





## A Proof-of-Concept Performance Portable SYCL-based Fast Fourier Transformation Library

10<sup>th</sup> International Workshop on OpenCL and SYCL

Vincent R. Pascuzzi Mehdi Goli (Codeplay)



May 11, 2022

## **DOE HPC Systems**

**2013** ~10 PFLOPS







2016

~10 PFLOPS

**2018** ~100 PFLOPS

**2021** ~100 PFLOPS **2022** ~1000 PFLOPS















### Homogeneous

### Heterogenous

# **DOE HPC Systems**

2013 2016 2018 2021 2022+ ~10 PFLOPS ~10 PFLOPS ~100 PFLOPS ~100 PFLOPS ~1000 PFLOPS ersc Perlmutter "EASY" SAK RIDG "HARD" EHEK (Ø 1NT **ENERG o** = AMD 3 **Heterogenous** 

Homogeneous

## HEP Center for Computational Excellence (CCE)

- A Department of Energy High-Energy Physics program investigating:
  - Performance portability
  - I/O
  - Complex workflows
  - Event generators\*

Portable Parallelization Strategies (PPS) effort focuses on performance and portability solutions for current and future HEP software

> Select among the participating experiments a number of x86-based 'testbeds' and rewrite the codes in various programming models

\* Event generator software is written and maintained by theorists.





Limited number of developers (we are physicists) and there are numerous platforms with various architectures

- Large codebases which need to stand the 'test of time" (~decade)
- We cannot afford to support and maintain multiple codebases
- We need to utilize leadership computing facilities.
- We need portability and attain a fair level of performance.

 LUMI (Finland)
 Leonardo (Italy)

 MD CPU, AMD GPU
 Ital (CPU, NVIDIA GPU)

Swiss National Supercomputing Center (CSCS), 2023 NVIDIA/ARM CPU, NVIDIA GPU

AL PS CSCS ETH Zürich Riken Center for Computational Science, 2021 ARM CPU







### Decouple

- from vendors providing single architecture/platform libraries
- boiler-plate code when using interoperability

### Readability, maintainability and sustainability

- backwards compatibility (*e.g.*, APIs constantly evolving ⇒ maintain interoperable implementations)
- source polluted with macros e.g., #ifdef to deal with many backends
- foster inclusivity: community-driven open-source projects

### And more...

- provide template metaprogramming for users to optimize for their target hardware (we don't care about what they want to run on, only that it does run)
- auto-tuning mechanisms: query available devices and configure/tune implementation for automatically optimizing parameters for a given device
- elevate reproducibility as first-class requirement

## **Performance portability**

"An application is performance portable if it achieves a consistent ratio of the actual time to solution to either the best-known or the theoretical best time to solution on each platform with minimal platform specific code required." <sup>1</sup>

#### Performance

- It runs: {Yes, No}
- It runs efficiently with respect to some baseline

### Portability

- Can execute on multiple systems
- Adaptable to varying architectures and platforms

$$\mathcal{P}(a, p; H) = \begin{cases} \frac{|H|}{\sum_{i \in H} \frac{1}{e_i(a, p)}} & \text{if } i \text{ is supported } \forall_i \in H \\ 0 & \text{otherwise} \end{cases}$$

### Useful metric should<sup>2</sup>:

- Be measured specific to a set of platforms of interest H
- Be independent of the absolute performance across H
- Be zero if a platform in H is unsupported, and approach zero as the performance of platforms in H approach zero
- Increase if performance increases on any platform in H
- Be directly proportional to the sum of scores across H

#### Productivity

- SLoC, maintainability, sustainability
- Port/migration/translation

### Reproducibility

- Crucial for most science
- Results (required precision) cannot depend on hardware

 <sup>&</sup>lt;sup>1</sup> (definition of) Performance Portability, <u>2016 Department of Energy Center of Excellence Meeting</u>.
 <sup>2</sup> Pennycook *et al.* (2019).

