# Pitfalls of UPMEM kernel development

**Heterogeneous Memory Software Lab, Poland**
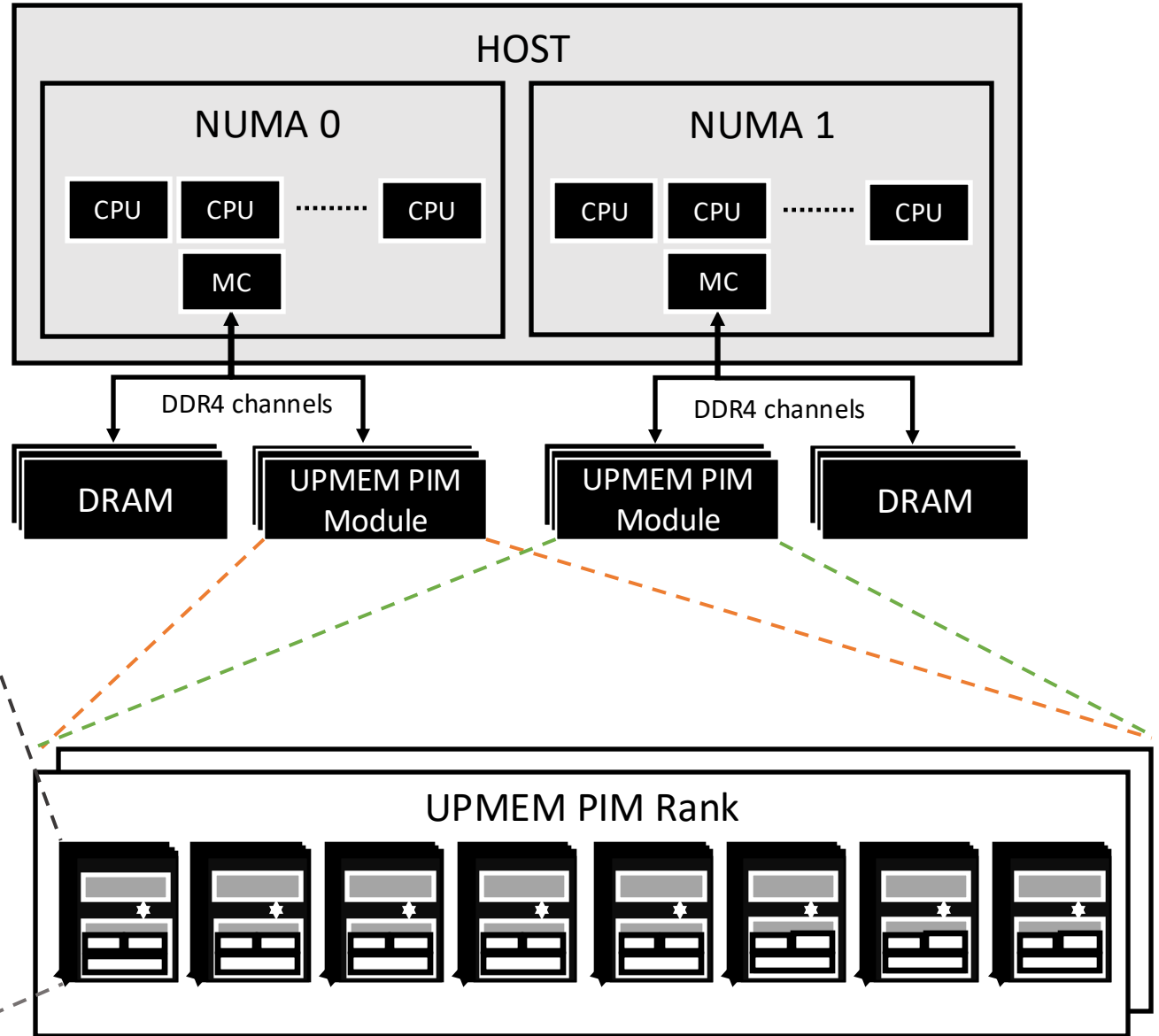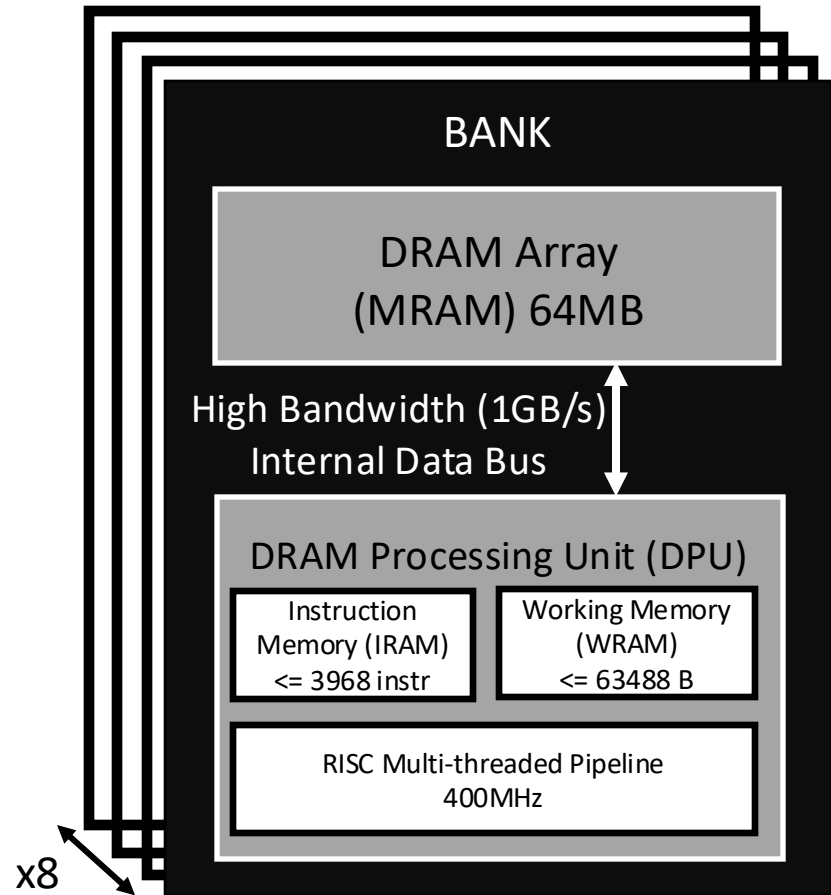
Krystian Chmielewski, Tadeusz Kobus, Paweł Piotrowicz, Jarosław Ławnicki, Uladzislau Lukyanau, Maciej Maciejewski, Zhang Hongjun
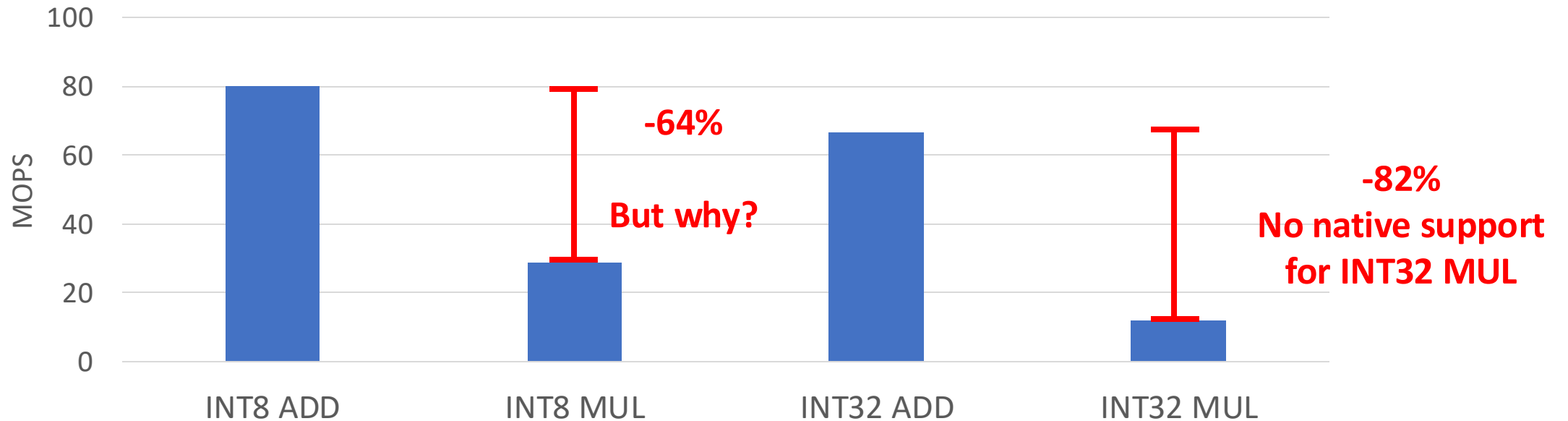
krystian.chmielewski@huawei.com

