

# Embedded Systems (ECE340)

## Project 3

Χριστόδουλος Ζερδαλής & Τσιαντός Δημήτριος  
(3531 & 3796)

# **x86 Software implementation**

# Base implementation

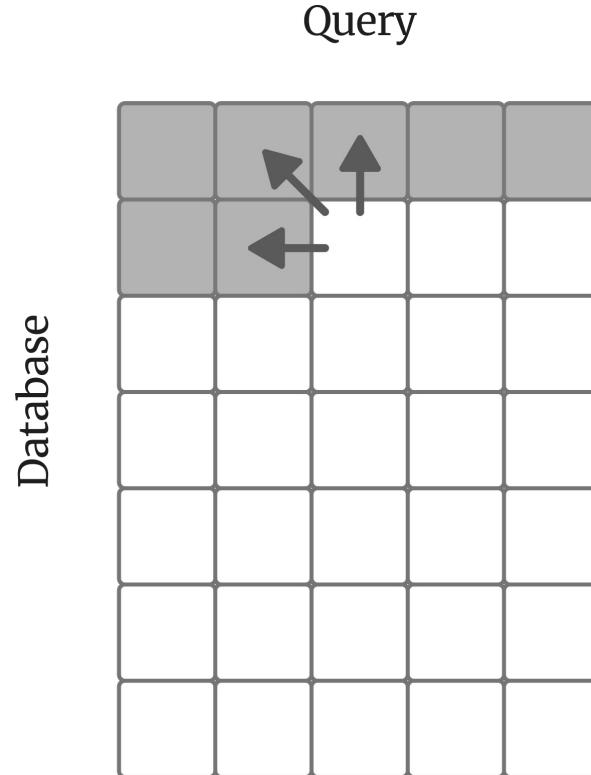
---

**Algorithm 1** Smith-Waterman Algorithm base

---

```
1: for  $i \leftarrow 1$  to  $M$  do
2:   for  $j \leftarrow 1$  to  $N$  do
3:     match  $\leftarrow$   $\begin{cases} \text{MATCH}, & \text{if } Q[j] = D[i] \\ \text{MISS}, & \text{else} \end{cases}$ 
4:     North  $\leftarrow S[i - 1][j] - 1$ 
5:     West  $\leftarrow S[i][j - 1] - 1$ 
6:     Northwest  $\leftarrow S[i - 1][j - 1] + \text{match}$ 
7:      $S[i][j] \leftarrow \max(0, \text{North}, \text{West}, \text{Northwest})$ 
8:     if  $S[i][j] > \text{MAX}$  then
9:       MAX  $\leftarrow S[i][j]$ 
10:      MAX_POS  $\leftarrow (i, j)$ 
11:    end if
12:  end for
13: end for
```

---



# OpenMP implementation

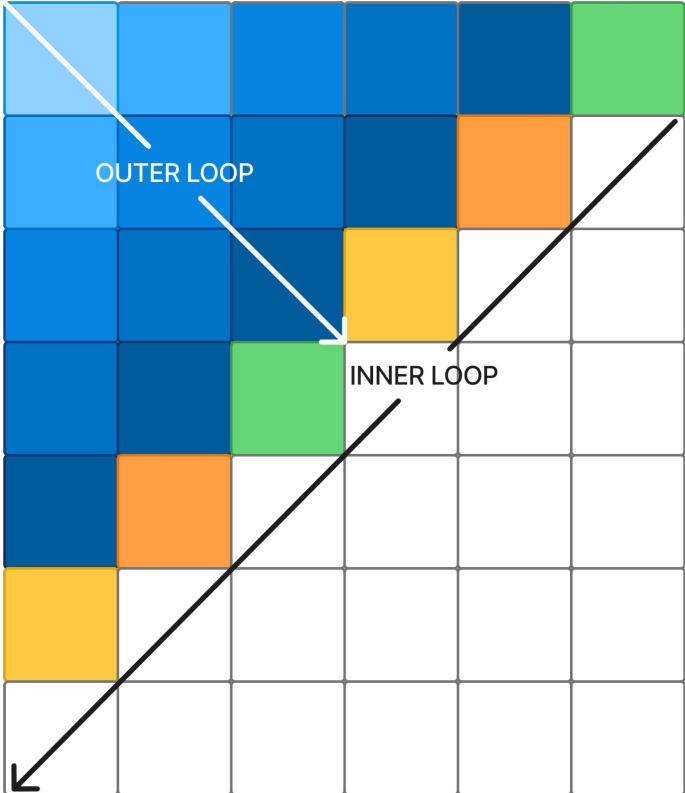
---

## Algorithm 2 Parallel Smith-Waterman with OpenMP

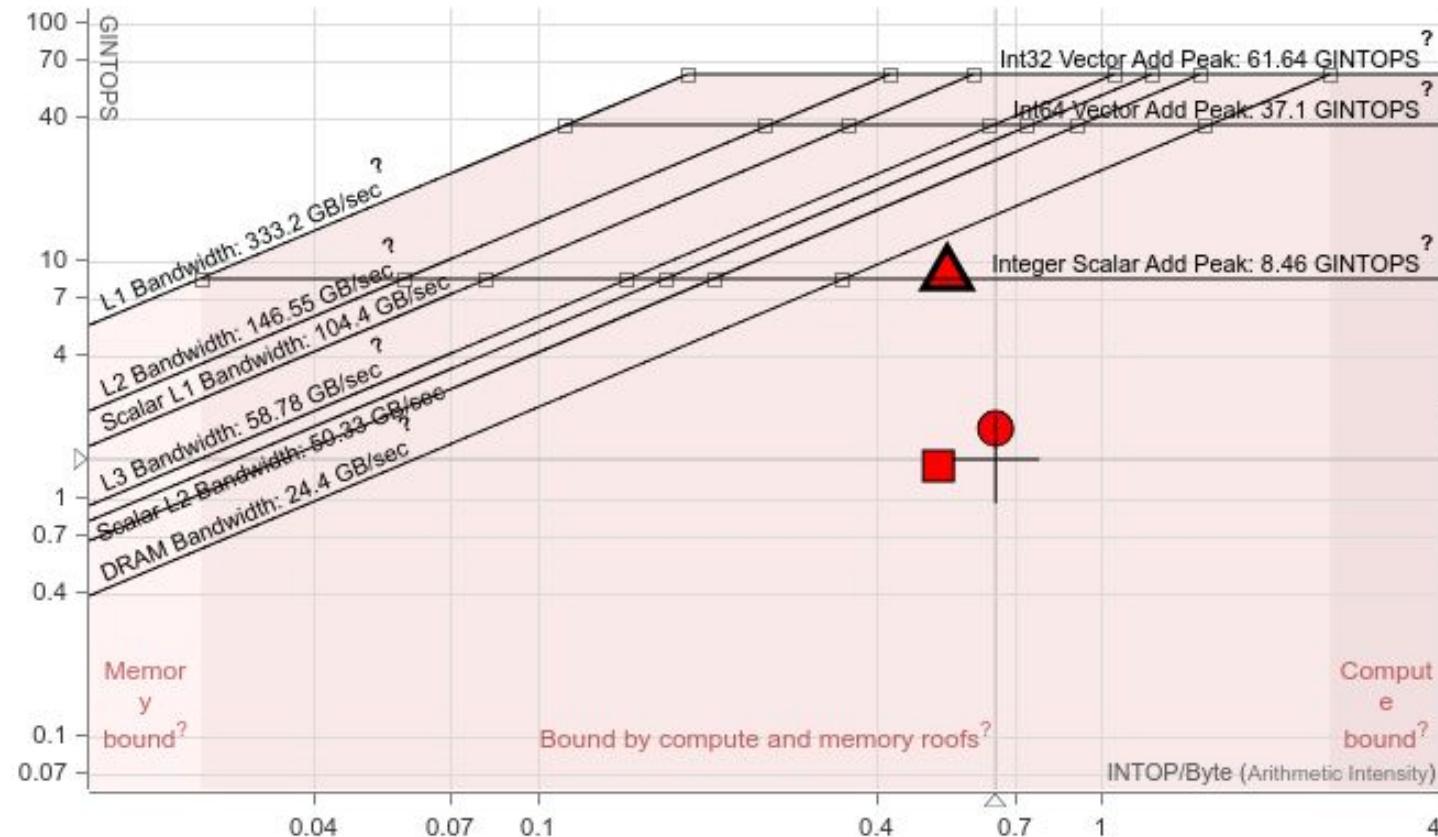
---

```
1: #pragma omp parallel
2: for sum ← 0 to  $M + N - 2$  do
3:   Set  $k \leftarrow \max(1, \text{sum} - N + 1)$ 
4:   #pragma omp for
5:   for  $i \leftarrow k$  to  $N + k - 1$  do
6:     Set  $j \leftarrow \text{sum} - i$ 
7:     ...
8:     if MAX > local_max then
9:       L_MAX ← MAX
10:      L_MAX_index ←  $(i, j)$ 
11:    end if
12:  end for
13: end for
14: #pragma omp critical
15: if L_MAX > G_MAX then
16:   G_MAX ← L_MAX
17:   MAX_index ← L_MAX_index
18: end if
```