## SYCL-FFT

Based on Cooley-Tukey (see backup)

Implements radix-{2,4,8} algorithms

### Header-only

- Functor templated class, three template arguments
- WG\_FACTOR depends on input sequence length, determined *a priori*

### Proof-of-concept

- Limited to 1D
- C2C up to 2<sup>11</sup> length
- Computed out-of-place

```
template <typename T, size_t WG_SIZE,</pre>
      int WG_FACTOR, int SYCL_LANGUAGE_VERSION>
    class fft_1d {
     public:
      using float2 = sycl::float2;
      if constexpr (SYCL_LANGUAGE_VERSION < 202000) {</pre>
         using size_accessor =
           sycl::accessor<size_t, 1,</pre>
             sycl::access::mode::read,
             sycl::access::target::global_buffer>;
11
         using read_accessor =
           sycl::accessor<T, 1,</pre>
12
13
             sycl::access::mode::read.
14
             sycl::access::target::global_buffer>;
15
         using write_accessor =
           sycl::accessor<T, 1,</pre>
16
17
             sycl::access::mode::write,
18
             sycl::access::target::global_buffer>;
19
      } else {
20
         using stage_accessor =
21
           sycl::accessor<size_t, 1,</pre>
             sycl::access::mode::read,
22
23
             sycl::access::target::device>;
24
         using read_accessor =
25
           sycl::accessor<T, 1.</pre>
26
             sycl::access::mode::read.
27
             sycl::access::target::device>;
28
         using write_accessor =
29
           sycl::accessor<T, 1,</pre>
30
             sycl::access::mode::write,
31
             sycl::access::target::device>;
32
```

# $$\begin{split} X_k &= E_k + \omega_N^k O_k + \omega_N^{3k} O_k' \\ X_{k+N/2} &= E_k - (\omega_N^k O_k + \omega_N^{3k} O_k') \\ X_{k+N/4} &= E_{k+N/4} - i(\omega_N^k O_k - \omega_N^{3k} O_k') \\ X_{k+3N/4} &= E_{k+N/4} + i(\omega_N^k O_k - \omega_N^{3k} O_k') \end{split}$$

33	<pre>using local_rw_accessor =</pre>
34	<pre>sycl::accessor<t, 1,<="" pre=""></t,></pre>
35	<pre>sycl::access::mode::read_write,</pre>
36	<pre>sycl::access::target::local&gt;;</pre>
37	<pre>if constexpr (SYCL_LANGUAGE_VERSION &gt;= 202000) {</pre>
38	[[sycl::reqd_work_group_size(WG_SIZE)]]
39	}
40	<pre>fft1d(size_accessor stage_sizes,</pre>
41	read_accessor inputs,
42	write_accessor outputs,
43	<pre>local_rw_accessor local_shared,</pre>
44	size_t work_group_size,
45	<pre>int direction);</pre>
46	<pre>void operator()(sycl::nd_item&lt;1&gt; item) const {]</pre>
47	<pre>inline void radix_2(sycl::nd_item&lt;1&gt; item,</pre>
48	<pre>size_t stage_mod,</pre>
49	<pre>float2* temp) const {}</pre>
50	<pre>inline void radix_4(sycl::nd_item&lt;1&gt; item,</pre>
51	<pre>size_t stage_mod,</pre>
52	<pre>float2* temp) const {}</pre>
53	<pre>inline void radix_8(sycl::nd_item&lt;1&gt; item,</pre>
54	<pre>size_t stage_mod,</pre>
55	<pre>float2* temp) const {}</pre>
56	};

## **Experimental setup**

Device (Architecture)	Maximum Work-Group Size	Backend	Compiler(s)	Native Library
ARM Neoverse-N1	4096	POCL 1.9	ComputeCpp 2.8.0	ARM PL 21.1
(ARMv8-A)		(pre-gde9b966b)	armclang 21.1	
Intel Xeon E3-1585 v5		OpenCL 3.0	ComputeCpp 2.8.0	oneMKI 2022 0.2
(x86_64)	0192	oneTBB 2022.0.2	sycl-nightly/20220223	
AMD MI-100	256	HIP 4.2.0	sycl-nightly/20220223	rocfft 4.2.0
(CDNA)	230		hipcc 4.2.21155	
NVIDIA A100	1024	PTX64	sycl-nightly/20220223	cufft 11.5.0
(Ampere)	1024		nvcc 11.5.0	
Intel Iris P580		$O_{\text{man}} \subset I_{2} O$	Compute Cpp 2.8.0	
(Gen9)*	230	Opener 5.0		

Table 1: Device hardware and software versions for each platform considered in these studies. The sycl-nightly/x compilers refer to the specific branch of the Intel LLVM compiler project. All systems run openSUSE 15.3, kernel version 5.3.18. No native libraries exist for Intel Iris P580 iGPU (marked by asterisk); performance plots from this device are presented in App. A.

## **Computational Performance Ampere and CDNA**



(a) Smallest total (kernel dispatch + execution) execution times.

(b) Smallest kernel execution times.

