# V1-24: High-Productivity Computing in **Heterogeneous Systems**



#### SHREC Annual Workshop (SAW23-24)







FLORIDA

January 17-18. 2024

Faculty: Wu Feng

Students: Patrick Dewey, Atharva Gondhalekar, Ritvik Prabhu, Rashmi Ravindranath, Paul Sathre, Frank Wanye, Mukund Yadav

Number of requested memberships  $\geq 6$ 

## Goal

- Enable high-productivity computing in heterogeneous computing systems: CPU + {cpu, GPU, FPGA, TPU, ... }
  - Similar to DARPA <u>High-Productivity Computing Systems</u> program for <u>homogeneous</u> systems (e.g., Chapel, Fortress, X10) but for heterogeneous systems (e.g., Chapel, OpenCL, SYCL, oneAPI, VITIS)
  - Preferred Vehicle: Modern, Open Standard Languages & Runtimes
  - Case Studies: Applications and Benchmarks, e.g., Berkeley Dwarfs 

    OpenDwarfs (@VT)



#### Write ONCE, run ANYWHERE!





# **Background & Motivation**

 Extend our R&D to create and analyze an ecosystem of high-productivity tools, environments, and benchmarks for heterogeneous computing



- Challenges: How to productively ...
  - Program an application so it runs on many platforms?
  - Evaluate a processor architecture & compare it to others?
  - Develop back-end optimizations & know that they will work well?

Application-dependent





# **Recent Work: Performance & Productivity**

Sobel Filter on Intel Arria 10, AMD Alveo U250, and NVIDIA RTX 3090



\* Evaluated the same OpenCL kernel written for Arria 10 on U250 without any vendor-specific optimizations

 Rigorous Performance & Productivity Evaluation of Representative Apps (FFT, Jaccard similarity, biconjugate gradient stabilized method – BiCGSTAB, and graph algorithms) in Different Languages on Different Devices (CPUs, GPUs, and FPGAs)

(intel)

**XEON**<sup>®</sup>



OpenC

SYCL. PYNQ CHAPEL 🗯 PYNQ XILINX

oneAPT

BYU

FLORIDA

• University of

VIRGINIA TECH

Pittsburgh

# Approach

- Implement a diverse set of *application benchmarks* 
  - Regular vs irregular. Floating point vs integer. CPU- vs memory-intensive
- Characterize the *productivity* of a heterogeneous system
  - Kernel Development Time (KDT)  $\rightarrow$  Wall Clock Time
  - Source Lines Of Code (SLOC) and Code Convergence (CC)
- Characterize the *performance-vs-productivity* tradeoff
  - Performance Portability  $(\mathbf{P})$
  - Performance-Productivity Product (  $\Pi$  )
- Identify the best platform and associated ecosystem for productivity, performance, or both (across many apps)
- Enable further high-productivity research: automated co-scheduling at runtime and auto-generation of translators \_\_\_\_\_\_ Open Source \_\_\_\_\_\_ Closed Source \_\_\_\_\_\_ Closed Source

Open**CL** 







CHAPEL PYNQ

AMD RADEON PRO W7700





oneAPI

CUDA

Stratix 10

(2+1)

(2+1)

(2+1)

(0+2)

(0+2)

# **Proposed Tasks for V1-24**

(Memberships: Mandatory + Optional, e.g., 2+1)

- Task 1: FPGA Productivity and Performance (Speed/Space/Power)
- Task 2: GPU Productivity and Performance (Speed/Power)
- Task 3: Chapel GPU Parallel & Distributed Programming: Study, Analysis, Outreach
- Task 4: Simultaneous Co-scheduling of Heterogeneous Devices: CPU+GPU+FPGA
- Task 5: Auto-Generation of Source-to-Source Translators
  - One of {CUDA-to-OpenCL, CUDA-to-OpenMP, CUDA-to-Chapel, CUDA-to-SYCL, ... }
  - Baseline: Update of CU2CL: <u>CU</u>DA-<u>to</u>-Open<u>CL</u> Source-to-Source Translation







Mission-Critical Computing NSF CENTER FOR SPACE, HIGH-PERFORMANCE, AND RESILIENT COMPUTING (SHREC)

UF

FLORIDA

## Task 1: FPGA Productivity and Performance (Speed/Space/Power)

- Motivation
  - Multiple design-entry options (e.g., OpenCL, VITIS HLS, VITIS AI)
  - Performance & productivity tradeoffs of FPGA languages and tools relatively unknown
- Approach
  - Develop target apps via open-source & vendor toolchains
  - Evaluate performance & productivity of FPGA design options
- Tasks 1a, 1b, 1c: Productivity & performance of ...
  - a) Fast Fourier transform (FFT) with custom-length input via VITIS HLS
    - Target: Snickerdoodle board (ARM Cortex-A9 CPU + AMD Artix-7 FPGA)
  - b) Existing reinforcement learning (RL) model via VITIS AI toolchain
    - Target: Alpha Data PA100 (2x ARM Cortex-A72, 2x ARM Cortex-R5, 400x AI Engines)
  - c) FFT & RL via open-source toolchains, e.g., SYCL (currently unsupported)
  - d) "To-be-provided" RL model for inference