---



# Roofline model



# **ARM Software implementation**

# Roofline model

Roof y-axis:

$$\text{Peak INTOPS/s} = f_{\text{CPU}} \times \left( \frac{\#\text{Integer Ops}}{\text{Cycle}} \right)$$

Roof x-axis:

$$\text{RAM Bandwidth} = \text{Bus Width (Bytes)} \times f_{\text{RAM}} \times 2$$

# Roofline model

Roof y-axis:

$$\text{Peak INTOPS} = 0.667 \times 2 = 1.334 \text{ GINTOPS}$$

Roof x-axis:

$$\text{Bandwidth} = 8 \text{ bytes} \times 867 \times 10^6 \text{ Hz} \times 2 = 13.872 \times 10^9 \text{ bytes/s} = 13.872 \text{ GB/s}$$

# Roofline model

Program y-axis:

$$\text{Operational Intensity} = \frac{\#\text{Integer Instructions}}{\#\text{Memory Loads} \times \text{Bus Width (Bytes)}}$$

Program x-axis:

$$\text{GINTOPS} = \frac{\#\text{Integer Instructions}}{10^9 \times \text{Total Time (seconds)}}$$

# Roofline model

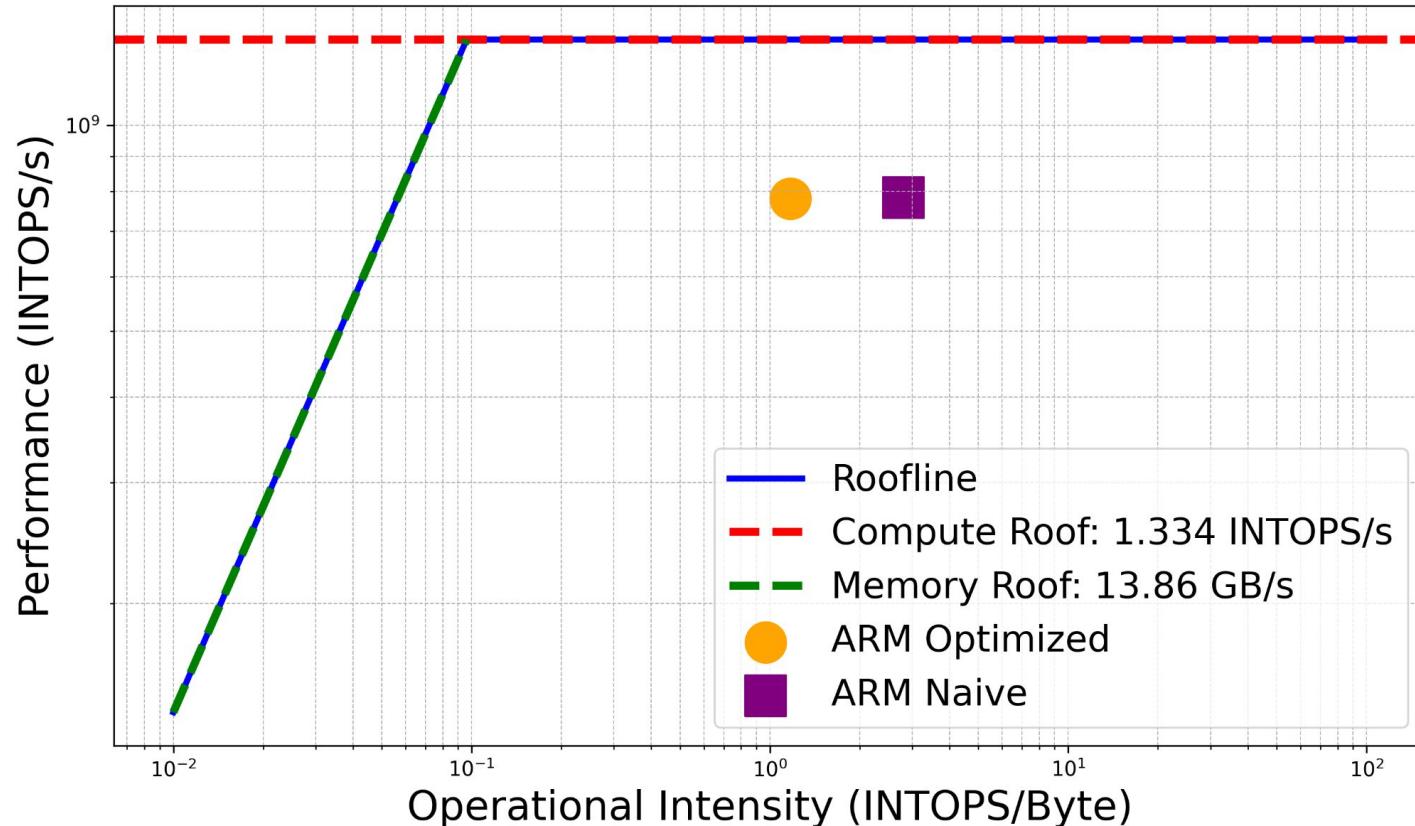
Program y-axis:

$$OI = \frac{\# \text{Instructions}}{\# \text{Loads} \times \text{Bus Width}} = \frac{3,099,341,072}{662,432,510 \times 8} = 0.585 \text{ INTOPS/Byte}$$

