# EXTENDING MAGMA PORTABILITY

## Final Presentation

Anna Fortenberry, UNT
UTK RECSEM REU 2022

Mentors: Dr. Stan Tomov, UTK and Dr. Kwai Wong, UTK

# CONTENTS

- Problem Overview
- Software and Hardware
- Methodology
- CUDA to DPC++ Translation
- Porting MAGMA SGEMM
- Hardware Usage
- Performance
- Conclusion

# 1 PROBLEM OVERVIEW

○ Supercomputers provide the computational power necessary to resolve problems in a vast number of important domains

**data science**

**quantum information science**

**applied mathematics**

**high performance computing**

**cybersecurity**

**artificial intelligence research**

[1], [2], [3]

4

## EVOLUTION OF SUPERCOMPUTER SYSTEM DESIGN

- ○ NVIDIA opened a new door for supercomputing (SC) capabilities with the invention of the GPU in 1999
- ○ NVIDIA Tesla K20X GPU powered the first successful hybrid SC system in 2012
- ○ SC Systems are continually increasing in diversity

[4], [5]

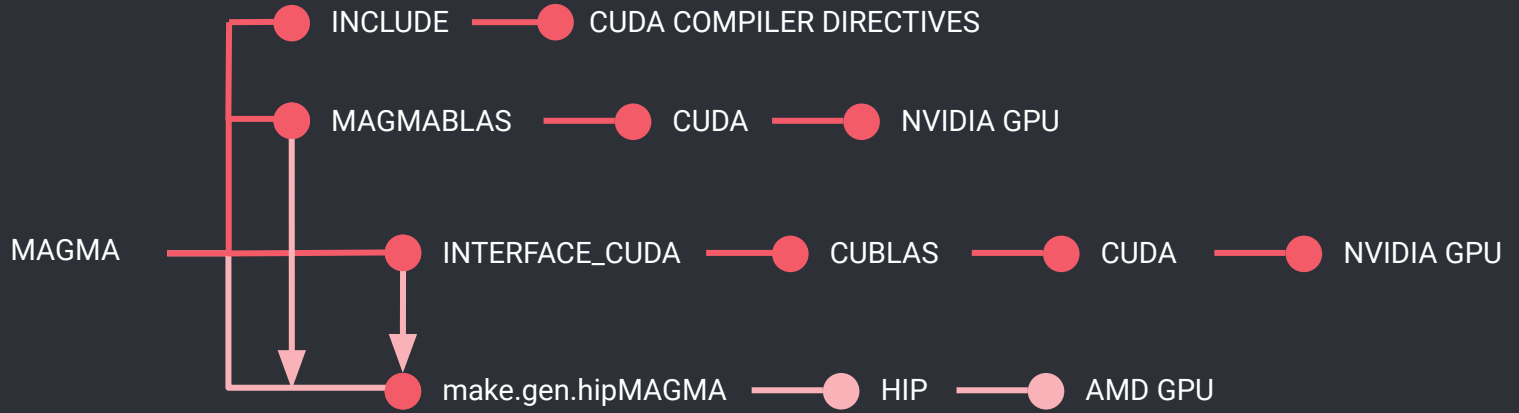| TOP500 The List | | | |
|---|---|---|---|
| JUNE 2022 | CPU/ Accelerator | JUNE 2019 | CPU/ Accelerator |
| Frontier | AMD, AMD | Summit | IBM, NVIDIA |
| S.C. Fugaku | Fugaku | Sierra | IBM, NVIDIA |
| LUMI | AMD, AMD | Sunway TaihuLight | Sunway |
| Summit | IBM, NVIDIA | Tianhe-2A | Intel |
| Sierra | IBM, NVIDIA | Frontera | Intel |
| Sunway TaihuLight | Sunway | Piz Daint | Intel, NVIDIA |
| Perlmutter | AMD, NVIDIA | Trinity | Intel |
| Selene | AMD, NVIDIA | ABCI | Intel, NVIDIA |
| Tianhe-2A | Intel, NUDT | SuperMUC-NG | Intel |
| Adastra | AMD, AMD | Lassen | IBM, NVIDIA |

[6], [7]

