# AAML FINAL PROJECT

## GROUP 12

312832011 曾郁涵
312832003 林子泰
312833015 温峻揚
311511022 邱政岡
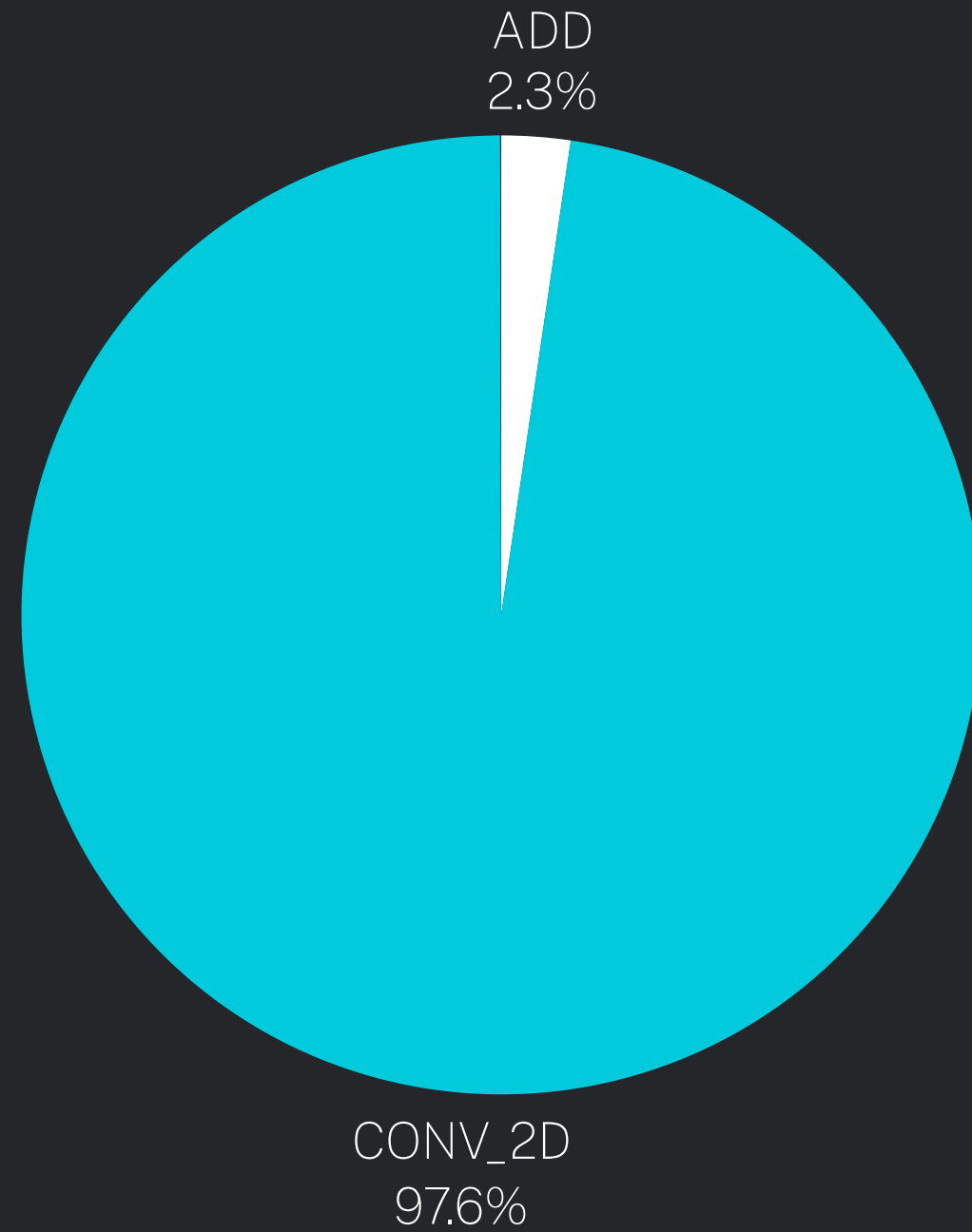
# Performance