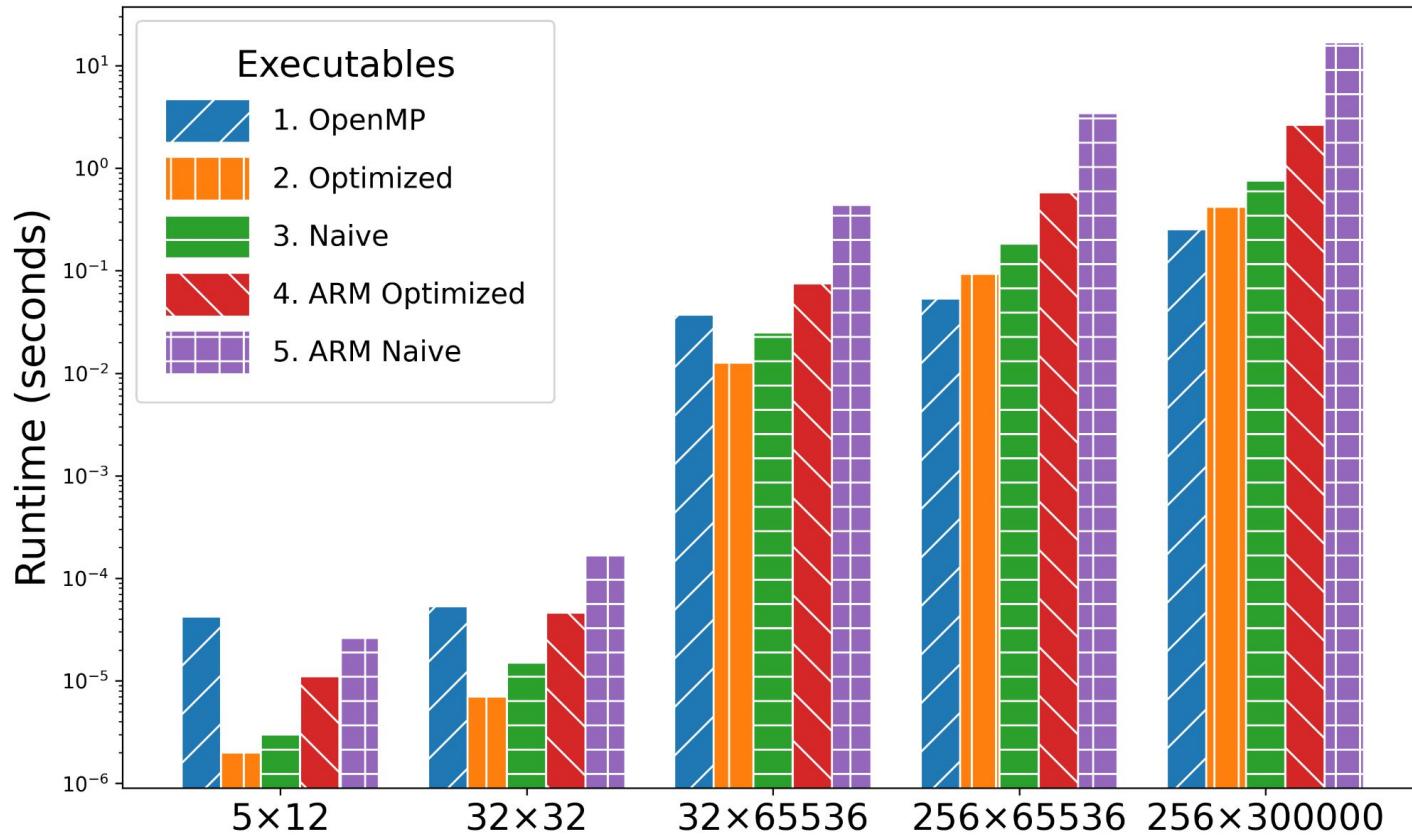
Program x-axis:

$$\text{GINTOPS} = \frac{\# \text{Instructions}}{10^9 \times \text{Time}} = \frac{3,099,341,072}{10^9 \times 3.965} \approx 0.782 \text{ GINTOPS}$$

# Roofline model



# Execution time comparisons



# **Hardware implementation**

# Base implementation on FPGA

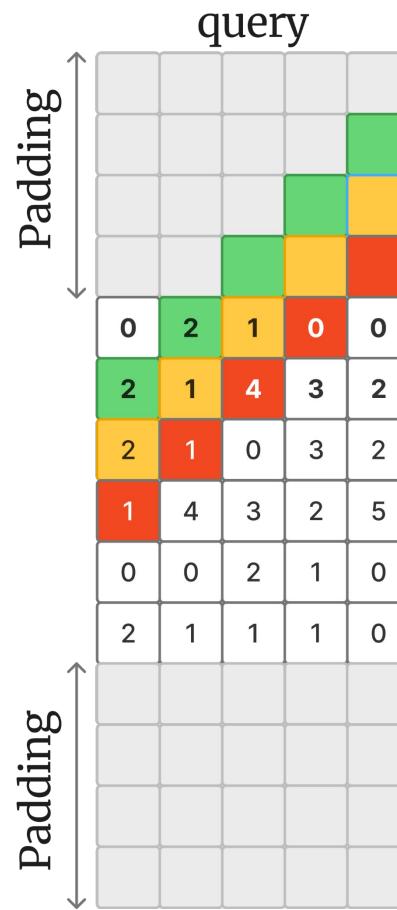
---

<b>Latency on FPGA (ms)</b>	4562.5097
<b>Loop Name</b>	OUTER_LOOP
<b>Min Latency (cycles)</b>	297088813
<b>Max Latency (cycles)</b>	435367289
<b>Latency Range</b>	4529 – 6637
<b>Iteration Latency</b>	–
<b>Initiation Achieved</b>	–
<b>Initiation Target</b>	–
<b>Trip Count</b>	65597
<b>Pipelined</b>	no

---

# First hardware optimizations

- Remove **similarity matrix**
- Add **buffers**
- Add **padding**
- **Memory** optimizations
- **Pipeline** and **unrolling**



# First hardware optimizations

```
1: for  $k \leftarrow 0$  to  $M - N$  do (Outer_loop)
2:   #pragma HLS PIPELINE
3:   for  $i \leftarrow 0$  to  $N - 1$  do (Inner_loop)
4:      $j \leftarrow N - 1 - i$ 
5:     match  $\leftarrow (\text{query\_buff}[j] = \text{database\_buff}[i + k]) ? \text{MATCH} : \text{MISS\_MATCH}$ 
6:     North  $\leftarrow \text{diag\_array\_2}[i] - 1$ 
7:     Northwest  $\leftarrow \text{match} + \text{diag\_array\_1}[i + 1]$ 
8:     West  $\leftarrow \text{diag\_array\_2}[i + 1] - 1$ 
9:     ...
10:    Direction_buff[i]  $\leftarrow \text{Direction}$ 
11:    Diag_array_3[i]  $\leftarrow \text{Max\_value}$ 
12:    ...
13:  end for
14:  Copy diag_array_2 to diag_array_1
15:  Copy diag_array_3 to diag_array_2
16:  for  $i \leftarrow 0$  to  $N - 1$  do (Dir_loop)
17:     $j \leftarrow N - 1 - i$ 
18:    direction_matrix[( $i + k$ )  $\times N + j] \leftarrow \text{direction\_buff}[i]$ 
19:  end for
20: end for
```