Anticipated for release in late 2022, Intel hopes to enter the supercomputer GPU vendor domain by powering the Aurora supercomputer at Argonne National Laboratory
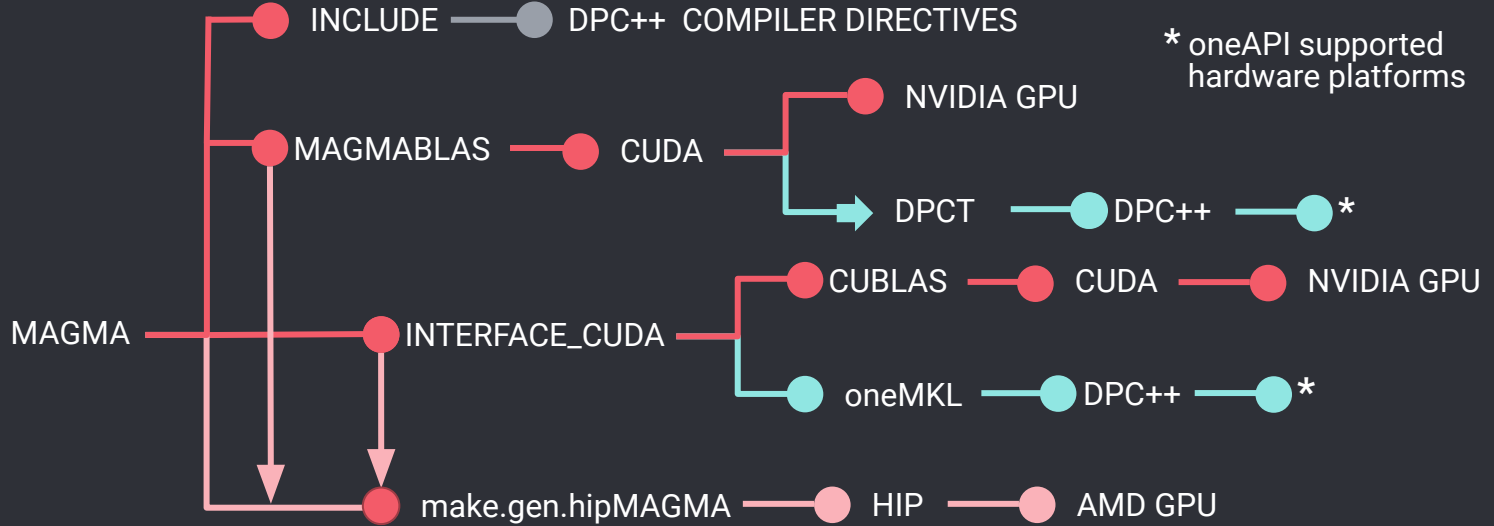
[8]

- Intel recently released a new programming model called **oneAPI**
- Applications that take advantage of oneAPI gain portability to all supported hardware platforms
  - CPUs (Scalar Architecture)
  - GPUs (Vector Architecture)
  - FPGAs (Spatial Architecture)
  - Other Accelerators (Matrix Architecture)

[9], [10]

INCLUDE — DPC++ COMPILER DIRECTIVES

* oneAPI supported hardware platforms

MAGMABLAS — CUDA

NVIDIA GPU

DPCT — DPC++ — *

CUBLAS — CUDA — NVIDIA GPU

MAGMA — INTERFACE_CUDA

oneMKL — DPC++ — *

make.gen.hipMAGMA — HIP — AMD GPU

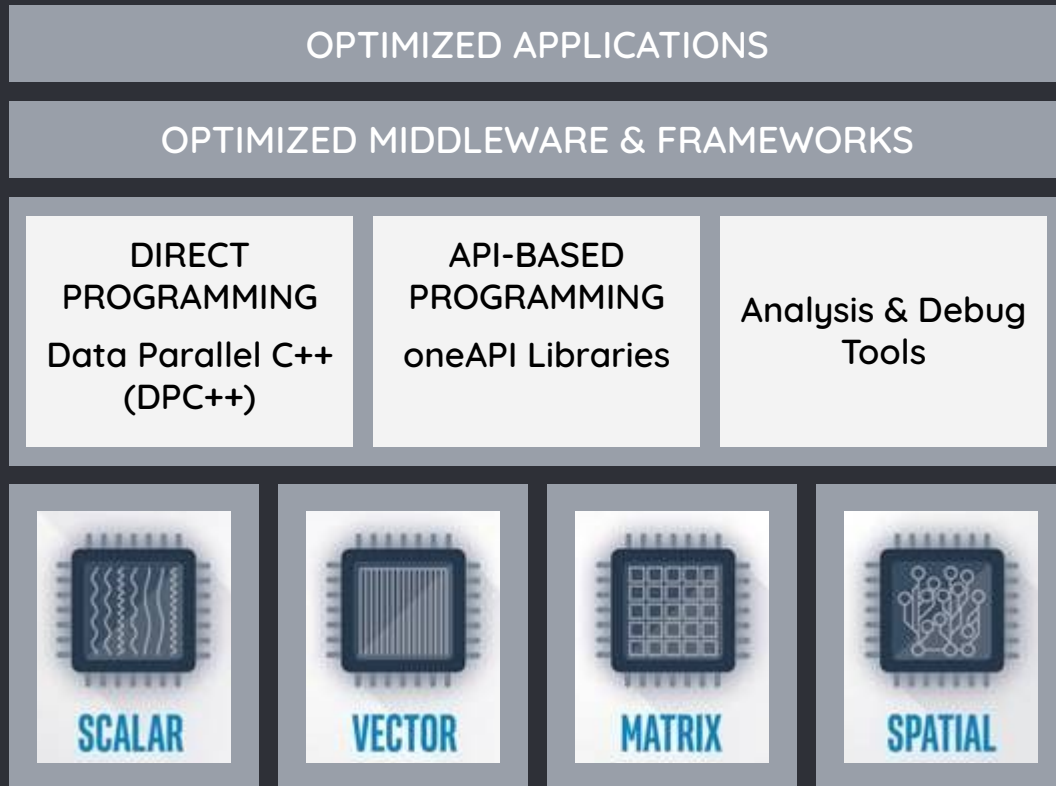○ oneAPI includes tools for adopting the model
  ○ Data Parallel C++ (DPC++) Translation Tool (DPCT)
  ○ oneAPI Math Kernel Library (oneMKL)

- How well does the DPCT tool translate CUDA code to DPC++ code?
- What are the common translation errors?
- Can this tool be used to translate MAGMA?
- Is DPC++ portable to Nvidia and AMD GPUs, and multicore CPUs?
- What is the performance of DPC++ on each of these accelerators comparative to CUDA?