# What is UPMEM?



BANK

DRAM Array
(MRAM) 64MB

High Bandwidth (1GB/s)
Internal Data Bus

DRAM Processing Unit (DPU)

Instruction
Memory (IRAM)
<= 3968 instr

Working Memory
(WRAM)
<= 63488 B

RISC Multi-threaded Pipeline
400MHz

x8

HOST

NUMA 0

CPU  CPU  ........  CPU

MC

NUMA 1

CPU  CPU  ........  CPU

MC

DDR4 channels

DDR4 channels

DRAM

UPMEM PIM
Module

UPMEM PIM
Module

DRAM

UPMEM PIM Rank
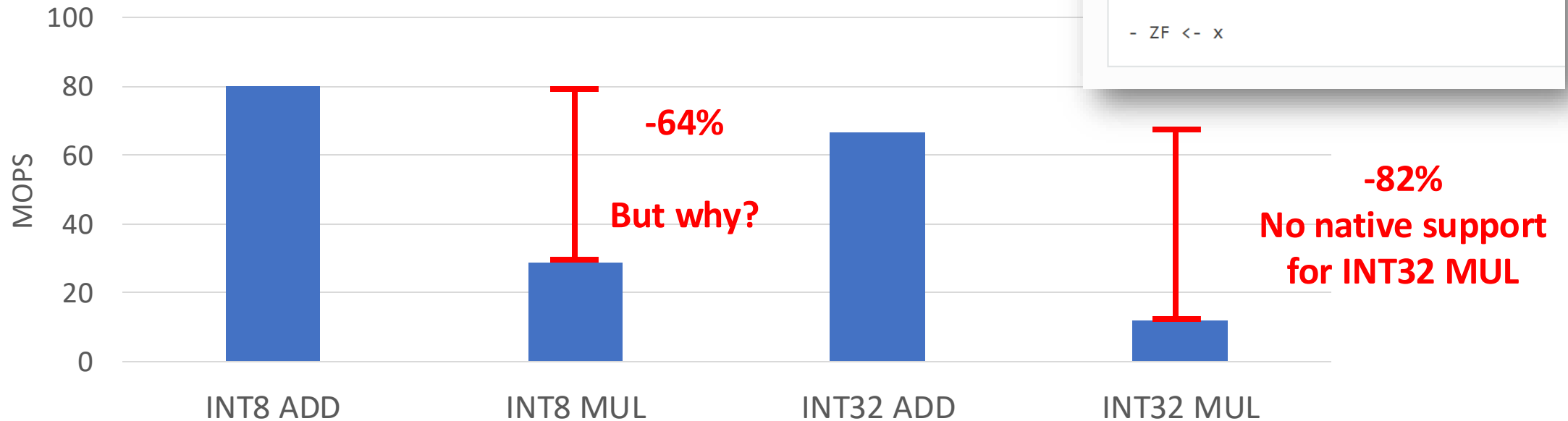
# ADD and MUL performance analysis

# Basic operations on UPMEM

```
void update(T *buffer, T scalar) {
    for (unsigned int i = 0; i < BLOCK_SIZE / sizeof(T); i += 1) {
#ifdef ADD
        bufferA[i] += scalar;
#elif MUL
        bufferA[i] *= scalar;
#endif
    }
}
```



-64%

But why?