Tasks: Baseline & Optional (2+1)

Write Once, Run Anywhere? OpenCL and SYCL

Write once, run anywhere

8



## Task 2: GPU Productivity and Performance (Speed/Power)

- Motivation
  - High Performance *but* Low Productivity: Low-level C/C++-based languages → CUDA, OpenCL, SYCL
  - High Productivity but Less Performant: (1) Vendor libraries & tools. (2) High-Productivity Lang.
    - Issues: (1) Non-portable libraries, e.g., CUDA libraries not portable to AMD hardware. (2) Lack of hardwaresoftware interface control to achieve high performance. (3) Rigidity of pipelines and data abstractions
- Approach

ssion-Critical Computing

#### Think AMD MI300 CPU+GPU and NVIDIA GH200 CPU+GPU

 $js(A,B) = \frac{|A \cap B|}{|A \cup B|}$ 

A O B

Fourier

- Evaluate 1<sup>st</sup> & 3<sup>rd</sup>-party solutions to reach more devices with less code (without sacrificing performance)
- First-party solutions? Translators, e.g., Intel's dpct (CUDA $\rightarrow$ oneAPI/SYCL) or AMD's hipify (CUDA $\rightarrow$ HIP)
- Third-party? Languages, e.g., Chapel, OpenCL, SYCL. Libraries, e.g., PyTorch, Apache Arrow, C++ std::par
- Tasks 2a, 2b, 2c: Productivity & performance of openCL (SYCL)
  - a) Jaccard similarity & other graph apps (e.g., Gunrock)  $\rightarrow$  parallel & distributed  $\bigcirc$
  - b) Biconjugate gradient stabilized method (BiCGSTAB) OpenCL (SYCL.
  - c) Community detection  $\rightarrow$  graph clustering (see V3-24)  $\int_{OREAPI}$
  - d) FFT with custom-length (non-power-of-two) input (SYCL.

Tasks: Baseline & Optional (2+1) Write Once, Run Anywhere?

<sup>9</sup> OpenCL, SYCL, Chapel, oneAPI

# Task 3: Study, Analysis, and Outreach for Chapel GPU

- Motivation
  - Chapel 1.0 (2009) to Chapel 1.26 (2022)
    - High-Productivity Computing for Parallel & Distributed **CPU** Computing → Alternative to MPI or MPI+OpenMP
  - Rise of GPU Computing (since 2007)
    - Low-Level Languages (CUDA, OpenCL, ... ) + High-Level but Rigid/Non-Portable (PyTorch, cuBLAS, ... ) Libraries
- Approach
  - Study, analyze, & educate community on Chapel GPU for parallel & distributed CPU+GPU computing
- Tasks 3a, 3b, and 3c
  - a) Stress-test and evaluate updated and new features of the Chapel 1.33+ releases
  - b) Profile & analyze performance of Chapel 1.33+ on NVIDIA GPU with Jaccard similarity & OpenDwarfs
    - Chapel runtime overheads, Intra-node parallelism and **inter-node parallelism** (*a la* OpenMP/OpenACC/CUDA and MPI, respectively)

10

- Conduct on AMD/Intel CPU and AMD/Intel GPU, as time & access permits
- c) Testing, deployment, and education with Chapel 1.33+ ... for parallel computation class(es) at VT



Tasks: Baseline & Optional (2 + 1)

Write Once, Run Anywhere? Chapel







(intel)





## **Task 5: Auto-Generation of Source-to-Source Translators**

## Motivation

- Too many heterogeneous programming languages & devices in legacy scientific code. Non-portable.
- Alternative to "write once, run anywhere" languages? Source-to-source translation
  - Examples: CUDA to OpenCL; OpenMP to CUDA; etc., but such translators are MANUALLY coded.
- Goal: Automatically generate such automated source-to-source translators Write Once, Run Anywhere via Source-to-Source Xlation
- Approach
  - Leverage abstract syntax trees (ASTs) to "structure" CUDA code, ensuring that translations respect domain knowledge and best practices.
  - Combine neural machine translation (NMT) with neurosymbolic AI to decipher underlying logic of domainspecific CUDA code, enabling more accurate translation.



Instead of manually creating CU2CL (from V1-12) and manually updating it when the clang ecosystem changes, *auto-generate the CU2CL source-to-source translator*.