# 2 SOFTWARE AND HARDWARE

| OPTIMIZED APPLICATIONS | | |
|---|---|---|
| **OPTIMIZED MIDDLEWARE & FRAMEWORKS** | | |
| **DIRECT PROGRAMMING** Data Parallel C++ (DPC++) | **API-BASED PROGRAMMING** oneAPI Libraries | Analysis & Debug Tools |



SCALAR  VECTOR  MATRIX  SPATIAL

- DPC++ is a oneAPI implementation of the Khronos standard **SYCL**
- SYCL is an accelerator language that allows code reuse across hardware targets
- SYCL adds data parallelism and heterogeneous programming to standard ISO C++

[10], [11]

# SOFTWARE OVERVIEW

## DPC++ Compatibility Tool (DPCT)

oneAPI tool to assist with migrating CUDA code to DPC++ code; translates with high accuracy

## oneAPI Math Kernel (oneMKL)

set of math routines for use in high performance computing on a variety of computational devices

## Compute Unified Device Architecture (CUDA)

NVIDIA parallel computing platform for harnessing power of GPUs

## DPC++-LLVM (CLang-LLVM)

LLVM-based compiler project that supports SYCL language

## DPC++ LLVM NVIDIA*

CLANG-LLVM build on Linux with CUDA NVIDIA support; allows DPC++ to port to NVIDIA GPUs

## Intel DevCloud

Remote development environments that grant access to Intel hardware for testing oneAPI projects*

[12], [13], [14], [15], [20]

# CENTRAL PROCESSING UNITS

## AMD EPYC 7742 PROCESSOR

| | |
|---|---|
| Cores: | 64 |
| Base Clock: | 2.25 Ghz |
| # of Threads: | 128 |
| Cache: | 256 MB |

## INTEL® XEON® PROCESSOR E5-2698 V4

| | |
|---|---|
| Cores: | 20 |
| Base Clock: | 2.20 Ghz |
| # of Threads: | 40 |
| Cache: | 50 MB |

[16], [17]

## GRAPHICS PROCESSING UNITS

**NVIDIA GEFORCE RTX 3060 (Discrete)**

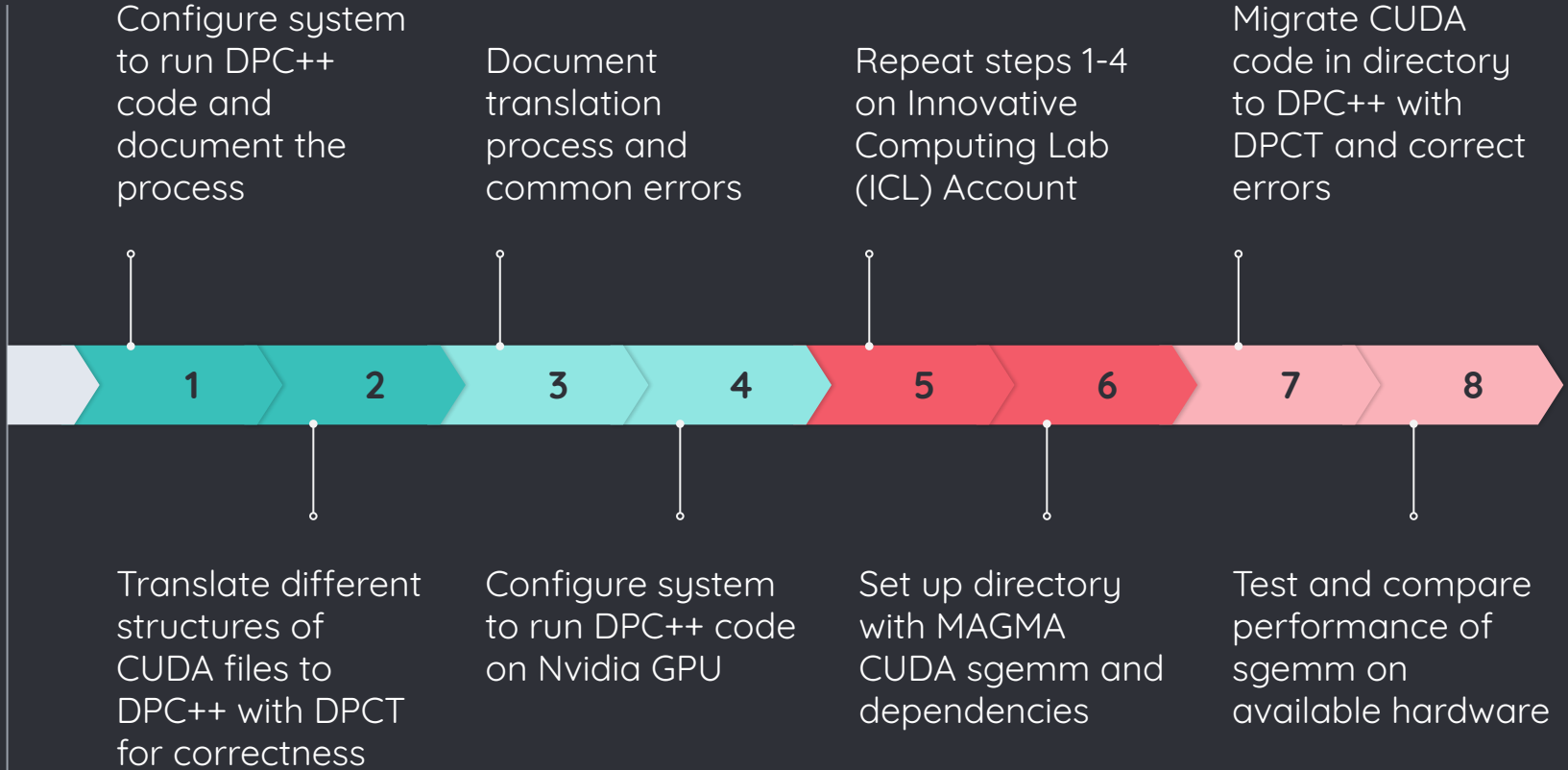GPU Cores:      3584
Base Clock:     1320 MHz
Memory Size:    12 GB