-82%
No native support
for INT32 MUL

# Basic operations on UPMEM

```
void update(T *buffer, T scalar) {
    for (unsigned int i = 0; i < BLOCK_SIZE / sizeof(T); i +=
#ifdef ADD
        bufferA[i] += scalar;
#elif MUL
        bufferA[i] *= scalar;
#endif
    }
}
```

**mul_sl_sl**

**mul_sl_sl:rrr**

*mnemonic*

```
mul_sl_sl rc:wr32 ra:r32 rb:wr32
```

*behavior*

```
- x = (ra[0:7] * rb[0:7]):S32

- rc = x

- ZF <- x
```



**-64%**

**But why?**

**-82%**
**No native support**
**for INT32 MUL**

MOPS — chart: INT8 ADD ≈ 80, INT8 MUL ≈ 28, INT32 ADD ≈ 66, INT32 MUL ≈ 12

# Issues with INT8 MUL

- Compiler sometimes uses
  **INT32 MUL** *__mulsi3* function
    (SHIFT&ADD, ~9 cycles)
  instead of
  **INT8 MUL** *mul_sl_sl* (1 cycle) instruction
- Often loops are not being unrolled at all
- Unrolling through *#pragma unroll* is often not possible – not enough IRAM

```
void update(char* x, char a) {
    for (int i = 0; i < N; ++i){
        char tmp = x[i];
        tmp *= a;
        x[i] = tmp;
    }
}

int main() {
    update(X, A);
    return 0;
}
```

```
    move r16, 0
.LBB0_1:
.Ltmp1:
    sub r17, r15, r16
    lbs r0, r17, 0
    move r1, r14
    call r23, __mulsi3
.Ltmp2:
    sb r17, 0, r0
.Ltmp3:
    add r16, r16, -1
.Ltmp4:
    jneq r16, -120, .LBB0_1
```

# __mulsi3 – SHIFT & ADD

## mul_step

## mul_step:rrrici

*mnemonic*

```
mul_step dc:wr64 ra:r32 db:wr64 shift:u5 boot_cc:cc pc:pc16
```
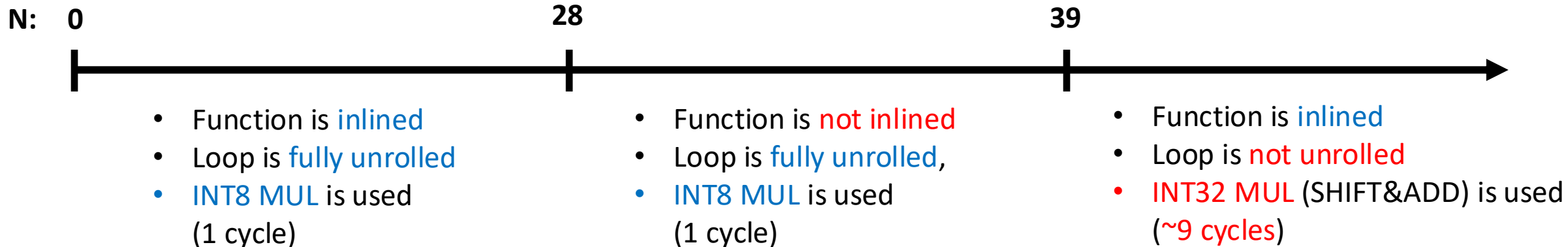
*behavior*

```
- cc = x = dbe >> 1

- cc = (dbe & 0x1) - 1

- if (const_cc_zero cc) then

-      dco = (dbo + (ra << shift))

- dce = x

- if (boot_cc cc) then

-      jump @[pc]

- ZF <- x
```

```
 1       jgtu r1, r0, __mulsi3_swap
 2       move r2, r0
 3       move r0, r1, true, __mulsi3_start
 4  ∨ __mulsi3_swap:
 5       move r2, r1
 6       move r0, r0
 7  ∨ __mulsi3_start:
 8       move r1, zero
 9       mul_step d0, r2, d0, 0, z, __mulsi3_exit
10       mul_step d0, r2, d0, 1, z, __mulsi3_exit
11       ...
12       mul_step d0, r2, d0, 31, z, __mulsi3_exit
13  ∨ __mulsi3_exit:
14       move r0, r1
15       jump r23
```

# Changing instructions on function inline

- Compilation with –O2 or –O3 will result in inlining based on N (if it's compile time constant)

