

Enabling Efficient Intra-Warp Communication for Fourier Transforms in a Many-Core Architecture

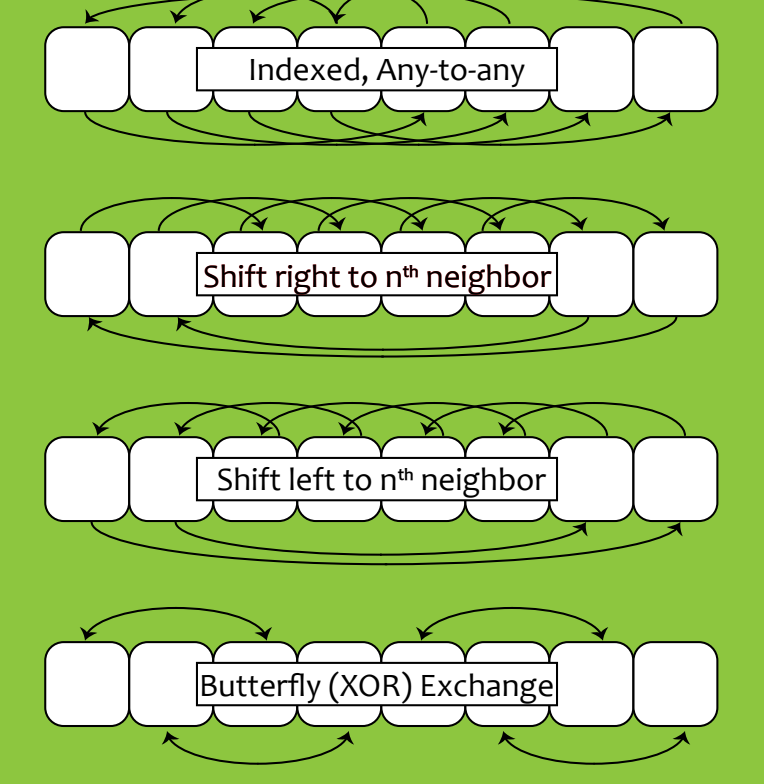
Student:

Carlo del Mundo (cdel@vt.edu)

Faculty:

Wu-chun Feng (feng@cs.vt.edu)

This work was supported in part by NSF I/UCRC IIP-0804155 via the NSF Center for High-Performance Reconfigurable Computing