**INTEL UHD GRAPHICS P630 [0x3e96] (Integrated)**

GPU Cores:      192
Base Clock:     350 MHz
Memory Size:    Shared System

[18], [19]

# 3 METHODOLOGY

METHODOLOGY

Configure system to run DPC++ code and document the process

Document translation process and common errors

Repeat steps 1-4 on Innovative Computing Lab (ICL) Account

Migrate CUDA code in directory to DPC++ with DPCT and correct errors

**1** **2** **3** **4** **5** **6** **7** **8**

Translate different structures of CUDA files to DPC++ with DPCT for correctness

Configure system to run DPC++ code on Nvidia GPU

Set up directory with MAGMA CUDA sgemm and dependencies

Test and compare performance of sgemm on available hardware

18

# 4 CUDA TO DPC++ TRANSLATION

## SIMPLE KERNEL TRANSLATION

```
__global__ void VectorAddKernel(float* A, float* B, float* C)
{
    A[threadIdx.x] = threadIdx.x + 1.0f;
    B[threadIdx.x] = threadIdx.x + 1.0f;
    C[threadIdx.x] = A[threadIdx.x] + B[threadIdx.x];
}


void VectorAddKernel(float* A, float* B, float* C, sycl::nd_item<3> item_ct1)
{
    A[item_ct1.get_local_id(2)] = item_ct1.get_local_id(2) + 1.0f;
    B[item_ct1.get_local_id(2)] = item_ct1.get_local_id(2) + 1.0f;
    C[item_ct1.get_local_id(2)] =
    A[item_ct1.get_local_id(2)] + B[item_ct1.get_local_id(2)];
}
```

# TEST 1: ISOLATED FILE

- Translated files for CUDA vector addition and vector-matrix multiplication
- 100% compilation and execution accuracy
- CUDA error handling dead code clean up for file readability

```
/* DPCT1003:30: Migrated API does not return error code. (*, 0) is inserted.
You may need to rewrite this code. */

checkCudaErrors((h_C = (float *)sycl::malloc_host(mem_size_C,
    dpct::get_default_queue()),0));

-> h_C = (float *)sycl::malloc_host(mem_size_C,dpct::get_default_queue());
```

- Matrix-matrix multiplication file with six included headers
- 98.7% compilation accuracy and 98.0% execution accuracy in the main file
- 10% of the code needed dead code touchups
- Header files had 100% compilation accuracy and execution accuracies ranging from 75%-100%

```
cudaGetDeviceCount(&device_count);            device_count =
                                              dpct::dev_mgr::instance().device_count()


while (current_device < device_count)         while (current_device < device_count)
{                                             {
    cudaGetDeviceProperties                       dpct::dev_mgr::instance()
      (&deviceProp, current_device);                .get_device(current_device)
                                                    .get_device_info(deviceProp);


    if (deviceProp.computeMode !=                 if (true)
        cudaComputeModeProhibited)                { . . . }
    { . . . }
    else {                                        else {
        devices_prohibited++;                         devices_prohibited++;
    }                                             }


    current_device++;                             current_device++;
}                                             }
```

**5** PORTING MAGMA SGEMM

# C = A B
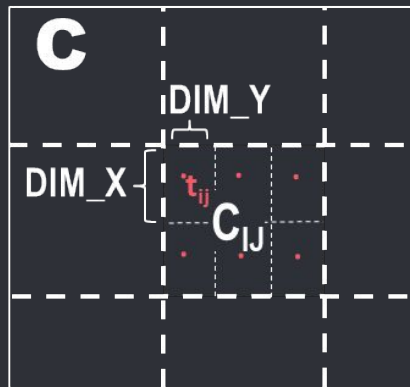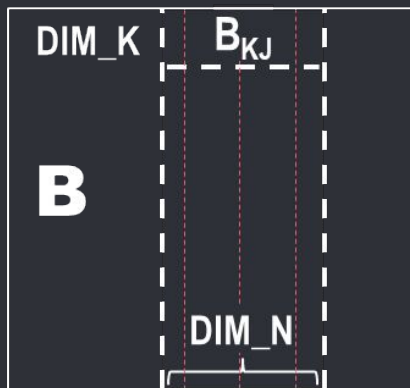
```
template < DIM_X, DIM_Y, DIM_M, DIM_N,
DIM_K, DIM_XA, DIM_YA, DIM_XB, DIM_YB>

For I = 1 .. M step DIM_M
    For J = 1 .. N step DIM_N
        For K = 1 .. K step DIM_K
            C_IJ += A_IK B_KJ
```