# Memory optimizations

```
#pragma HLS ARRAY_PARTITION diag_array_1 dim=1 factor=32 cyclic
#pragma HLS BIND_STORAGE diag_array_1 type=ram_t2p impl=bram

#pragma HLS ARRAY_PARTITION diag_array_2 dim=1 factor=32 cyclic
#pragma HLS BIND_STORAGE diag_array_2 type=ram_t2p impl=bram

#pragma HLS ARRAY_PARTITION diag_array_3 dim=1 factor=32 cyclic
#pragma HLS BIND_STORAGE diag_array_3 type=ram_t2p impl=bram

#pragma HLS ARRAY_PARTITION query_buff dim=1 factor=16 cyclic
#pragma HLS BIND_STORAGE query_buff type=ram_1wnr impl=bram

#pragma HLS ARRAY_PARTITION database_buff dim=1 factor=16 cyclic
#pragma HLS BIND_STORAGE database_buff type=ram_1wnr impl=bram

#pragma HLS ARRAY_PARTITION direction_buff dim=1 factor=16 cyclic
#pragma HLS BIND_STORAGE direction_buff type=ram_t2p impl=bram
```

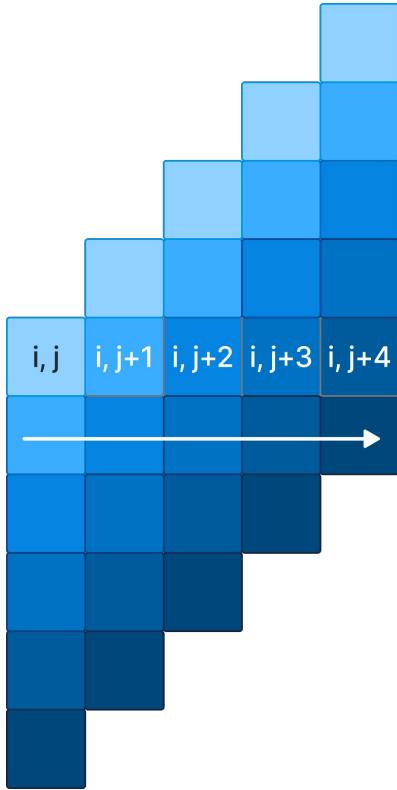
# First hardware optimizations

---

<b>Latency on FPGA (ms)</b>	21.641
<b>Loop Name</b>	OUTER_LOOP
<b>Min Latency (cycles)</b>	2,098,215
<b>Max Latency (cycles)</b>	2,098,215
<b>Iteration Latency</b>	104
<b>Initiation Achieved</b>	32
<b>Initiation Target</b>	1
<b>Trip Count</b>	65,567
<b>Pipelined</b>	yes

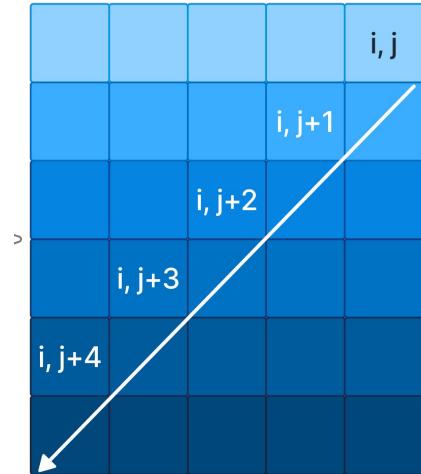
---

# Flattening the direction array



```
for  $i \leftarrow 0$  to  $N - 1$  do (Dir_loop)  
     $j \leftarrow N - 1 - i$   
    direction_matrix[ $(i + k) \times N + j$ ]  $\leftarrow$  direction_buff[i]  
end for
```

# Flattening the direction array



```
for  $j \leftarrow 0$  to  $N - 1$  do (Dir_loop)  
    direction_matrix[ $i * N + j$ ]  $\leftarrow$  direction_buff[ $j$ ]  
end for
```

# Flattening the direction array

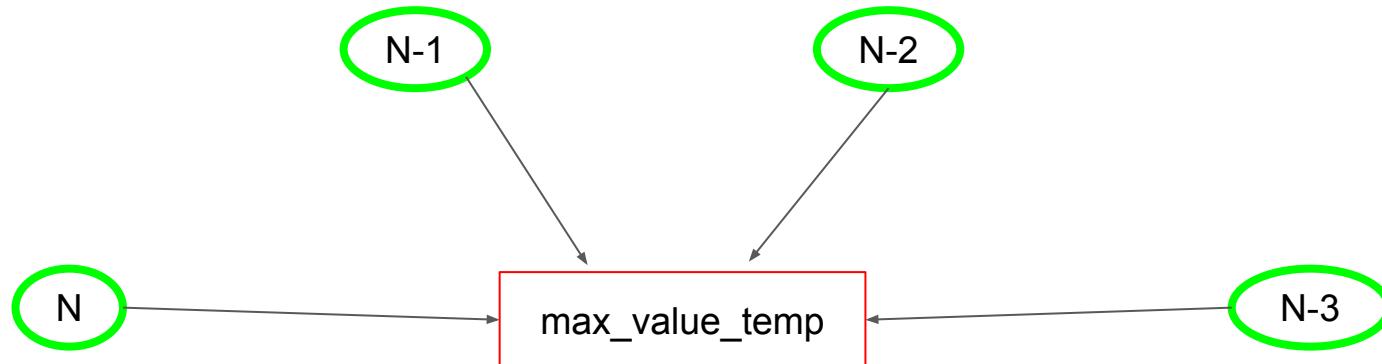
---

<b>Latency on FPGA (ms)</b>	12.8369
<b>Loop Name</b>	OUTER_LOOP
<b>Min Latency (cycles)</b>	1,049,075
<b>Max Latency (cycles)</b>	1,049,075
<b>Iteration Latency</b>	20
<b>Initiation Achieved</b>	16
<b>Initiation Target</b>	1
<b>Trip Count</b>	65,567
<b>Pipelined</b>	yes

