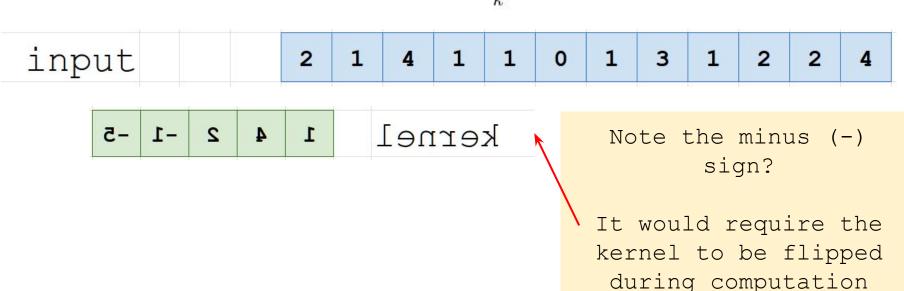
This is the **discrete convolution** formula

$$y[i] = input[i] * kernel[i] = \sum_{k} input[k] \cdot kernel[i - k]$$

This is the **discrete convolution** formula

$$y[i] = input[i] * kernel[i] = \sum_k input[k] \cdot kernel[i - k]$$

| input | | | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| kernel | | 1 | 4 | 2 | -1 | -5 |
|---|---|---|---|---|---|---|

This is the **discrete convolution** formula

$$y[i] = input[i] * kernel[i] = \sum_{k} input[k] \cdot kernel[i - k]$$

| input | | | | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| kernel | | 1 | 4 | 2 | -1 | -5 |
|---|---|---|---|---|---|---|

Note the minus (-) sign?

This is the **discrete convolution** formula

$$y[i] = input[i] * kernel[i] = \sum_{k} input[k] \cdot kernel[i - k]$$

input

| 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|

kernel

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

Note the minus (-) sign?

It would require the kernel to be flipped during computation

$$y[i] = input[i] * kernel[i] = \sum_k input[k] \cdot kernel[i + k]$$

| input | | | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| kernel | | | 1 | 4 | 2 | -1 | -5 |
|---|---|---|---|---|---|---|---|

Note the plus sign

Kernel does not have to be flipped

$$y[i] = input[i] * kernel[i] = \sum_k input[k] \cdot kernel[i+k]$$

input

| 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|

kernel

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

Outside of signal processing, we will call cross-correlations as convolutions.

It is fine…
Everyone does it

For a given output y[i], let say y[2]

$$y[2] = input[2] * kernel[2] = \sum_{k} input[k] \cdot kernel[2 + k]$$

```
k     input[k]      kernel[2+k]
-10   ??            ??
```

| input | | | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|-------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| kernel | 1 | 4 | 2 | -1 | -5 |
|--------|---|---|---|----|----|

For a given output y[i], let say y[2]

$$y[2] = input[2] * kernel[2] = \sum_k input[k] \cdot kernel[2+k]$$

```
k      input[k]      kernel[2+k]
-10    0                 0
```

Accessing outside
of the arrays = 0

| input | | | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|-------|--|--|---|---|---|---|---|---|---|---|---|---|---|---|

| kernel | 1 | 4 | 2 | -1 | -5 |
|--------|---|---|---|----|----|

$$y[2] = input[2] * kernel[2] = \sum_{k} input[k] \cdot kernel[2 + k]$$

```
k      input[k]        kernel[2+k]
-10 0                  0
-3  0                  0
```

| input | | | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| kernel | 1 | 4 | 2 | -1 | -5 |
|---|---|---|---|---|---|

$$y[2] = input[2] * kernel[2] = \sum_k input[k] \cdot kernel[2 + k]$$

| k | input[k] | kernel[2+k] |
|---|---|---|
| -10 | 0 | 0 |
| -3 | 0 | 0 |
| -2 | 0 | 1 |

Value of kernel[2-2]=kernel[0]

| input | | | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| kernel | 1 | 4 | 2 | -1 | -5 |
|---|---|---|---|---|---|

$$y[2] = input[2] * kernel[2] = \sum_k input[k] \cdot kernel[2 + k]$$

| k | input[k] | kernel[2+k] |
|---|----------|-------------|
| -10 | 0 | 0 |
| -3 | 0 | 0 |
| -2 | 0 | 1 |
| -1 | 0 | 4 |

| input | | | | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|-------|--|--|--|---|---|---|---|---|---|---|---|---|---|---|---|

| kernel | | 1 | 4 | 2 | -1 | -5 |
|--------|--|---|---|---|----|----|

$$y[2] = input[2] * kernel[2] = \sum_{k} input[k] \cdot kernel[2+k]$$

| k | input[k] | kernel[2+k] |
|---|----------|-------------|
| -10 | 0 | 0 |
| -3 | 0 | 0 |
| -2 | 0 | 1 |
| -1 | 0 | 4 |
| 0 | 2 | 2 |