- Implementation is templated with **9 parameters**
- Computation is done with thread blocks of size **[ DIM_X , DIM_Y ]**
- Thread $t_{ij}$ computes [ DIM_M / DIM_X, DIM_N / DIM_Y ] elements of $C_{IJ}$
- $A_{IK}$ gets loaded in **shared memory** by [ DIM_XA , DIM_YA ] threads
- $B_{KJ}$ gets loaded in **shared memory** by [ DIM_XB , DIM_YB ] threads
- $C_{IJ}$ is held and computed in **registers**

25

## MAGMA SGEMM TRANSLATION PROCESS

- Collected MAGMA SGEMM CUDA code and dependencies in one directory
- Used DPCT to recursively migrate CUDA code to DPC++
- Translated header files that did not migrate independently in a separate directory and then copied them into the MAGMA SGEMM directory
- Implemented compiler directives as needed

# 6 HARDWARE USAGE

# MULTICORE CPUS



```
user1@REU1901-HP-Z800-Workstation: ~/anna/mtxMtxMulCnvt/one/dp...        user1@REU1901-HP-Z800-Workstation: ~/anna/mtxMtxMulCnvt/one/dp...

1  [|||||||||||||||100.0%]    4  [||||||||||||||||||||100.0%]    7  [|||||||||||||||||||||100.0%]   10 [||||||||||||||||||||||100.0%]
2  [|||||||||||||||100.0%]    5  [||||||||||||||||||||100.0%]    8  [|||||||||||||||||||||100.0%]   11 [||||||||||||||||||||||100.0%]
3  [|||||||||||||||100.0%]    6  [|||||||||||||||||||||7.24G/47.1G] 9 [|||||||||||||||||||||100.0%]   12 [||||||||||||||||||||||100.0%]
Mem[|||||||||||||||||||||||||||||||||||||||7.24G/47.1G]          Tasks: 207, 1038 thr; 12 running
Swp[                                        0K/2.00G]            Load average: 13.02 13.23 12.82
                                                                 Uptime: 23 days, 05:45:41


  PID USER      PRI  NI  VIRT   RES   SHR S CPU% MEM%   TIME+  Command
3602506 user1    20   0 13.1G 1628M  267M R 1178  3.4 15h41:40 ./intelCpuExec -wA=8192 -wB=8192 -hA=8192 -hB=8192
3602527 user1    20   0 13.1G 1628M  267M R 99.4  3.4 1h18:57 ./intelCpuExec -wA=8192 -wB=8192 -hA=8192 -hB=8192
3602528 user1    20   0 13.1G 1628M  267M R 99.4  3.4 1h18:57 ./intelCpuExec -wA=8192 -wB=8192 -hA=8192 -hB=8192
3602523 user1    20   0 13.1G 1628M  267M R 98.8  3.4 1h19:05 ./intelCpuExec -wA=8192 -wB=8192 -hA=8192 -hB=8192
3602526 user1    20   0 13.1G 1628M  267M R 98.8  3.4 1h18:44 ./intelCpuExec -wA=8192 -wB=8192 -hA=8192 -hB=8192
3602515 user1    20   0 13.1G 1628M  267M R 95.5  3.4 1h16:11 ./intelCpuExec -wA=8192 -wB=8192 -hA=8192 -hB=8192
3602521 user1    20   0 13.1G 1628M  267M R 99.4  3.4 1h18:58 ./intelCpuExec -wA=8192 -wB=8192 -hA=8192 -hB=8192
3602519 user1    20   0 13.1G 1628M  267M R 98.8  3.4 1h18:47 ./intelCpuExec -wA=8192 -wB=8192 -hA=8192 -hB=8192
3602534 user1    20   0 13.1G 1628M  267M R 98.8  3.4 1h19:03 ./intelCpuExec -wA=8192 -wB=8192 -hA=8192 -hB=8192
3602518 user1    20   0 13.1G 1628M  267M R 96.2  3.4 1h18:50 ./intelCpuExec -wA=8192 -wB=8192 -hA=8192 -hB=8192
3602516 user1    20   0 13.1G 1628M  267M R 97.5  3.4 1h18:09 ./intelCpuExec -wA=8192 -wB=8192 -hA=8192 -hB=8192
3602525 user1    20   0 13.1G 1628M  267M R 99.4  3.4 1h18:23 ./intelCpuExec -wA=8192 -wB=8192 -hA=8192 -hB=8192
3603925 user1    20   0 11708  5128  3220 R  2.6  0.0 0:16.55 htop
```