- Loop unrolling also depends on N

- Inside the main function, INT32 MUL (SHIFT&ADD) is used

```
void update(char* x, char a) {
    for (int i = 0; i < N; ++i){
        char tmp = x[i];
        tmp *= a;
        x[i] = tmp;
    }
}

int main() {
    update(X, A);
    return 0;
}
```
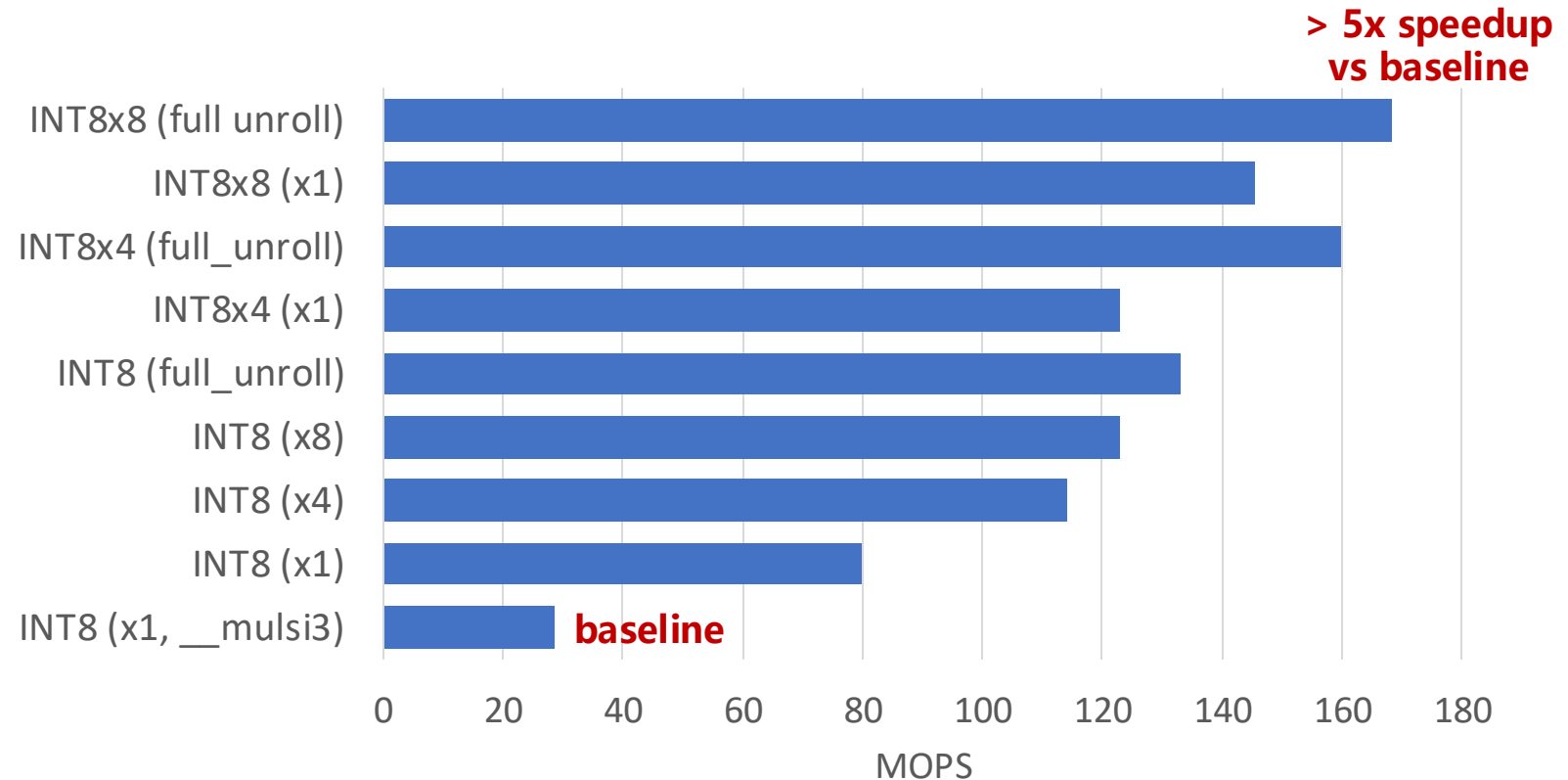
**N:**  **0**                              **28**                              **39**

- Function is inlined
- Loop is fully unrolled
- INT8 MUL is used (1 cycle)

- Function is not inlined
- Loop is fully unrolled,
- INT8 MUL is used (1 cycle)

- Function is inlined
- Loop is not unrolled
- INT32 MUL (SHIFT&ADD) is used (~9 cycles)

8

# INT8/INT32 MUL Optimization

# Optimized INT8 MUL

- Ensuring the use of the correct instruction for INT8
- Unrolling:
  - Full (compiler decides)
    *#pragma unroll*
  - Manual, e.g., x4
    *#pragma unroll 4*
- Loading data by 4 or 8 bytes (INT8x4, INT8x8)

**> 5x speedup vs baseline**

Chart — MOPS (x-axis: 0, 20, 40, 60, 80, 100, 120, 140, 160, 180):