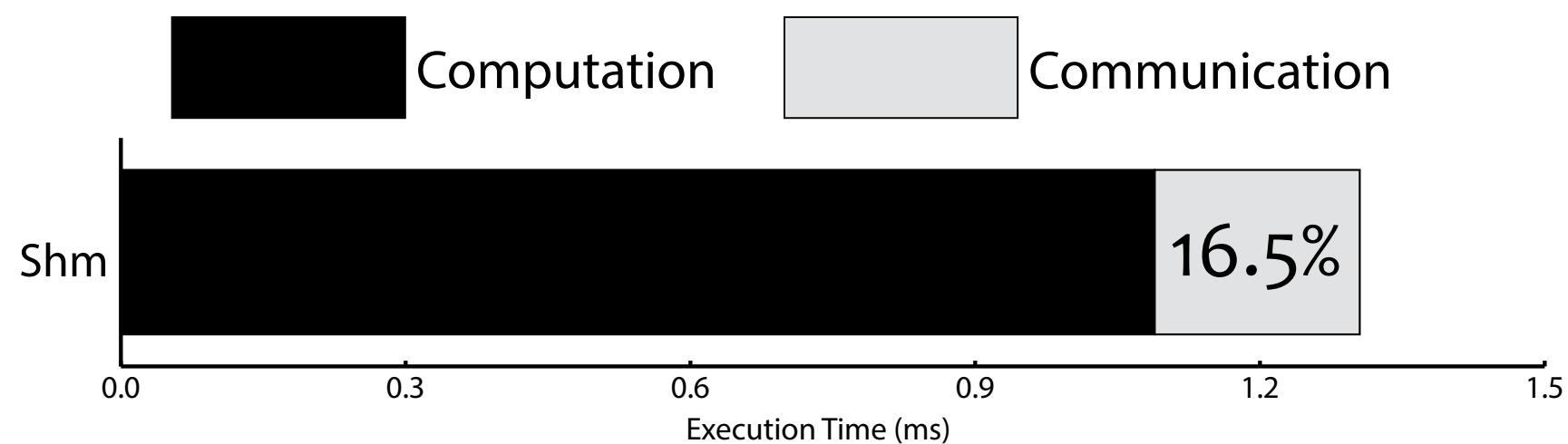


Shuffle Access Patterns

Goal >> Accelerate intra-warp communication for FFT using shuffle instructions <<

1. Motivation

- Communication oftentimes contributes to the cost of an algorithm.



Execution profile of a 256-pt FFT

- A new hardware mechanism (shuffle) allows for fast, intra-warp communication without shared memory. However, there is little success in accelerating applications using shuffle outside of vendor claims.

2. Background

- Shared memory:** traditional memory space for scratchpad communication for threads in a threadblock.
- Shuffle:** a new communication mechanism that allows for intra-warp register to register exchange. It was introduced in NVIDIA SM 3.0+. It supports arbitrary indexed references for data reads within a warp and has the advantage of carrying out a store-and-load operation in one step [1].
- CUDA local memory:** a per-thread spill-over memory resident in the L1 cache or global memory. Accessing CUDA local memory is slow.
- CUDA SELP instructions:** PTX instructions that allow for predicated stores. Generally, these instructions reduce (but do not eliminate) divergence. Instruction are of the form: $a = (\text{predicate}) ? b : c$;

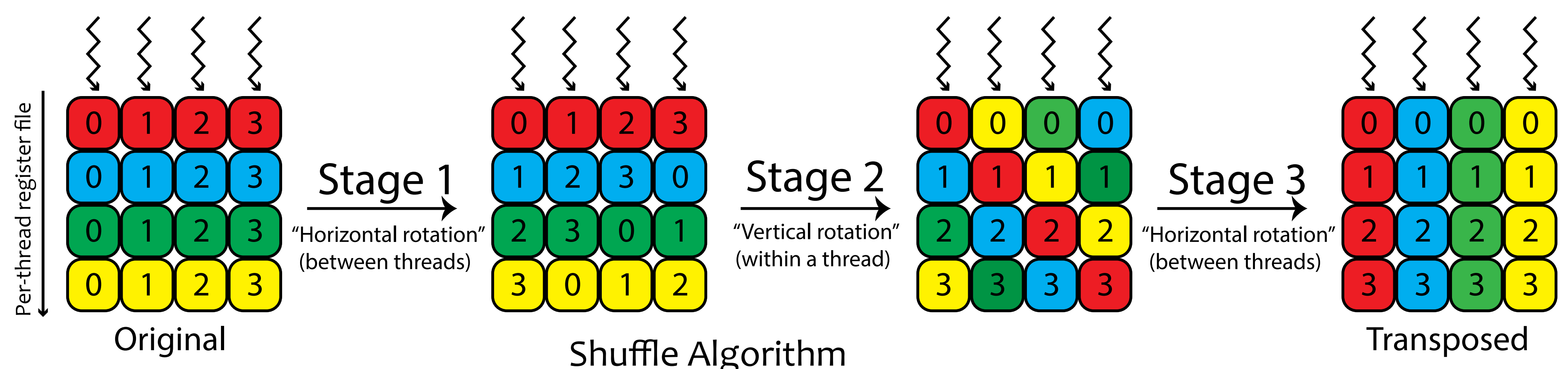
[1] NVIDIA's Next Generation CUDA™ Compute Architecture: Kepler™ Compute Architecture: Kepler™ v1.0, Available <http://www.nvidia.com/content/PDF/kepler/NVIDIA-Kepler-GK110-Architecture-Whitepaper.pdf>.