Intel(R) Xeon(R) CPU X5650 @ 2.67GHz

# MULTICORE CPUS



```
 27 [||||||||||||||||||||  100.0%]   91 [||||||||||||||||||||  100.0%]  155[||||||||||||||||||||  100.0%]  219[||||||||||||||||||||  100.0%]
 28 [||||||||||||||||||||  100.0%]   92 [||||||||||||||||||||  100.0%]  156[||||||||||||||||||||  100.0%]  220[||||||||||||||||||||  100.0%]
 29 [||||||||||||||||||||  100.0%]   93 [||||||||||||||||||||  100.0%]  157[||||||||||||||||||||  100.0%]  221[||||||||||||||||||||  100.0%]
 30 [||||||||||||||||||||  100.0%]   94 [||||||||||||||||||||  100.0%]  158[||||||||||||||||||||  100.0%]  222[||||||||||||||||||||  100.0%]
 31 [||||||||||||||||||||  100.0%]   95 [||||||||||||||||||||  100.0%]  159[||||||||||||||||||||  100.0%]  223[||||||||||||||||||||  100.0%]
 32 [||||||||||||||||||||  100.0%]   96 [||||||||||||||||||||  100.0%]  160[||||||||||||||||||||  100.0%]  224[||||||||||||||||||||  100.0%]
 33 [||||||||||||||||||||  100.0%]   97 [||||||||||||||||||||  100.0%]  161[||||||||||||||||||||  100.0%]  225[||||||||||||||||||||  100.0%]
 34 [||||||||||||||||||||  100.0%]   98 [||||||||||||||||||||  100.0%]  162[||||||||||||||||||||  100.0%]  226[||||||||||||||||||||  100.0%]
 35 [||||||||||||||||||||  100.0%]   99 [||||||||||||||||||||  100.0%]  163[||||||||||||||||||||  100.0%]  227[||||||||||||||||||||  100.0%]
 36 [||||||||||||||||||||  100.0%]  100[||||||||||||||||||||  100.0%]  164[||||||||||||||||||||  100.0%]  228[||||||||||||||||||||  100.0%]
 37 [||||||||||||||||||||  100.0%]  101[||||||||||||||||||||  100.0%]  165[||||||||||||||||||||  100.0%]  229[||||||||||||||||||||  100.0%]
 38 [||||||||||||||||||||  100.0%]  102[||||||||||||||||||||  100.0%]  166[||||||||||||||||||||  100.0%]  230[||||||||||||||||||||  100.0%]
 39 [||||||||||||||||||||  100.0%]  103[||||||||||||||||||||  100.0%]  167[||||||||||||||||||||  100.0%]  231[||||||||||||||||||||  100.0%]
 40 [||||||||||||||||||||  100.0%]  104[||||||||||||||||||||  100.0%]  168[||||||||||||||||||||  100.0%]  232[||||||||||||||||||||  100.0%]
 41 [||||||||||||||||||||  100.0%]  105[||||||||||||||||||||  100.0%]  169[||||||||||||||||||||  100.0%]  233[||||||||||||||||||||  100.0%]
 42 [||||||||||||||||||||  100.0%]  106[||||||||||||||||||||  100.0%]  170[||||||||||||||||||||  100.0%]  234[||||||||||||||||||||  100.0%]
 43 [||||||||||||||||||||  100.0%]  107[||||||||||||||||||||  100.0%]  171[||||||||||||||||||||  100.0%]  235[||||||||||||||||||||  100.0%]
 44 [||||||||||||||||||||  100.0%]  108[||||||||||||||||||||  100.0%]  172[||||||||||||||||||||  100.0%]  236[||||||||||||||||||||  100.0%]
 45 [||||||||||||||||||||  100.0%]  109[||||||||||||||||||||  100.0%]  173[||||||||||||||||||||  100.0%]  237[||||||||||||||||||||  100.0%]
 46 [||||||||||||||||||||  100.0%]  110[||||||||||||||||||||  100.0%]  174[||||||||||||||||||||  100.0%]  238[||||||||||||||||||||  100.0%]
 47 [||||||||||||||||||||  100.0%]  111[||||||||||||||||||||  100.0%]  175[||||||||||||||||||||  100.0%]  239[||||||||||||||||||||  100.0%]
 48 [||||||||||||||||||||  100.0%]  112[||||||||||||||||||||  100.0%]  176[||||||||||||||||||||  100.0%]  240[||||||||||||||||||||  100.0%]
 49 [||||||||||||||||||||  100.0%]  113[||||||||||||||||||||  100.0%]  177[||||||||||||||||||||  100.0%]  241[||||||||||||||||||||  100.0%]
```

AMD EPYC 7742 64-Core Processor

NVIDIA GPUS
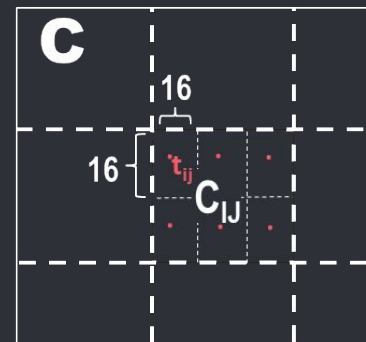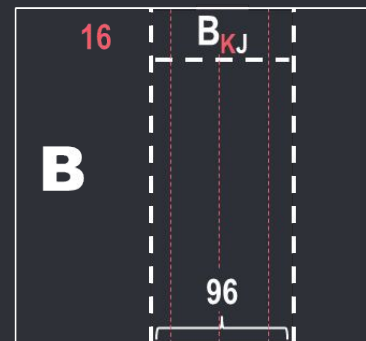
NVIDIA GeForce GTX 1650