---

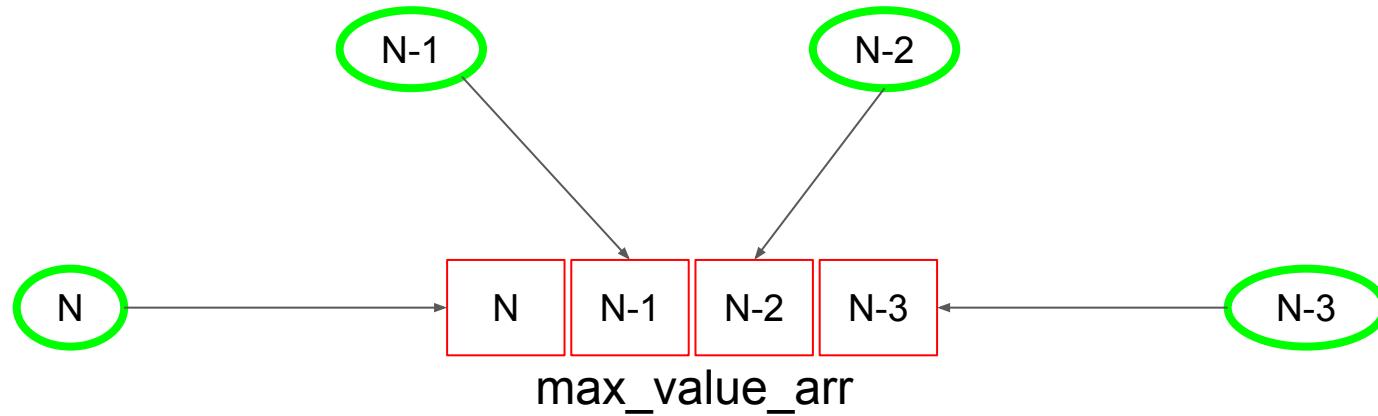
# Adding max value array

```
if max_value > max_value_temp then  
    max_value_temp ← max_value  
    max_index_temp ← ((i + k) * N + j)  
end if
```



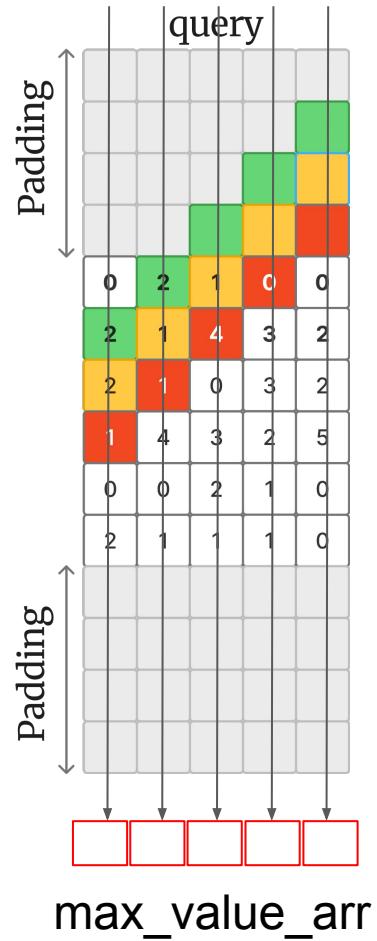
# Adding max value array

```
if max_value > max_value_arr[i] then  
    max_value_arr[i] ← max_value  
    max_index_arr[i] ← ((i + k) * N + j)  
end if
```



# Adding max value array

```
if max_value > max_value_arr[i] then  
    max_value_arr[i] ← max_value  
    max_index_arr[i] ← ((i + k) * N + j)  
end if
```



## Adding max value array

```
for  $i \leftarrow 0$  to  $N - 1$  do (MAX_LOOP)
    if  $\text{max\_value\_arr}[i] > \text{max\_value\_temp}$  then
         $\text{max\_value\_temp} \leftarrow \text{max\_value\_arr}[i]$ 
         $\text{max\_idx\_temp} \leftarrow \text{max\_index\_arr}[i]$ 
    end if
end for
* $\text{max\_index} \leftarrow \text{max\_idx\_temp}$ 
```

Only one write to main memory

# Adding max value array

```
#pragma HLS ARRAY_PARTITION variable=max_value_arr dim=1 factor=16 cyclic
#pragma HLS BIND_STORAGE variable=max_value_arr type=ram_t2p impl=bram

#pragma HLS ARRAY_PARTITION variable=max_index_arr dim=1 factor=16 cyclic
#pragma HLS BIND_STORAGE variable=max_index_arr type=ram_t2p impl=bram
```

# Adding max value array

---

<b>Latency on FPGA (ms)</b>	4.96728
<b>Loop Name</b>	OUTER_LOOP
<b>Min Latency (cycles)</b>	262,268
<b>Max Latency (cycles)</b>	262,268
<b>Iteration Latency</b>	4
<b>Initiation Achieved</b>	4
<b>Initiation Target</b>	1
<b>Trip Count</b>	65,567
<b>Pipelined</b>	yes

---

# Optimizing if statements

```
1: max_value ← north
2: direction ← NORTH
3: if northwest > max_value then
4:     max_value ← northwest
5:     direction ← NORTH_WEST
6: end if
7: if west > max_value then
8:     max_value ← west
9:     direction ← WEST
10: end if
11: if max_value ≤ 0 then
12:     max_value ← 0
13:     direction ← CENTER
14: end if
```

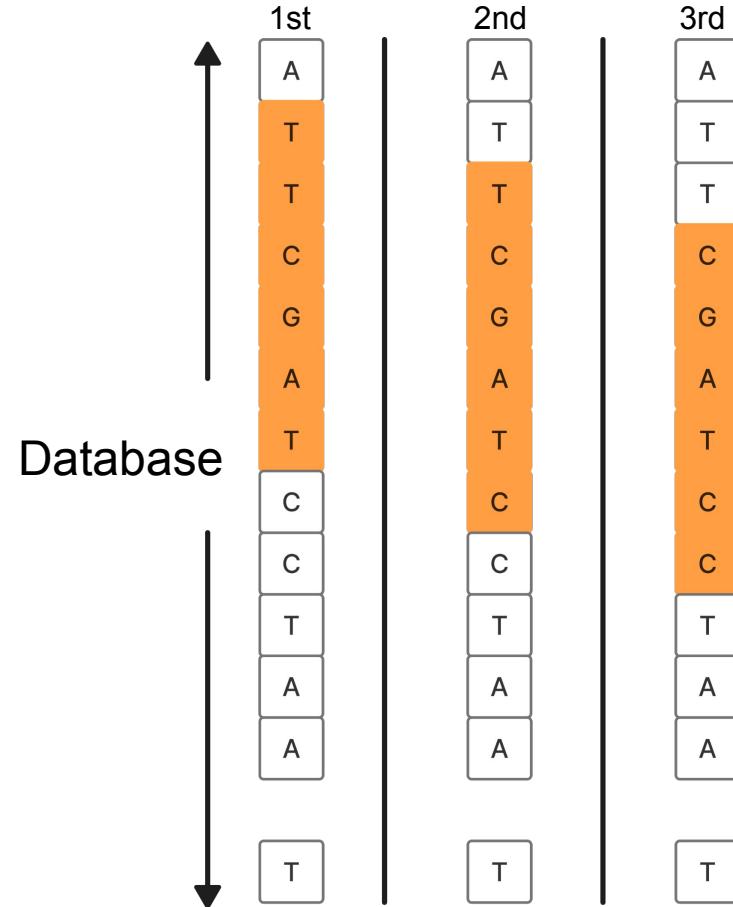
Many if statements

# Optimizing if statements

```
1: if north < northwest and northwest < west and west ≠ -1 then
2:     max_value ← west
3:     direction ← WEST
4: else if north < northwest and northwest ≠ -1 then
5:     max_value ← northwest
6:     direction ← NORTH_WEST
7: else if north < west and west ≠ -1 then
8:     max_value ← west
9:     direction ← WEST
10: else if north = -1 then
11:     max_value ← 0
12:     direction ← CENTER
13: else
14:     max_value ← north
15:     direction ← NORTH
16: end if
```