3. Approach

- Apply the shuffle mechanism to accelerate matrix transpose. Matrix transpose is the fundamental communication stage in the Fast Fourier Transform (FFT).

Application & Experimental Testbed

- Batched 256-pt FFT (128 MB of data). Each FFT contains a transpose 16x16 in size.
- NVIDIA Tesla K20c, 8192 thread blocks, 64 threads per block
- Our algorithm is shown for a 4x4 transpose. It is composed of horizontal and vertical rotations. Note: our algorithm is scalable up to 32x32 in size.



4. Characterization

- Kernels shown: 4x4 transpose. Actual results use the 16x16 transpose versions embedded inside a 256-pt FFT kernel
- We divide our implementations between the shared memory version ("shm") and several incremental shuffle versions ("naive", "div", "selp ip", "selp oop").

4.2 Shuffle

- Shuffle implementations bottlenecked by stage two (vertical rotations). This step is responsible for rearranging data within a thread's private register data.
- Our optimization process focuses on accelerating stage two.

4.1 Shared Memory

Shm

- Traditional mechanism to perform a transpose operation for threads in a threadblock.

```

#define WIDTH 4
void __global__ shm(int *input_ary, int *output_ary)
{
    int src_registers[4];
    // Initialization + Load from global memory to registers
    // ...
    // Store from registers to global memory
}
    
```

Naive

- CUDA local memory allocation and usage causes poor performance

```

#define WIDTH 4
void __global__ naive(int *input_ary, int *output_ary)
{
    int src_registers[4];
    // Initialization + Load from global memory to registers
    // ...
    // Store from registers to global memory
}
    
```

Div

- Refactored stage two by unrolling (but introduces divergence)
- No CUDA local memory

```

#define WIDTH 4
void __global__ div(int *input_ary, int *output_ary)
{
    int src_registers[4];
    // Initialization + Load from global memory to registers
    // ...
    // Store from registers to global memory
}
    
```

SELP (IP)

- Replaced divergence with SELP
- No CUDA local memory

```

#define WIDTH 4
void __global__ selp_ip(int *input_ary, int *output_ary)
{
    int src_registers[4];
    // Initialization + Load from global memory to registers
    // ...
    // Store from registers to global memory
}
    
```

SELP (OOP)

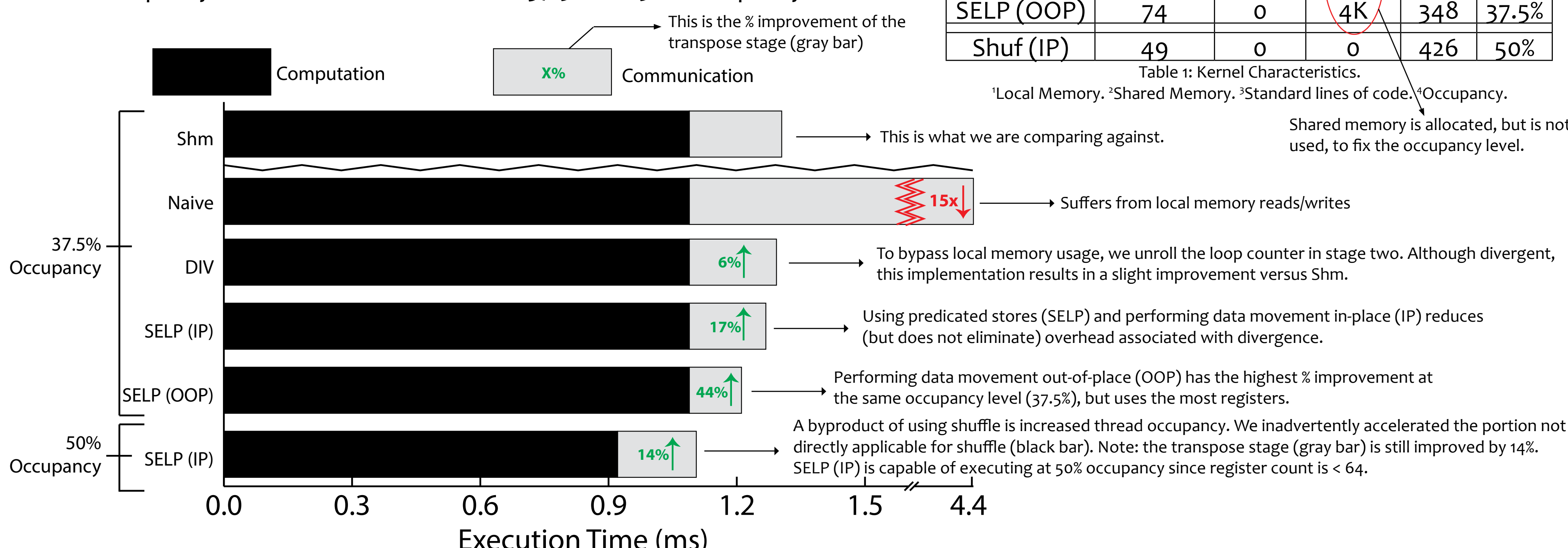
- Perform movement out of place
- No CUDA local memory

```

#define WIDTH 4
void __global__ selp_oop(int *input_ary, int *output_ary)
{
    int src_registers[4];
    int dst_registers[4];
    // Initialization + Load from global memory to registers
    // ...
    // Store from registers to global memory
}
    
```