Original   100%

Our Method   13.8%

ADD
2.3%

CONV_2D
97.6%

# Execution Time of Layers in Ticks

in MLPerf Tiny

# **Conv2D** Layer Speedup

# Method

**1**    **2**    **3**    **4**

**SIMD**    **Software Optimization**    **Postprocess on FPGA**    **Systolic Array with im2col**

# SIMD

```
                  7 bits
            +--------------+
funct7 =    | (bool) reset |
            +--------------+


            int8_t            int8_t            int8_t            int8_t
            +--------------+--------------+--------------+--------------+
  in0 =     | input_data[0] | input_data[1] | input_data[2] | input_data[3] |
            +--------------+--------------+--------------+--------------+


            int8_t            int8_t            int8_t            int8_t
            +--------------+--------------+--------------+--------------+
  in1 =     | filter_data[0] | filter_data[1] | filter_data[2] | filter_data[3] |
            +--------------+--------------+--------------+--------------+


                                    int32_t
            +---------------------------------------------------------------+
output =    | output + (input_data[0, 1, 2, 3] + offset) * filter_data[0, 1, 2, 3] |
            +---------------------------------------------------------------+
```
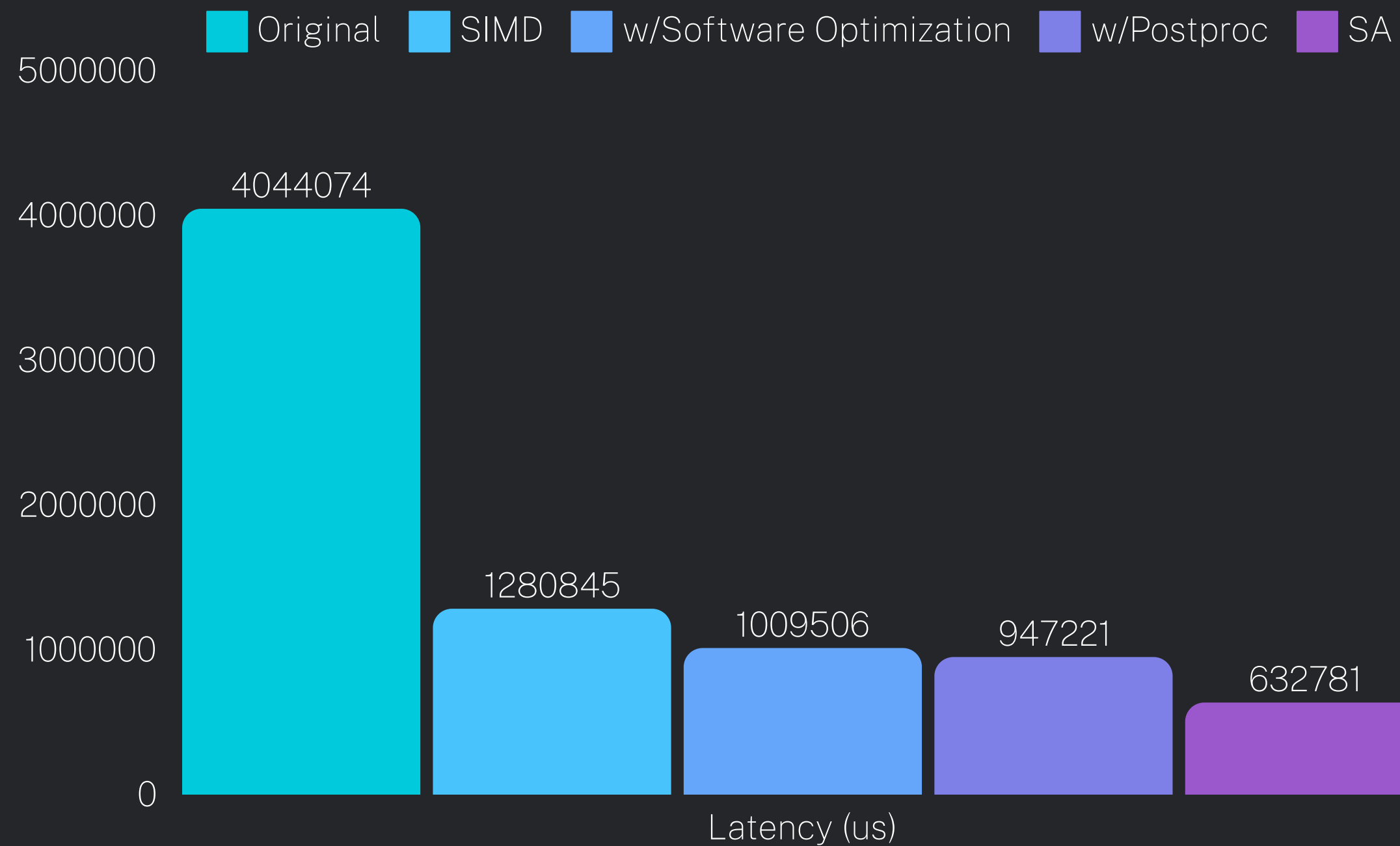
# Software Optimization