One if-else-if chain

# Database Buffer Shift



```
1: for  $k \leftarrow 0$  to  $M - N$  do
2:   for  $i \leftarrow 1$  to  $N - 1$  do
3:     database_buf[i - 1]  $\leftarrow$  database_buf[i]
4:   end for
5:   database_buf[N - 1]  $\leftarrow$  database[k + N - 1]
6:   for  $i \leftarrow 0$  to  $N - 1$  do
7:     // ... inner loop body
8:   end for
9: end for
```

# Database Buffer Shift

---

<b>Latency on FPGA (ms)</b>	3.78776
<b>Loop Name</b>	OUTER_LOOP
<b>Min Latency (cycles)</b>	196,841
<b>Max Latency (cycles)</b>	196,841
<b>Iteration Latency</b>	144
<b>Initiation Achieved</b>	3
<b>Initiation Target</b>	1
<b>Trip Count</b>	65,567
<b>Pipelined</b>	yes

---

# Why did the Iteration Interval dropped?

<b>Latency on FPGA (ms)</b>	4.96728	<b>Latency on FPGA (ms)</b>	3.78776
<b>Loop Name</b>	OUTER_LOOP	<b>Loop Name</b>	OUTER_LOOP
<b>Min Latency (cycles)</b>	262,268	<b>Min Latency (cycles)</b>	196,841
<b>Max Latency (cycles)</b>	262,268	<b>Max Latency (cycles)</b>	196,841
<b>Iteration Latency</b>	4	<b>Iteration Latency</b>	144
<b>Initiation Achieved</b>	4	<b>Initiation Achieved</b>	3
<b>Initiation Target</b>	1	<b>Initiation Target</b>	1
<b>Trip Count</b>	65,567	<b>Trip Count</b>	65,567
<b>Pipelined</b>	yes	<b>Pipelined</b>	yes

# Why did the Initiation dropped?

||=4

⚠ [HLS 200-880]

[LINK](#)

The II Violation in module 'compute\_matrices' (loop 'OUTER\_LOOP'): Unable to enforce a carried dependence constraint (II = 1, distance = 1, offset = 1) between 'store' operation ('diag\_array\_2[0].addr\_1.write\_ln84', diag\_dir\_max/[lsal.cpp:84](#)) of variable 'diag\_array\_3.load\_0', diag\_dir\_max/[lsal.cpp:84](#) on array 'diag\_array\_2[0]', diag\_dir\_max/[lsal.cpp:17](#) and 'load' operation ('diag\_array\_2[0].load', diag\_dir\_max/[lsal.cpp:65](#)) on array 'diag\_array\_2[0]', diag\_dir\_max/[lsal.cpp:17](#).

⚠ [HLS 200-880]

[LINK](#)

The II Violation in module 'compute\_matrices' (loop 'OUTER\_LOOP'): Unable to enforce a carried dependence constraint (II = 2, distance = 1, offset = 1) between 'store' operation ('diag\_array\_2[0].addr\_1.write\_ln84', diag\_dir\_max/[lsal.cpp:84](#)) of variable 'diag\_array\_3.load\_0', diag\_dir\_max/[lsal.cpp:84](#) on array 'diag\_array\_2[0]', diag\_dir\_max/[lsal.cpp:17](#) and 'load' operation ('diag\_array\_2[0].load', diag\_dir\_max/[lsal.cpp:65](#)) on array 'diag\_array\_2[0]', diag\_dir\_max/[lsal.cpp:17](#).

⚠ [HLS 200-880]

[LINK](#)

The II Violation in module 'compute\_matrices' (loop 'OUTER\_LOOP'): Unable to enforce a carried dependence constraint (II = 3, distance = 1, offset = 1) between 'store' operation ('diag\_array\_2[0].addr\_1.write\_ln84', diag\_dir\_max/[lsal.cpp:84](#)) of variable 'diag\_array\_3.load\_0', diag\_dir\_max/[lsal.cpp:84](#) on array 'diag\_array\_2[0]', diag\_dir\_max/[lsal.cpp:17](#) and 'load' operation ('diag\_array\_2[0].load', diag\_dir\_max/[lsal.cpp:65](#)) on array 'diag\_array\_2[0]', diag\_dir\_max/[lsal.cpp:17](#).

||=3

⚠ [HLS 200-880]

[LINK](#)

The II Violation in module 'compute\_matrices' (loop 'OUTER\_LOOP'): Unable to enforce a carried dependence constraint (II = 1, distance = 1, offset = 1) between 'store' operation ('diag\_array\_2[0].addr\_1.write\_ln86', storage/[lsal.cpp:86](#)) of variable 'diag\_array\_3.load\_0', storage/[lsal.cpp:86](#) on array 'diag\_array\_2[0]', storage/[lsal.cpp:17](#) and 'load' operation ('diag\_array\_2[0].load', storage/[lsal.cpp:80](#)) on array 'diag\_array\_2[0]', storage/[lsal.cpp:17](#).

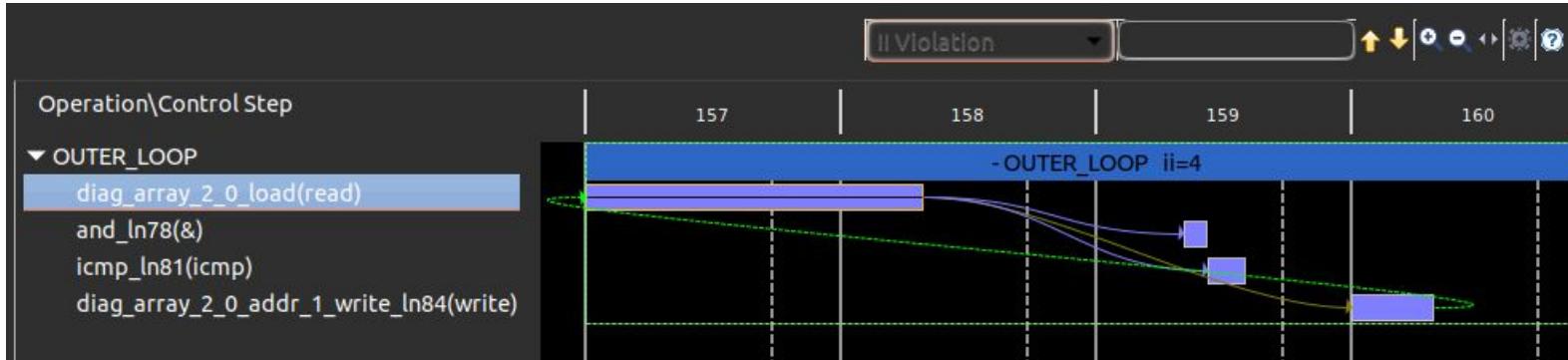
⚠ [HLS 200-880]

[LINK](#)