| Category | Approx MOPS |
| --- | --- |
| INT8x8 (full unroll) | ~168 |
| INT8x8 (x1) | ~145 |
| INT8x4 (full_unroll) | ~160 |
| INT8x4 (x1) | ~122 |
| INT8 (full_unroll) | ~132 |
| INT8 (x8) | ~122 |
| INT8 (x4) | ~114 |
| INT8 (x1) | ~80 |
| INT8 (x1, __mulsi3) | ~28 **baseline** |

```
uint32_t w = *((uint32_t*) &bufferA[i]);
int8_t temp;
__builtin_mul_sl_sl_rrr(temp, w, scalar);
bufferA[i]     = temp;
__builtin_mul_sh_sl_rrr(temp, w, scalar);
bufferA[i + 1] = temp;
w >>= 16;
__builtin_mul_sl_sl_rrr(temp, w, scalar);
bufferA[i + 2] = temp;
__builtin_mul_sh_sl_rrr(temp, w, scalar);
bufferA[i + 3] = temp;
```
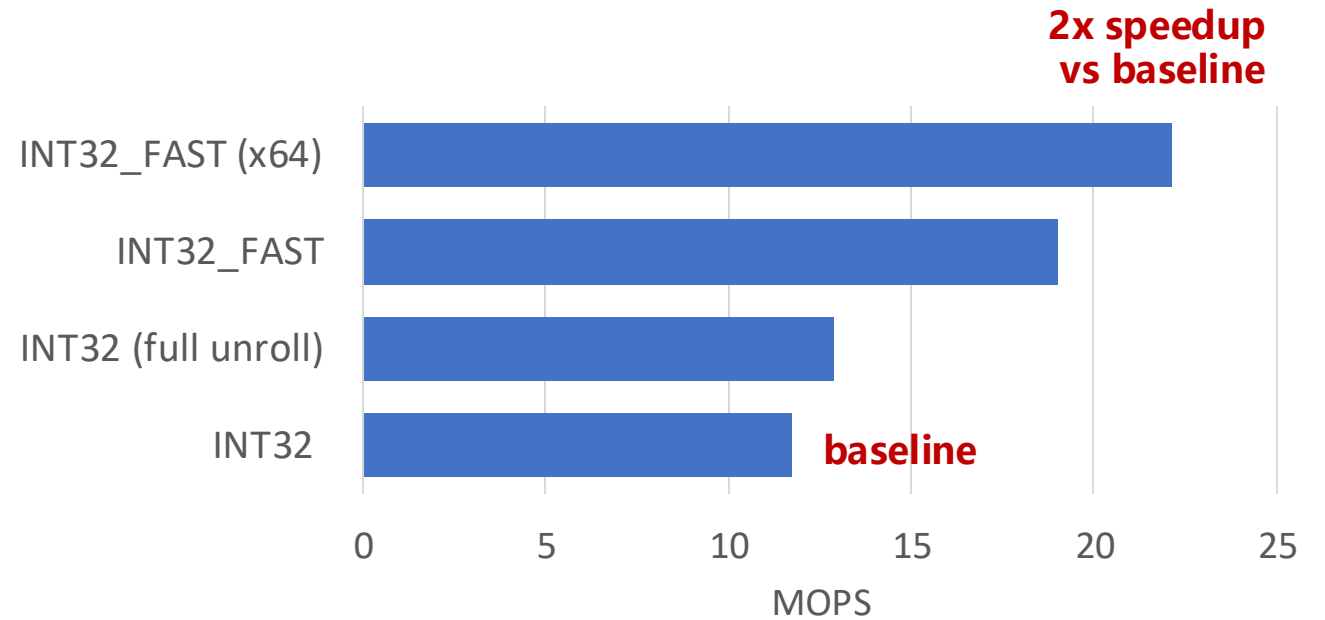
# Optimized INT32 MUL

- Replacing _*mulsi3*
  (SHIFT&ADD, ~30 cycles) with
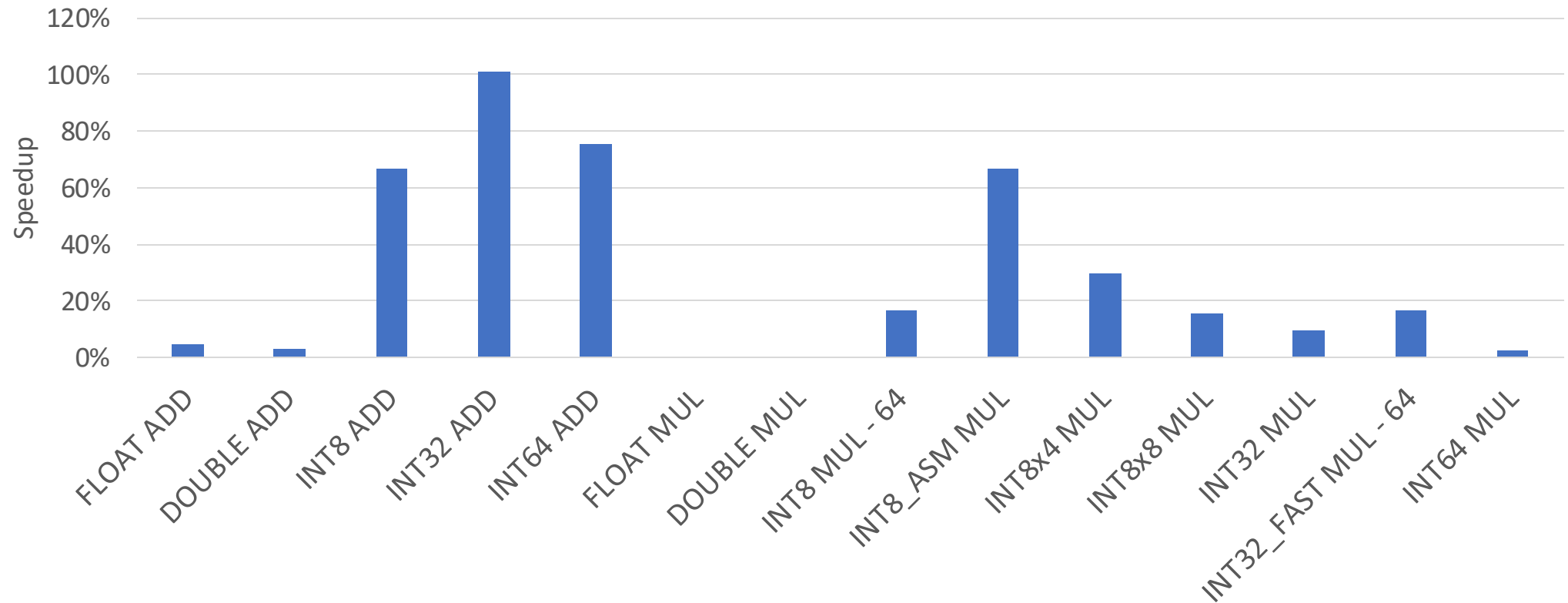  a custom implementation that
  uses INT8 MUL