- Task 5a, 5b, 5c (CUDA OpenCL OpenMCL OpenMCL OpenMCL OpenMCL OpenMCL OpenMCL OpenMCL OpenMCL
  - a) Update CU2CL source-to-source translator from V1-12
  - b) Train NLP to auto-generate CUDA-to-OpenCL translator
  - c) Evaluate accuracy of original CU2CL to auto-generated CU2CL on pre-existing suite of application codes
  - d) Assess efficacy of auto-generator for other languages





Tasks: Baseline & Optional (0+2)

12

# Milestones, Deliverables, and Budget

- Major Milestones (Tasks: T1-T6)
  - T1: FPGA Productivity: **FFT via VITIS HLS+PYNQ and RL via VITIS AI**
  - T2: GPU Productivity: Jaccard Similarity & Other Graphs via OpenCL/SYCL
  - T3: Chapel 1.32: Synthesis and analysis of Jaccard similarity parallel and distributed
  - T4: Simultaneous Co-scheduling of Heterogeneity: CPU+GPU co-scheduling of irregular app
  - T5: Auto-Generation of Source-to-Source Translation: Al-generated CUDA-to-OpenCL translator

## Deliverables

- Monthly progress reports, along with mid-year and end-of-year full reports
- 2-3 publications at top-tier conference venues or journals

## Recommended Budget

- Minimum: 6 memberships (300 votes)
- Maximum: 13 memberships (650 votes)









## Conclusion

- Enable high-productivity computing in heterogeneous computing systems: CPU + {cpu, GPU, FPGA, TPU, ... } via open standards: OpenCL, SYCL, Chapel, and to a lesser degree, oneAPI
- Evaluate performance & productivity of representative apps (OpenDwarfs, FFT, Jaccard similarity, biconjugate gradient stabilized method – BiCGSTAB, and graph algorithms) on different devices (CPUs, GPUs, and FPGAs)

# **Member Benefits**



- Direct influence over processors & frameworks studied and apps & datasets used
- Direct benefit from new methods, tools, datasets, codes, models, and insights created as well as new metrics of evaluation
- Direct insights from R&D and analysis





# **APPENDIX**

- 1. FPGA: Proprietary (RTL & HLS) vs. C/C++ (OpenCL/SYCL/oneAPI)
- 2. GPU: Low-Level SIMD (CUDA/OpenCL/SYCL) vs. High-Level SIMD (Chapel)
- 3. Open-Source Approach Leveraging Open Source: Evaluation Metrics
- 4. Open-Source Approach Leveraging Open Source
  - Enabling Further High-Productivity Research
- 5. Details for Task 2c: Graph Clustering via oneAPI/SYCL (Related to V2-24)
- 6. Details for Task 3: Study, Analysis, and Outreach for Chapel GPU
- 7. Details for Task 4: Simultaneous Co-scheduling of Heterogeneity (CPU+GPU+FPGA)
- 8. Details for Task 4: Simultaneous Co-scheduling of Heterogeneity (CPU+GPU+FPGA): Irregular Apps





# FPGA: Proprietary (RTL & HLS) VS. C/C++ (OpenCL ...)

## Low-Level RTL (Verilog/VHDL)

#### Pros

- Direct control over timing requirements and design
- Maximum performance via via HW-specific optimization

#### Cons

- SLOW development time
- Tedious to read & understand

"machine language"

- More error-prone
- NOT portable Verilog.



CPU GPU FPGA

Mid-Level HLS (VITIS HLS & AI)

#### **Pros** and **Cons**

Fairly slow development time

VITIS AL

- Less tedious to read & understand
- Fairly errorprone
- **NOT portable** 
  - CPI **Still requires** low-level knowledge to prepare and invoke the kernel

High-Level C/C++ (OpenCL/SYCL)

#### Pros



- Easier to read and understand
- Less error-prone
- Portable (CPU, GPU, FPGA)



#### Cons

**FPGA** 

GPL

- Limited control over timing requirements
- Performance loss due to overhead of high-level abstractions
- Mixed vendor support for open standards

"high-level language"





"assembly language"

## **GPU: CUDA/OpenCL/SYCL vs. Chapel** (Parallel) (Parallel+Distributed)

#### Low-Level SIMD (CUDA/OpenCL/SYCL)

#### Pros

 Direct control over mapping. registers, and memory



 Maximum performance via hardware-specific optimization

#### Cons

Slower development time



- Tedious to read & understand
- Manual device/thread decomposition
- Manual distribution via MPI



#### High-Level SIMD+MIMD (Chapel)

#### Pros

Faster development time



- Easier to read and understand
- Automatic array/domain decomposition
- Natively extensible from CPU-parallel to heterogeneous-parallel and distributed

