## **Toward Making HiCMA Hardware-Agnostic with SYCL**

<u>H. Ltaief<sup>1</sup></u>, S. Abdulah<sup>1</sup>, R. Alomairy<sup>1</sup>, A. Dabah<sup>1</sup>, Y. Hong<sup>1</sup>, H. Ibeid<sup>2</sup>, D. Keyes<sup>1</sup>, M. Moawad<sup>3</sup>, A. Nasr<sup>3</sup>

<sup>1</sup>Extreme Computing Research Center, KAUST <sup>2</sup>Intel Corporation <sup>3</sup>Brightskies

SIAM Conference on Computational Science and Engineering (CSE23) February 26 - March 3, 2023

RAI Congress Centre | Amsterdam, The Netherlands

### Acknowledgments

- ECRC: S. Abdulah, R. Alomairy, A. Dabah, M. Genton, Y. Hong, D. Keyes, L. Qu, M. Ravasi, Y. Sun (Linear Algebra, Statistics, and Seismic)
- Mines Paristech: W. Bader, Y. Mesri (Mesh Deformations)
- Paris Observatory France, Australian National University Research School of Astronomy & Astrophysic (ANU), National Astronomical Observatory of Japan:
   D. Gratadour, Jesse Cranney, O. Guyon (Computational Astronomy)
- o ICL@UTK: G. Bosilca, Q. Cao, J. Dongarra, Y. Pei (Runtime Systems)
- INRIA Bordeaux / LABRI: M. Faverge, S. Thibault (Runtime Systems)
- HPC Resources: KAUST, HLRS, ORNL, RIKEN (Machines)
- Vendors: A. Esposito (HPE/Cray), L. Gatineau (NEC), I. Said / S. Jones (NVIDIA),
   M. Sabony / A. Lashab (AMD), A. Al-Jeshi / A. Alabduljabbar (Intel), F. Dupros (ARM),
   T. Hoshiya (Fujitsu) (Chips)

### Hardware Landscape



### Intel #97



### IBM/NVIDIA CPU/GPU #4



AMD #27



Fujitsu ARM A64FX #2

### Hardware Landscape









AMD #27







GRAPHCORE



Fujitsu ARM A64FX #2

IBM/NVIDIA CPU/GPU #4

# **Shaheen-3 Announcement**

### • HPE to build 100+ Pflops/s: fastest Supercomputer in the Middle East

- 25 liquid-cooled Cray EX4000 cabinets
  - 18 CPU-only cabinets, 4608 nodes, each with two AMD Epyc Genoa CPUs
  - 7 CPU-GPU tightly coupled cabinets, 704 nodes, each with four NVIDIA Grace Hopper
- Slingshot 11 Interconnect
- 50PB of HPE's Cray ClusterStor E1000 storage
- ~20X faster than Shaheen-2
- Fully operational in 2023
- System heterogeneity mirrors our large user base
  - Material science, Earth science, Physical science, Biological and Environmental Science, Biological and Environmental Science, Marine Science, Plant Science, Computational Statistics, Chemistry, etc...
  - $\bigcirc$  Al for all!
- Half of our faculty members simulates: "to outcompute is to outcompete!"

# Hardware Landscape



#### AMD Epyc Milan-X

High cache capacity High memory bandwidth x86 programming env Memory-bound workloads



#### NVIDIA Grace Hopper

High speed CPU-GPU interconnect Memory coherency Support for mixed precisions Compute-bound workloads



6

Graphcore IPU

Al-focused chip Flat memory hierarchy High SRAM bandwidth Inference

# Hardware Landscape



### AMD Epyc Milan-X

#### High cache capacity

High memory bandwidth x86 programming env Memory-bound workloads



#### NVIDIA Grace Hopper

High speed CPU-GPU interconnect Memory coherency Support for mixed precisions Compute-bound workloads



7

Graphcore IPU

#### Al-focused chip

Flat memory hierarchy High SRAM bandwidth Inference

### How do we reconcile this hostile environment with HPC scientific applications?

### **HPC Scientific Applications**



Fig. 4: **Left:** Soil moisture residuals at the topsoil of the Mississippi River basin. **Right**: Wind speed (m/s) in the Arabian Sea.

3D Geospatial Statistics



3D Mesh Deformations



### Seismic Imaging



3D Computational Electromagnetics



Wireless Communications



Computational Astronomy 8

### Revisiting the Hourglass



### Reshaping Linear Algebra for Massively Parallel Architectures

- Expose parallelism using task-based programming models
- Enhance user-productivity using layers of abstraction
- Ensure scalability using asynchronous executions
- Mitigate synchronization overheads using fine-grained computations
- Reduce data motion using mixed precisions
- Exploit data sparsity using low-rank approximations
- Maintain code portability using standard basic blocks

Are you willing to redesign your algorithm? One possible productive solution: Matricization



11

Available at http://github.com/ecrc/hicma

### **HPC Scientific Applications**



Fig. 4: Left: Soil moisture residuals at the topsoil of the Mississippi River basin. Right: Wind speed (m/s) in the Arabian Sea.