- Unrolling also helps



**2x speedup vs baseline**

X = (x3, x2, x1, x0)
Y = (y3, y2, y1, y0)

$$X * Y = \quad 2^0 \quad (x0 * y0)$$
$$+ \ 2^8 \quad (x0 * y1 + x1 * y0)$$
$$+ \ 2^{16} \ (x0 * y2 + x1 * y1 + x2 * y0)$$
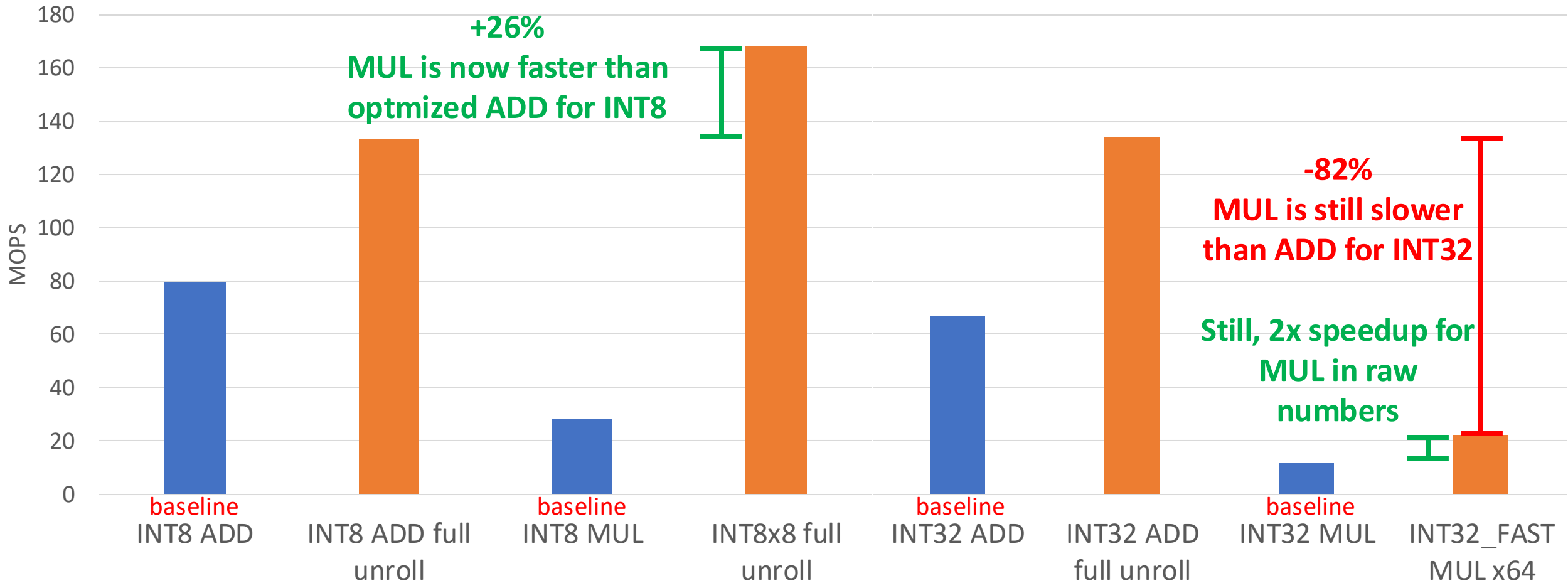$$+ \ 2^{24} \ (x0 * y3 + x1 * y2 + x2 * y1 + x3 * y0)$$
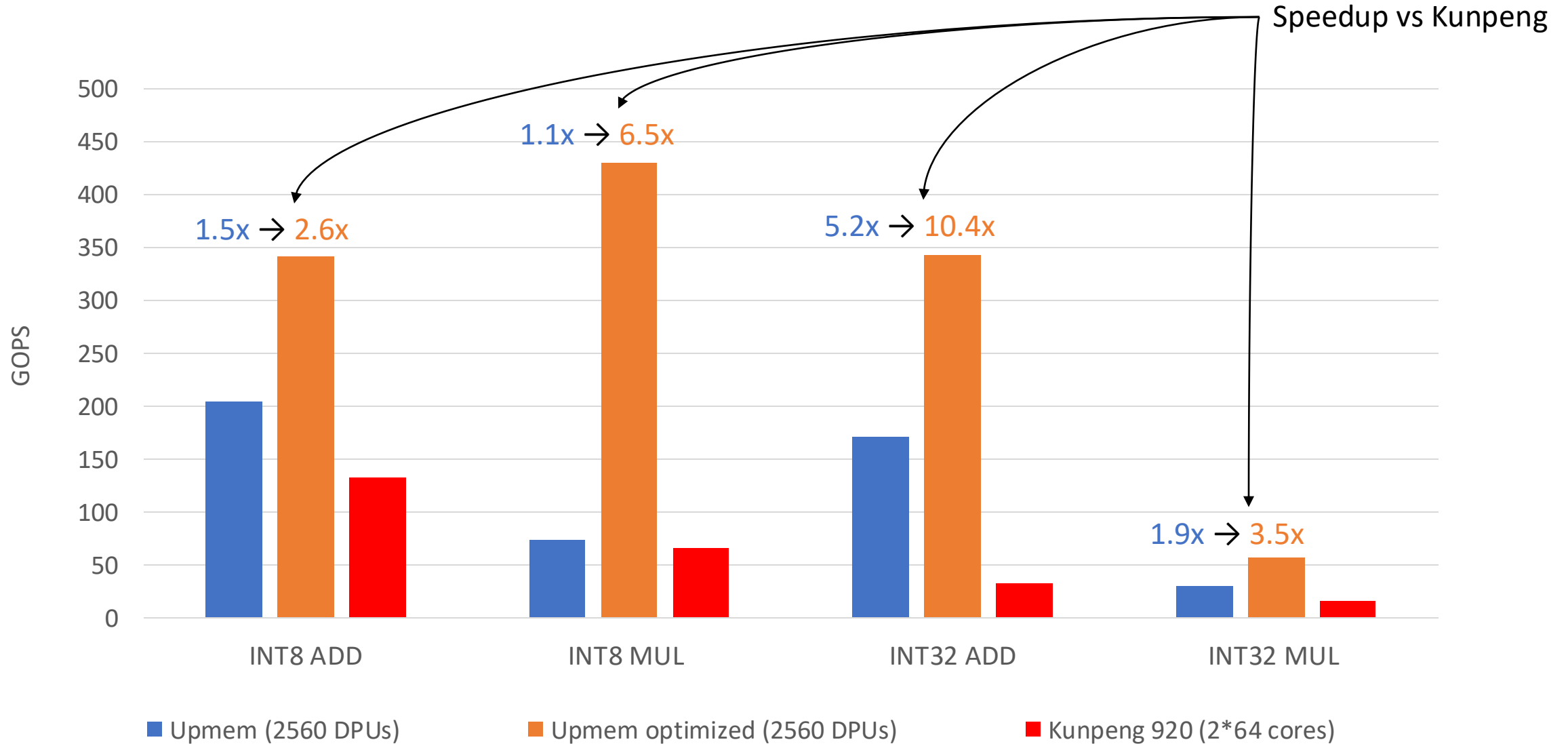
# Unrolling helps a lot