input

| 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|

kernel

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

$$y[2] = input[2] * kernel[2] = \sum_k input[k] \cdot kernel[2 + k]$$

| k | input[k] | kernel[2+k] |
|---|---|---|
| -10 | 0 | 0 |
| -3 | 0 | 0 |
| -2 | 0 | 1 |
| -1 | 0 | 4 |
| 0 | 2 | 2 |
| 1 | 1 | -1 |

input

| 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|

kernel

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

$$y[2] = input[2] * kernel[2] = \sum_k input[k] \cdot kernel[2 + k]$$

| k   | input[k] | kernel[2+k] |
|-----|----------|-------------|
| -10 | 0        | 0           |
| -3  | 0        | 0           |
| -2  | 0        | 1           |
| -1  | 0        | 4           |
| 0   | 2        | 2           |
| 1   | 1        | -1          |
| 2   | 4        | -5          |

input
| 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|

kernel
| 1 | 4 | 2 | -1 | -5 |
|---|---|---|----|----|

$$y[2] = input[2] * kernel[2] = \sum_k input[k] \cdot kernel[2 + k]$$

| k | input[k] | kernel[2+k] |
|---|---|---|
| -10 | 0 | 0 |
| -3 | 0 | 0 |
| -2 | 0 | 1 |
| -1 | 0 | 4 |
| 0 | 2 | 2 |
| 1 | 1 | -1 |
| 2 | 4 | -5 |
| 3 | 1 | 0 |
| 4 | 1 | 0 |
| 5 | 0 | 0 |

...

input

| 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|

kernel

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

$$y[2] = input[2] * kernel[2] = \sum_{k} input[k] \cdot kernel[2 + k]$$

Σ

| input[k] | | kernel[2+k] |
|---|---|---|
| 0 | X | 0 |
| 0 | X | 0 |
| 0 | X | 1 |
| 0 | X | 4 |
| 2 | X | 2 |
| 1 | X | -1 |
| 4 | X | -5 |
| 1 | X | 0 |
| 1 | X | 0 |
| 0 | X | 0 |
| ... | | |
| | ... | |

$$y[2] = input[2] * kernel[2] = \sum_k input[k] \cdot kernel[2+k]$$

input[k]      kernel[2+k]

$\sum$

0
0
0
0
**4**
**-1**
**-20**
0
0
0

$$y[2] = input[2] * kernel[2] = \sum_k input[k] \cdot kernel[2 + k]$$

$$y[2] = -17$$

**input[k]**          **kernel[2+k]**

$$\sum$$

0
0
0
0
**4**
**-1**
**-20**
0
0
0

But what are we doing visually?

$$\sum_j a[j] \cdot b[j]$$

input

| 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|

kernel

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

output          -17

# Perform the dot product of the areas that overlap

$$\sum_j a[j] \cdot b[j]$$



| | | | input | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |

kernel: 1 4 2 -1 -5

output: -17

# Slide the filter and repeat

| input | | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| kernel | 1 | 4 | 2 | -1 | -5 |
|---|---|---|---|---|---|

output   1

Now 4 elements overlap

| | | input | | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |

kernel: 1  4  2  -1  -5

output: 1

# Now 5 elements overlap

| input | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |

| kernel | | | | |
|---|---|---|---|---|
| 1 | 4 | 2 | -1 | -5 |

output  8

input

| 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|

kernel

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

output

18

input

| 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|

kernel

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

output

5

And so on…
On GPU: each thread calculates one output value

| input | | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |

| kernel | | 1 | 4 | 2 | -1 | -5 |

output    -11

| input | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| kernel | 1 | 4 | 2 | -1 | -5 |
|---|---|---|---|---|---|

| output | -17 | 1 | 8 | 18 | 5 | -11 | -5 | -1 | -3 | -11 | 9 | 18 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

# The "RADIUS"?

input

| 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|

kernel

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

output  -17

# The "RADIUS" -> padding with zeros

input    | 0 0 | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |

kernel    | 1 | 4 | 2 | -1 | -5 |

output    -17

The "RADIUS" -> padding with zeros
Must do this at the beginning and at the end

input | 0 0 | 2 1 4 1 1 0 1 3 1 2 2 4 | 0 0

kernel | 1 4 2 -1 -5

output | 18

| input | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |

kernel

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

output

| 8 | 18 | 5 | -11 | -5 | -1 | -3 | -11 |
|---|---|---|---|---|---|---|---|

# How about shared memory?

blockDim.x = 8

input    0  0  | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 |  ...

Each block must allocate the
blockDim.x + 2 times RADIUS

How about shared memory?

blockIdx.x = **0**
blockDim.x = 8

input   0  0  2  1  4  1  1  0  1  3  1  2  2  4  …

Each block must copy into shared memory
elements in the <u>orange</u> **and** <u>red</u> boxes

# How about shared memory?

blockIdx.x = **1**
blockDim.x = 8

input | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 | 1 | 2 | 2 | 4 | 2 | 1 | 4 | 1 | 1 | 0 | 1 | 3 |

Each block must copy into shared memory
elements in the <u>orange</u> **and** <u>red</u> boxes

# Carry On!