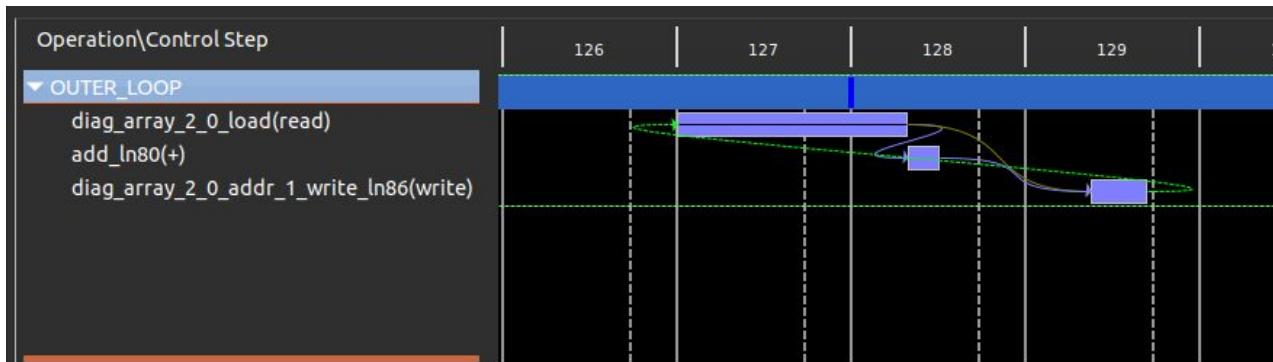
The II Violation in module 'compute\_matrices' (loop 'OUTER\_LOOP'): Unable to enforce a carried dependence constraint (II = 2, distance = 1, offset = 1) between 'store' operation ('diag\_array\_2[0].addr\_1.write\_ln86', storage/[lsal.cpp:86](#)) of variable 'diag\_array\_3.load\_0', storage/[lsal.cpp:86](#) on array 'diag\_array\_2[0]', storage/[lsal.cpp:17](#) and 'load' operation ('diag\_array\_2[0].load', storage/[lsal.cpp:80](#)) on array 'diag\_array\_2[0]', storage/[lsal.cpp:17](#).

# Why did the Initiation dropped?

||=4



||=3



# Diagonal rotation

```
1: for  $k \leftarrow 0$  to  $M - N$  do
2:   for  $i \leftarrow 0$  to  $N - 1$  do
3:     // ... inner loop body
4:   end for
5:   Copy diag_array_2 into diag_array_1
6:   Copy diag_array_3 into diag_array_2
7: end for
```

```
1: for  $k \leftarrow 0$  to  $M - N$  do
2:   Copy diag_array_2 into diag_array_1
3:   Copy diag_array_3 into diag_array_2
4:   for  $i \leftarrow 0$  to  $N - 1$  do
5:     // ... inner loop body
6:   end for
7: end for
```

# Diagonal rotation

```
if max_value > max_value_arr[i] then  
    max_value_arr[i] ← max_value  
    max_index_arr[i] ← ((i + k) * N + j)  
end if
```

N reads (32)

←

N writes (32)

We have ram\_t2p cyclic 16 : NOT ENOUGH

# Diagonal rotation

Copy contents of `diag_array_2` to `diag_array_1`  
Copy contents of `diag_array_3` to `diag_array_2`

...  
`match`  $\leftarrow$  (`query[j]` = `database[i + k]`)?MATCH : MISS\_MATCH

`North`  $\leftarrow$  `diag_array_2[i] - 1`

`Northwest`  $\leftarrow$  `match + diag_array_1[i + 1]`

`West`  $\leftarrow$  `diag_array_2[i + 1] - 1`

...  
`Diag_array_3[i]`  $\leftarrow$  Max\_value

N+1 reads (33)

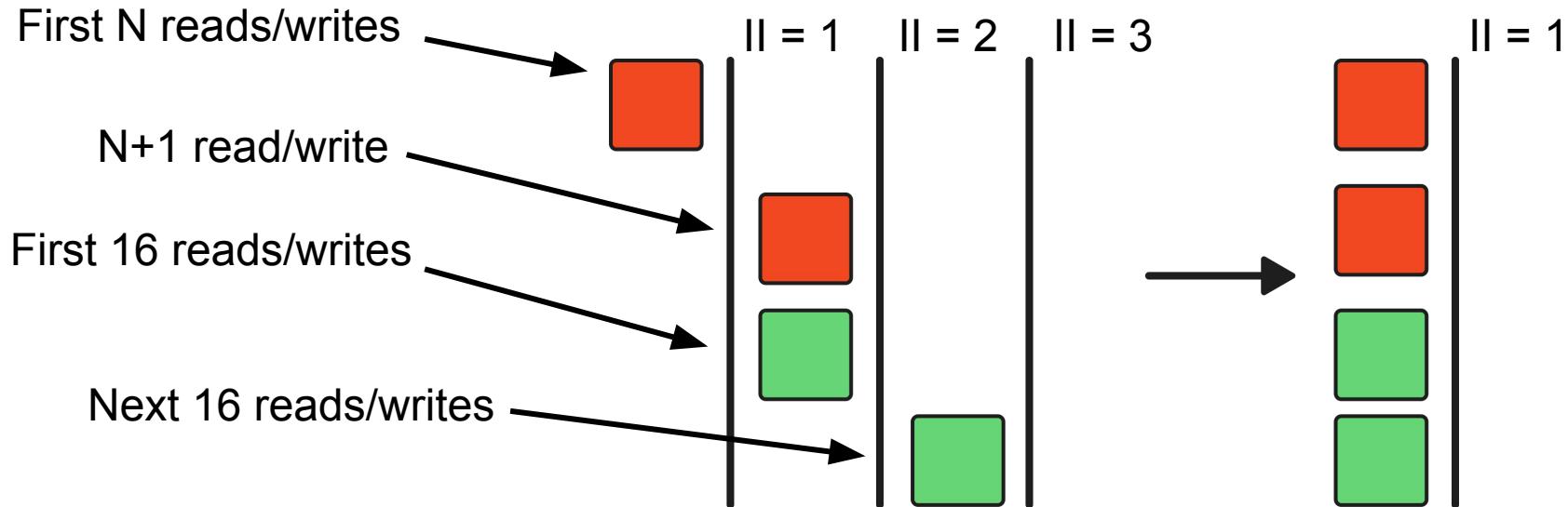
N writes (32)

We have `ram_t2p cyclic 32` : NOT ENOUGH

# Diagonal rotation

```
#pragma HLS ARRAY_PARTITION diag_array_3 dim=1 complete
```

```
#pragma HLS ARRAY_PARTITION max_value_arr dim=1 complete
```