- Loops are pretty costly, when we are not spending a lot of cycles on single iteration
- Use *#pragma unroll* to force compiler to unroll loops

# Basic operations on UPMEM revisited



+26%
**MUL is now faster than optmized ADD for INT8**

-82%
**MUL is still slower than ADD for INT32**

**Still, 2x speedup for MUL in raw numbers**

180
160
140
120
100
80
60
40
20
0

MOPS

baseline INT8 ADD | INT8 ADD full unroll | baseline INT8 MUL | INT8x8 full unroll | baseline INT32 ADD | INT32 ADD full unroll | baseline INT32 MUL | INT32_FAST MUL x64

# Memory transfer optimizations

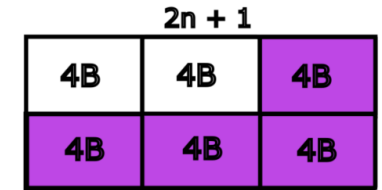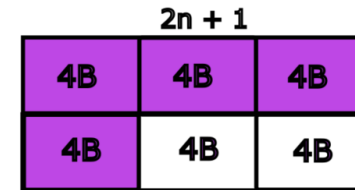# Memory transfers

**Host → PIM** (*dpu_prepare_xfer*, *dpu_push_xfer*):

- Transfer size and symbol offset need to be a multiple of 8B → error otherwise

- Extra care is needed to handle the remainder, if any

- Host buffers need to be aligned up to 8B

**MRAM → WRAM** (*mram_read, mram_write):*

- Transfer size needs to be a multiple of 8B (max 2048B)

- Read/write operations silently align down the address to 8B → no error is raised

- *mram_read_unaligned* and *mram_write_unaligned* functions require ~300 more instructions
  → GEMV where the matrix is stored contiguously and row size is odd

# Reading unaligned memory MRAM

- Address we read from/write to must be aligned to 8B
- Number of bytes read/written must be a multiple of 8B



```
1    int *buffer_wram = (int*)mem_alloc((BLOCK_SIZE + 2) * sizeof(float));
2    uint32_t buffer_addr = (uint32_t)(buffer_mram + offset);
3    int *buffer_access_ptr = NULL;
4 ∨  if (buffer_addr & 7) { // NOT aligned to 8B
5        mram_read(alignDownTo8(buffer_addr), buffer_wram, (BLOCK_SIZE + 2) * sizeof(int));
6        buffer_access_ptr = (buffer_wram + 1);
7 ∨  } else {
8        mram_read(buffer_addr, buffer_wram, BLOCK_SIZE * sizeof(int));
9        buffer_access_ptr = buffer_wram;
10   }
```
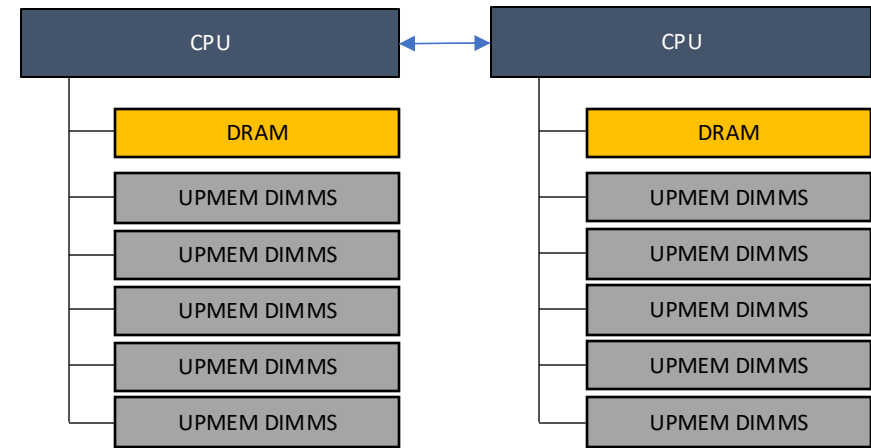
# Optimized data transfers

UPMEM server:

- DRAM can be the bottleneck (only 2 DDR4-3200 channels)
- Transposing data heavily engages the CPU