**7** PERFORMANCE

## TEST PARAMETERS

**cuda** =
-DMAGMA_TUNING
-DDIM_X=16
-DDIM_Y=16
-DBLK_M_nn=96
-DBLK_N_nn=96
-DBLK_K_nn=16
-DDIM_XA=32
-DDIM_YA=8
-DDIM_XB=8
-DDIM_YB=32

## C = A B

template $< 16, 16, 96, 96, 16, 32, 8, 8, 32>$

For I = 1 .. M **step 16**
   For J = 1 .. N **step 16**
      For K = 1 .. K **step 16**
        $C_{IJ} += A_{IK} \ B_{KJ}$



- Thread $t_{ij}$ computes [ **96 / 16** , **96 / 16** ] elements of $C_{IJ}$
- $A_{IK}$ gets loaded in shared memory by [ **32**, **8** ] threads
- $B_{KJ}$ gets loaded in shared memory by [ **8**, **32** ] threads
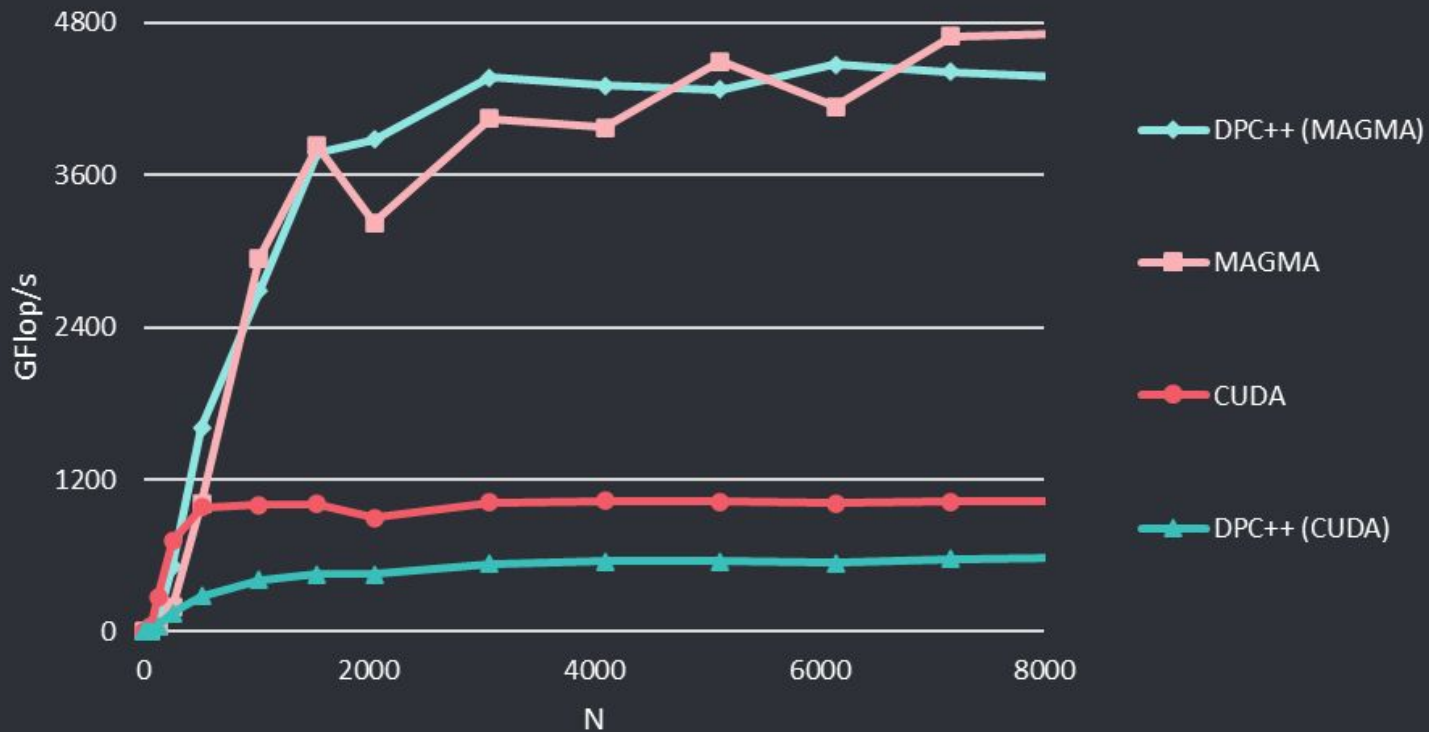- $C_{IJ}$ is held and computed in registers