3D Geospatial Statistics

South Barrison

3D Computational Electromagnetics



3D Mesh Deformations



Wireless Communications



### Seismic Imaging



Computational Astronomy Accelerating Multiple-Input Multiple-Output Wireless Communication Networks

- Implement the Sphere Decoder approach
- Matricize the tree traversal of the detection algorithm
- Design a multi-level approach to ensure high occupancy and increase accuracy
- Cast the computations in terms of batch GEMMs
- A. Dabah et al, accepted at ISC23

First multi-GPU version Multi-start





### Data Parallel C++ (DPC++)

### DPC++ = C++ and SYCL and Extensions.

Many of the DPC++ extensions (ex., Unified Shared Memory) are now merged in the new SYCL standard.

Why SYCL?

- Syntax for accelerators (device selection, offload, memory transfer).
- Single source (host and device code).
- Implicit or Explicit data-transfer.
- Stack based on standards and open specifications (CLANG, LLVM, Level Zero).

```
sycl::queue q(cpu_selector{});
auto A = sycl::malloc_shared<float>(n, q);
auto B = sycl::malloc_shared<float>(n, q);
but B = sycl::malloc_shared<float>(n, q);
auto B = sycl::malloc_shared<float>(n, q);
but B =
```

### oneMKL DPC++ APIs

sycl::queue Q{sycl::cpu\_selector{}}; sycl::queue Q{sycl::gpu\_selector{}}; sycl::queue Q{device};

Create queue attached to a given device or device type. All device execution goes through a queue object.

void \*mem = sycl::malloc\_shared(bytes, Q); void \*mem = sycl::malloc\_device(bytes, Q);

Allocate device-accessible memory. Malloc\_shared memory is also accessible from the host.

#### #include "oneapi/mkl.hpp"

#### int main() {

•••

## // Select device and create queue sycl::device dev = sycl::device(sycl::gpu\_selector()); sycl::queue queue(dev);

#### // Allocate arrays

a = sycl::malloc\_shared<double>(batch\_size \* sizeof(fp \*), queue);

•••

•••

•••

#### // Initialize array

sycl::event gemm batch done;

gemm\_batch\_done = oneapi::mkl::blas::gemm\_batch(main\_queue, ta, tb, m, n, k, alpha, (const fp \*\*) a, lda, (const fp \*\*) b, ldb, beta, c, ldc, group\_count, group\_size, gemm\_batch\_dependencies);

gemm\_batch\_done.wait();

### Intel Ponte Vecchio GPU Hardware Accelerator



### GEMM\_batch on PVC for MIMO Application



### GEMM\_batch on PVC for MIMO Application



### **HPC Scientific Applications**





3D Mesh Deformations



Wireless Communications



### Seismic Imaging



Computational Astronomy

### Tile Low-Rank as a Pragmatic Approach



*T. Mary*, PhD Dissertation, Block Low-Rank multifrontal solvers: complexity, performance, and scalability, 2017.

*C. Weisberger*, PhD Dissertation, Improving multifrontal solvers by means of algebraic Block Low-Rank representations, 2013.

### **Directed Acyclic Graph – Cholesky Factorization**



21

### Incremental performance improvements on Fugaku



### Most Time-Consuming Kernel: HCORE-TLR-GEMM

- Perform matrix-matrix multiplication on compressed tile low-rank data structures
- Call successive BLAS/LAPACK kernels, i.e., QR, RandSVD, GEMM
- Require recompression to restore the low rankness property of the tiles
- Encapsulate all these operations in a single HCORE-TLR-GEMM kernel call
- Limit opportunities for performance analysis and improvements

## Challenges with SYCL

- The new rank calculation is only known when the HCORE-TLR-GEMM kernel is executed, which blocks the asynchronous execution, until the rank is known.
- HCORE-TLR-GEMM needs to be free from intermediate memory allocations.
- A significant redesign of the main kernel is necessary to expose the internal function calls in favor of a batch mode of execution.
- This is required anyway for getting performance on hardware accelerators.
- For memory buffers, it not clear how this can be done optimally.

# Summary

Society for Industrial and Applied Mathematics

2022 SIAG/SC Initiatives

#### NEW CS PHDS (% of TOTAL) in the UNITED STATES by SPECIALITY, 2020

Source: CRA Taulbee Survey, 2021 | Chart: 2022 Al Index Report



## We are recruiting!