1. Hardcode the constant parameters

2. Loop Unrolling

3. Minimize invocation of *Offset* func

4. Remove redundant computations

# Postprocess on FPGA
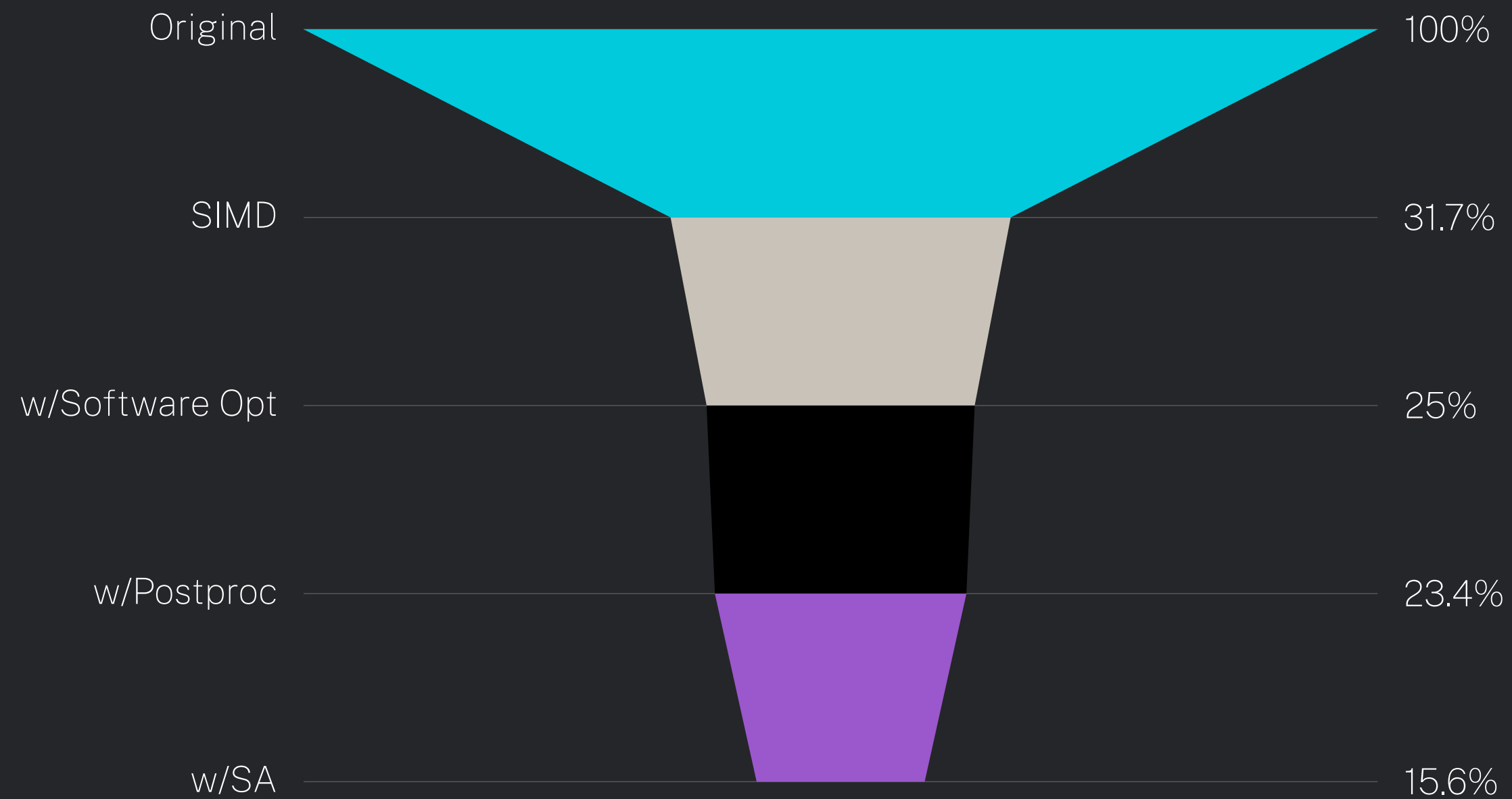
*MultiplyByQuantizedMultiplier*
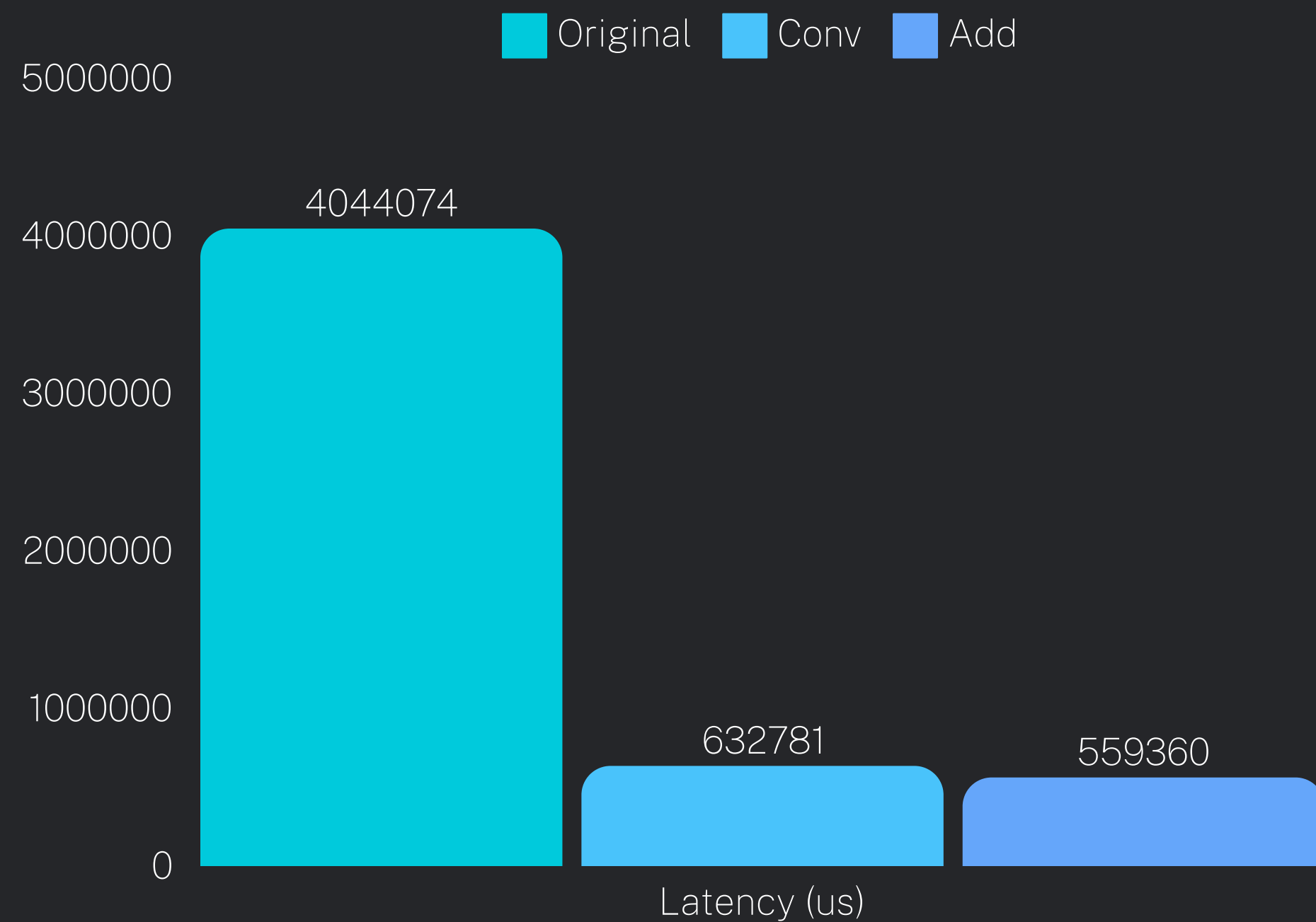
# Systolic Array with im2col

# Performance



| | |
|---|---|
| Original | 100% |
| SIMD | 31.7% |
| w/Software Opt | 25% |
| w/Postproc | 23.4% |
| w/SA | 15.6% |

# ADD Layer Speedup

*SIMD & Postprocess on FPGA*

# Performance



| | |
|---|---|
| Original | 100% |
| Conv | 15.5% |
| ADD | 13.8% |

# Fully Connected Layer Speedup

(Same as Conv2D )

# Performance

# Thank you!