5. Results

- Traditional shared memory FFTs are limited to 37.5% occupancy
- Shuffle implementations allocate, but do not use, shared memory to maintain a fixed occupancy. Results shown are for both 37.5% and 50% occupancy levels.



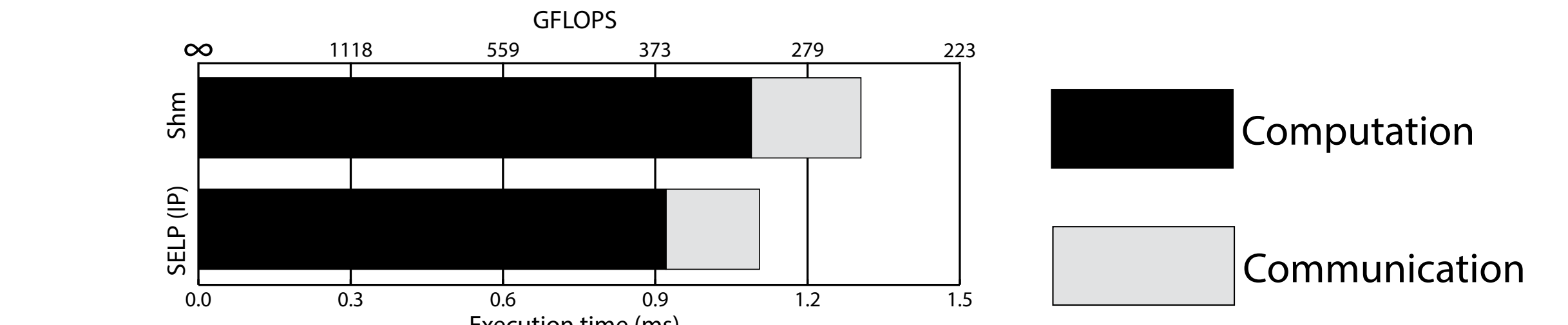
6. Insights & Conclusions

Insights

- Be wary of CUDA local memory allocation and usage. The CUDA compiler will allocate local memory if it cannot determine register indices at compile time. Accesses to local memory is slower than registers.

Conclusions

- SELP OOP improves the communication stage of FFT (gray bar) by 44%.
 - Overall speedup: 1.08x
 - Maximum speedup: 1.19x
- Since shared memory is eliminated, SELP IP is able to run at a higher occupancy level. We indirectly accelerated the computation portion of FFT in addition to communication.
 - Overall speedup: 1.17x
 - Maximum speedup: 1.19x



¹SELP OOP @ 37.5% occupancy. ²SELP IP @ 50% occupancy. ³Shm @ 37.5% occupancy.