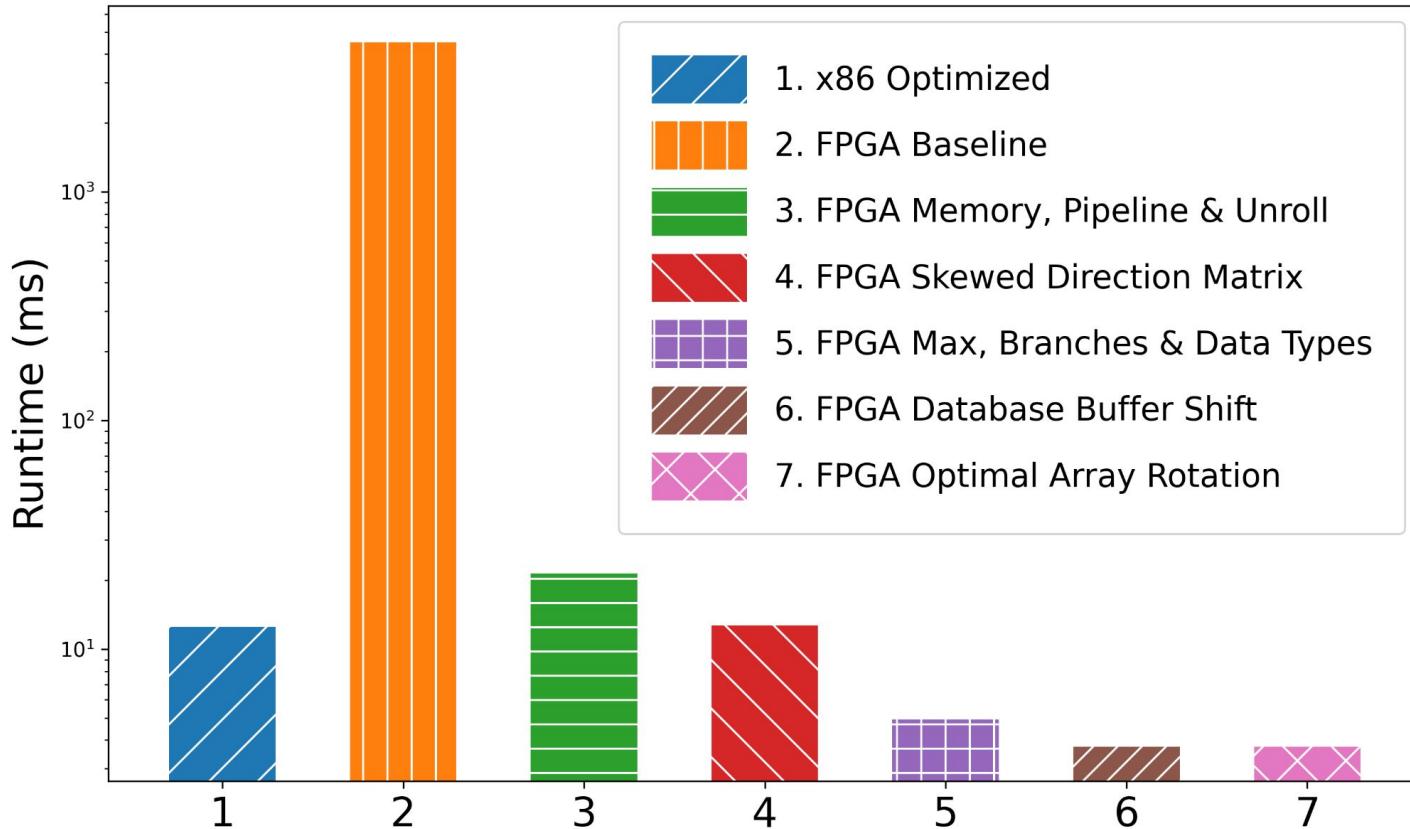
# Diagonal rotation

---

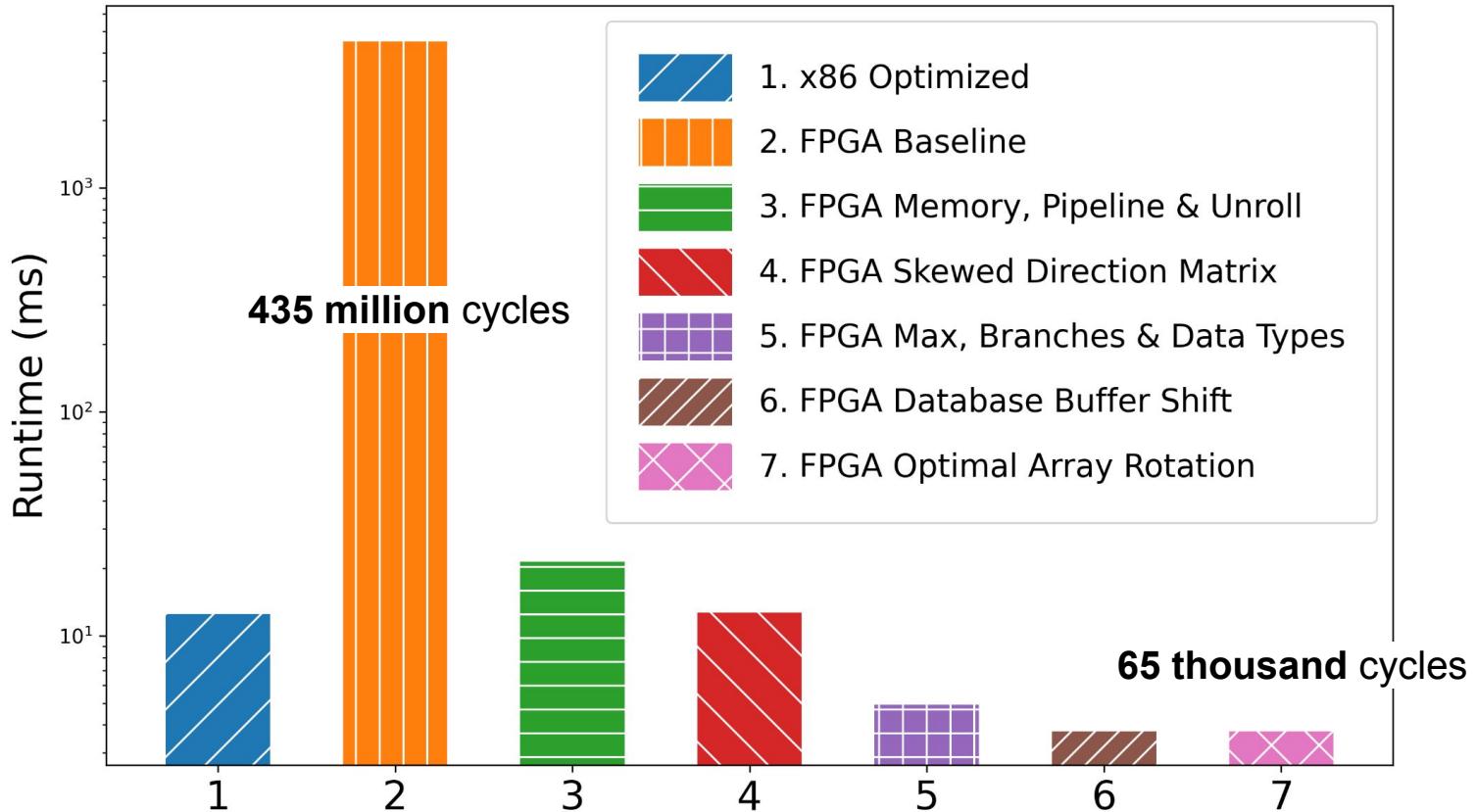
<b>Latency on FPGA(ms)</b>	3.77532
<b>Loop Name</b>	OUTER_LOOP
<b>Min Latency (cycles)</b>	65,740
<b>Max Latency (cycles)</b>	65,740
<b>Iteration Latency</b>	175
<b>Initiation Achieved</b>	1
<b>Initiation Target</b>	1
<b>Trip Count</b>	65,567
<b>Pipelined</b>	yes

---

# Time improvement with different optimizations



# Time improvement with different optimizations



# Time improvement with different optimizations