# Computational Performance RISC, x86 and Gen9



(a) Smallest total (kernel dispatch + execution) execution times.





(b) Smallest kernel execution times.

 $\chi^2_{\text{reduced}} = \sum_{i=1}^{N} \frac{(s_i - n_i)^2}{n_i} \frac{1}{\text{ndf}}$ 

# Reproducibility





(b) AMD MI-100



(a) ARM Neoverse

(b) Intel x86\_64

# Conclusions

## Help grow SYCL ecosystem with purely SYCL-based libraries

- Decouple from vendors; improve readability, maintainability and focus on reproducibility
- Foster inclusive community-driven open-source software projects

## PoC SYCL-FFT

- Limited to 1D, C2C, lengths up to 2<sup>11</sup>
- To our knowledge, first demonstration of intel/llvm HIP backend
- Kernel runtime competitive with vendor-optimized libraries
- Small kernels suffer from launch latencies, sporadic and highly fluctuate on some hardware

### Ongoing and future work

- Continued improvements to SYCL backend implementations and offloading mechanisms potential close gaps
- SYCL-FFT support for {2,3}D FFT, accommodate arbitrary input sizes
- Further analyses including hipSYCL (highly optimized HIP backend)

# Conclusions

## support and access to test machines! Help grow SYCL ecosystem with purely SYCL-based libraries

- Decouple from vendors; improve readability, maintainability and focus on reproducibility
- Foster inclusive community-driven open-source software projects ٠

## PoC SYCL-FFT

- Limited to 1D, C2C, lengths up to 2<sup>11</sup>
- To our knowledge, first demonstration of intel/llvm HIP backend
- Kernel runtime competitive with vendor-optimized libraries ٠
- Small kernels suffer from launch latencies, sporadic and highly fluctuate on some hardware

## Ongoing and future work

- Continued improvements to SYCL backend implementations and offloading mechanisms potential close gaps
- SYCL-FFT support for {2,3}D FFT, accommodate arbitrary input sizes
- Further analyses including hipSYCL (highly optimized HIP backend) ٠

Special 'thanks' to our ANL friends---

. Benjamin Allen, Kalyan Kumaran

and Kevin Harms---for all their

## Backup

## FFT 101

 Discretize input function (time/space), map to frequency domain

$$X_{k} = \sum_{n=0}^{N-1} x_{n} e^{-i2\pi kn/N} = \sum_{n=0}^{N-1} x_{n} \omega_{N}^{kn} \qquad x_{k} = \frac{1}{N} \sum_{n=0}^{N-1} x_{n} \omega_{N}^{-kn}$$

Cooley-Tukey: exploit periodicity, chunk length-*N* FFT into smaller ones

$$X_{k} = \sum_{n=0}^{N/2-1} x_{2n} \omega_{N/2}^{(2n)k/N} + \sum_{n=0}^{N/2-1} x_{2n+1} \omega_{N/2}^{(2n+1)k/N}$$
$$X_{k} = \sum_{n_{2}=0}^{N/2-1} x_{2n_{2}} \omega_{N/2}^{n_{2}k} + \sum_{n_{4}=0}^{N/4-1} \left( \omega_{N}^{k} x_{4n_{4}+1} \omega_{N/4}^{n_{4}k} + \omega_{N}^{3k} x_{4n_{4}+3} \omega_{N/4}^{n_{4}k} \right)$$

 $T(N) = N \log n$ with log *n* splits



Figure 1: Illustration of a radix-2 DIT on an DFT with input size N = 8. Intersecting vertical lines—displaying butterflylike patterns—among the input,  $\{x_i\}$ , correspond to combinations of additions and subtractions as per the twiddle factors,  $\omega_N^k$ ; see main text.

## **Auxiliary Figures**



Device	Compiler + Backend	Launch Latency $[\mu s]$	
ARM Neoverse-N1	ComputeCpp 2.8.0 + POCL 1.9	200-250	
Intel Xeon E3-1585 v5	ComputeCpp 2.8.0 + OpenCL 3.0	~ 50	
Intel Iris P580	ComputeCpp + OpenCL 3.0	650-800	
AMD MI-100	Intel LLVM + HIP 4.2.0	~ 80	
NVIDIA A100	Intel LLVM + CUDA 11.5.0	~ 40 (13)	



Figure 7: Distributions of 1000 combined kernel launch and execution times of SYCL-FFT across all hardware for an input sequence length of 1024.

Figure 8: Distributions of 1000 combined kernel launch and execution times of SYCL-FFT across all hardware for an input sequence length of 2048.