DPU allocation is not NUMA-aware:

- No way to allocate DPUs assigned to specific CPU
- A significant variance (~12%) in memory transfer speeds with default DPU allocation

Simple API extension to allow DPU allocation on specific NUMA node:

- Changes only in the user-space library
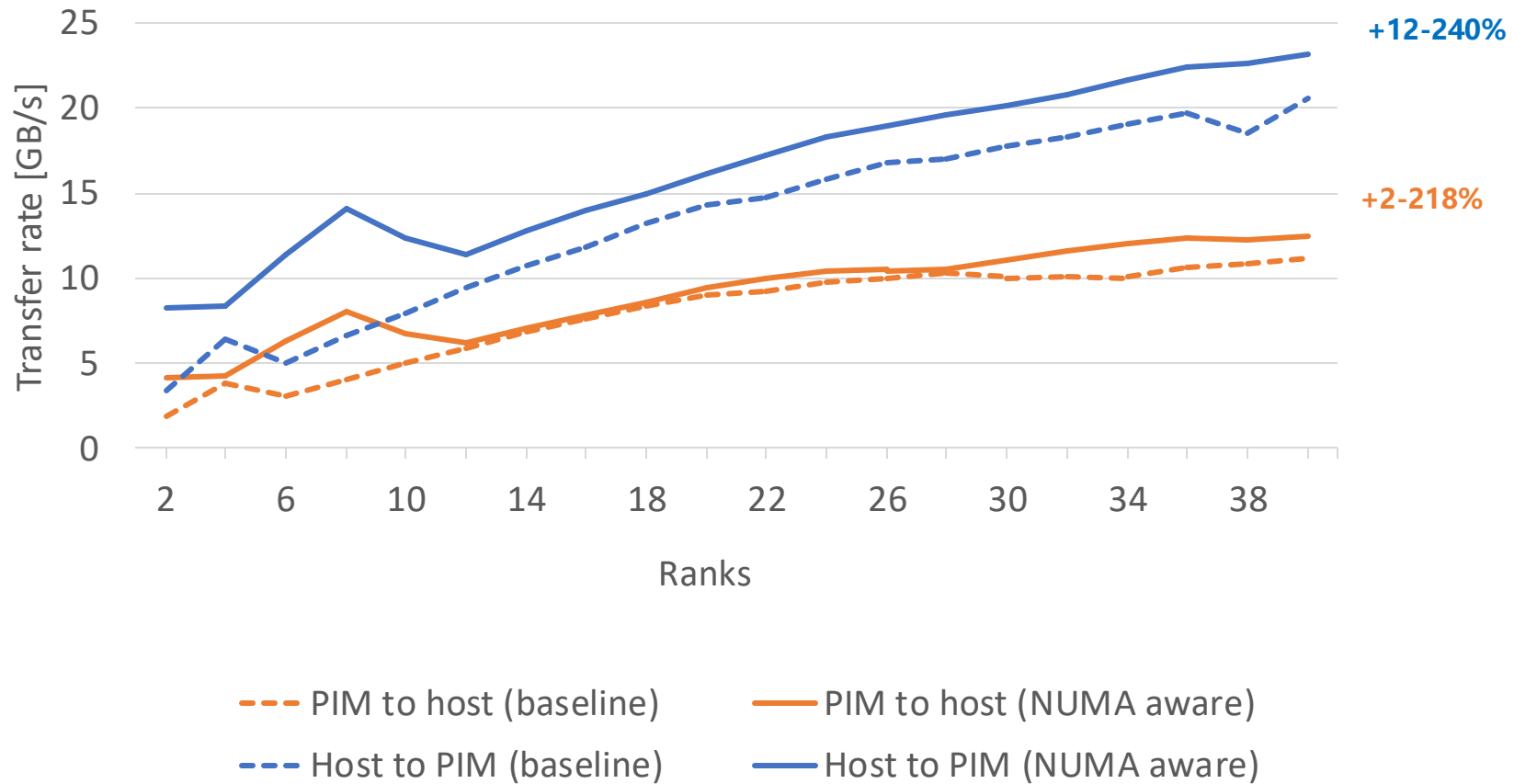- Simple filtering of available ranks during DPU allocations



```
uint64_t *buffer = alloc_buffer();
dpu_set_t *set;
DPU_ASSERT(dpu_alloc_ranks(ranks, NULL, set));
/* transfer data from buffer to set */
```

```
uint64_t *buffer[2];
buffer[0] = alloc_buffer_on_cpu(0);
buffer[1] = alloc_buffer_on_cpu(1);
dpu_set_t *set[2];
DPU_ASSERT(dpu_alloc_ranks(ranks/2, NULL, set[0], 0));
DPU_ASSERT(dpu_alloc_ranks(ranks/2, NULL, set[1], 1));
/* transfer data from buffer[0] to set[0] and from buffer[1] to set[1] */
```

# Optimized data transfers (parallel copy)



PIM to host (baseline) · · · · ·   PIM to host (NUMA aware) ———
Host to PIM (baseline) · · · · ·   Host to PIM (NUMA aware) ———

# Summary

UPMEM Compiler:

- Generates inefficient INT8 MUL (__*mulsi3* SHIFT&ADD)
- Rarely performs unrolling

DPU Optimizations:

- A little extra care with **INT8 MUL** provides huge gains
- **INT32 MUL** performance is also improved
- The achieved performance of **INT8/INT32 ADD/MUL** makes UPMEM **viable** for testing AI workloads (albeit with quantization, FP remains unpractical on UPMEM)

Optimized data transfers:

- **NUMA-awareness** can go a long way

Contact me at: krystian.chmielewski@huawei.com

# Thank you!