# NVIDIA GEFORCE RTX 3060

# INTEL UHD GRAPHICS P630 [0x3e96]

# ADDITIONAL TEST PARAMETERS

|       | DIM_X | DIM_Y | DIM_M | DIM_N | DIM_K | DIM_XA | DIM_YA | DIM_XB | DIM_YB |
|-------|-------|-------|-------|-------|-------|--------|--------|--------|--------|
| cuda  | 16    | 16    | 96    | 96    | 16    | 32     | 8      | 8      | 32     |
| ker2  | 16    | 16    | 64    | 64    | 8     | 32     | 8      | 8      | 32     |
| ker11 | 12    | 4     | 48    | 48    | 2     | 24     | 2      | 24     | 2      |

# INTEL UHD GRAPHICS P630 [0x3e96]

# 8 CONCLUSION

## SUMMARY

- oneAPI is a promising approach for parallel programming across various architectures
- DPCT tool can be used successfully for an initial port of CUDA code to DPC++
- Large numerical libraries like MAGMA, originally written in CUDA to support Nvidia GPUs, can be easily translated to DPC++ to provide functional portability to different vendor GPUs, as well as multicore CPUs

- Initial migrated code tuned for Nvidia GPUs performs well on multicore CPUs
- Initial migrated code tuned for Nvidia GPUs retains performance on Nvidia GPUs
- Initial migrated code tuned for Nvidia GPUs performs poorly on the available Intel GPU
  - Tuning is required, but optimal parameters are difficult to find without further knowledge on the hardware design

41

- Full translation of MAGMA
- ICL account configuration
- Finding near optimal parameters for the Intel integrated GPU
- Testing migrated code on discrete Intel GPU upon release

# ACKNOWLEDGEMENTS

Home University

University of North Texas

# REFERENCES

[1]     *Advancing computing and data capabilities for scientific discovery and continued U.S. technological leadership.* Oak Ridge National Lab. https://www.ornl.gov/directorate/ccsd

[2]     https://thenounproject.com/search/icons/?iconspage=1&q=quantum

[3]     *Computing at LLNL.* Lawrence Livermore National Laboratory. https://computing.llnl.gov/

[4]     *NVIDIA HISTORY.* Nvidia. https://www.nvidia.com/en-us/about-nvidia/corporate-timeline/

[5]     *New Titan Supercomputer Named Fastest in the World.* Department of Energy. https://www.energy.gov/articles/new-titan-supercomputer-named-fastest-world-0

[6]     *June 2019.* The Top 500 List. https://www.top500.org/lists/top500/2019/06/

[7]     *June 2022.* The Top 500 List. https://www.top500.org/lists/top500/2022/06/

# REFERENCES

[8]     *Aurora: HPC and AI at Exascale.* Intel.
         https://www.intel.com/content/www/us/en/high-performance-computing/
         supercomputing/exascale-computing.html

[9]     *Compare Benefits of CPUs, GPUs, and FPGAs for Different oneAPI Compute
         Workloads.* Intel.
         https://www.intel.com/content/www/us/en/developer/articles/technical/c
         omparing-cpus-gpus-and-fpgas-for-oneapi.html#gs.83gstn

[10]    *Intel oneAPI Programming Overview.* Intel.
         https://www.intel.com/content/www/us/en/develop/documentation/onea
         pi-programming-guide/top/introduction-to-oneapi-programming/intel-on
         eapi-programming-overview.html

[11]    *Data Parallel C++: the oneAPI Implementation of SYCL*.* Intel.
         https://www.intel.com/content/www/us/en/developer/tools/oneapi/data-
         parallel-c-plus-plus.html#gs.83xmmq

# REFERENCES

[12]    *Intel® DPC++ Compatibility Tool.* Intel.
        https://www.intel.com/content/www/us/en/developer/tools/oneapi/dpc-compatibility-tool.html#gs.83zp77

[13]    *oneMKL.* Intel.
        https://spec.oneapi.io/versions/latest/elements/oneMKL/source/index.html

[14]    *What Is CUDA?* NVIDIA.
        https://blogs.nvidia.com/blog/2012/09/10/what-is-cuda-2/

[15]    *Compiling SYCL\* for Different GPUs.* Intel.
        https://www.intel.com/content/www/us/en/developer/articles/technical/compiling-sycl-with-different-gpus.html

[16]    *AMD EPYC™ 7742.* AMD.
        https://www.amd.com/en/products/cpu/amd-epyc-7742

## REFERENCES

[17]  *Intel® Xeon® Processor E5-2698 v4.* Intel.
https://ark.intel.com/content/www/us/en/ark/products/91753/intel-xeon-processor-e52698-v4-50m-cache-2-20-ghz.html

[18]  *GEFORCE RTX 3060 FAMILY.* Nvidia.
https://www.nvidia.com/en-us/geforce/graphics-cards/30-series/rtx-3060-3060ti/

[19]  *Intel UHD Graphics P630.* TechPowerUp.
https://www.techpowerup.com/gpu-specs/uhd-graphics-p630.c3676

[20]  *Intel® DevCloud.* Intel. Intel® DevCloud

Presentation Template:
Catalina, J. (n.d.). Minimal business. Free PowerPoint Template
&amp; Google Slides theme. SlidesCarnival. Retrieved July 5, 2022, from
https://www.slidescarnival.com/eleanor-free-presentation-template/308#preview