#### Cons

- Limited control over mapping, registers, and memory
- Performance loss due to overhead of high-level abstractions
- Limited developer resources and community









Performance Portability  $(\mathbf{P})$ 

RESILIENT COMPUTING (SHREC)







## Task 2c: Graph Clustering via oneAPI/SYCL

- Motivation
  - Accurate graph clustering → computationally expensive ×



- Available implementations of statistically robust graph clustering algorithms? CPU-oriented ×
  - Frontier supercomputer: 99% of FLOPS from GPUs
- Observation
  - CPU-only implementations sacrifice runtime gains of massively parallel GPUs (and low-power gains of FPGAs)
- Approach
  - Translate C++ graph clustering code (V2-24) to oneAPI/SYCL to run heterogeneously on CPU+GPU
  - Evaluate and compare runtime performance of C++ optimized CPU code to oneAPI/SYCL translation on different architectures: CPU-only, GPU-only, CPU+GPU, CPU+GPU with co-scheduling)



# Networking/Finance Bioinformatics Image: Constraint of the second sec



Graph clustering  $\rightarrow$  use cases across many domains  $\checkmark$ 



Intrusion/Fraud Detection

Epidemiology & Rece Drug Discovery S

Recommender Systems



# Task 3: Study, Analysis, and Outreach for Chapel GPU

## Motivation

- Chapel 1.30 (2023): Intro of GPU Support
- Need for a *portable high-level* approach that provides flexibility, interoperability, and performance of C-based solutions with the more productive and intuitive syntax of library- and Python-based solutions
- Approach
  - Study, analyze, & educate on Chapel GPU for parallel & distributed CPU+GPU computing
    - Intra-node parallelism across CPUs and GPUs from multiple vendors (unlike CUDA/HIP)
    - Multi-node distribution of tasks without an add-on library like MPI (unlike CUDA/HIP/SYCL/OpenCL)
    - Intuitive parallel abstractions like forall and automatically promoted and distributed array ops

- Tasks 3a, 3b, and 3c
  - a) Stress-test and evaluate updated and new features of the Chapel 1.32 release
  - b) Profile & analyze performance of Chapel 1.32
    - Device: NVIDIA GPU (No support for other GPUs yet.)
    - Apps: Jaccard similarity → OpenDwarfs for Chapel (i.e., ChapelDwarfs)
    - Environments: Intra-node and inter-node Frontier Compute Blade (Two Nodes)





The intersect of A & B

The union of A & B

division

c) Testing, deployment, and education with Chapel
 1.32 for parallel computation class(es) at VT





## **V1**

## **Task 4: Simultaneous Co-scheduling of Heterogeneity** (CPU+GPU+FPGA)

- Motivation (in Detail)
  - CPU typically idle when GPU executes kernel
    - Why not distribute a portion of workload to CPU?!
- Our Prior Work
  - Automated distribution of workload of regular apps for better performance
  - Context: OpenMP OpenACC
    Directors for Accelerators
- Proposed Research
  - Automated distribution of workload for irregular apps for better performance
  - Context: 
     SYCL
  - Heterogeneity: CPU+GPU, CPU+FPGA, CPU+GPU+FPGA





## **Task 4: Simultaneous Co-scheduling of Heterogeneity** (CPU+GPU+FPGA): Irregular Apps

Graph Apps from IARPA AGILE Program [1]

- Jaccard similarity (JS) in graph datasets
  - Parallel JS computation across all vertex pairs that form edges of the graph



[1]: <u>https://www.iarpa.gov/images/PropsersDayPDFs/AGILE/AGILE\_Program\_Workflows\_FINAL.pdf</u>

- Triangle counting (TC) in graph datasets
  - Evaluate the total number of unique triangles formed by the edges in the graph



Image source: https://www.computer.org/csdl/journal/td/2017/12/08000612/13rRUxYINeZ

## Scientific Computing App

- Biconjugate gradient stabilized (BiCGSTAB)
  - Sparse linear system solver with irregular memory access pattern
  - Used in computational fluid dynamics

ion-Critical Computing





V1

## Task 5: Auto-Generation of Source-to-Source Translators

## Motivation

- Too many heterogeneous programming languages & devices in legacy scientific code. Non-portable.
- Alternative to "write once, run anywhere" languages? Source-to-source translation
  - Examples: CUDA to OpenCL; OpenMP to CUDA; etc., but such translators are MANUALLY coded.
- Goal: Automatically generate such automated source-to-source translators
- Approach
  - Leverage abstract syntax trees (ASTs) to "structure" CUDA code, ensuring that translations respect domain knowledge and best practices.
  - Combine neural machine translation (NMT) with neurosymbolic AI to decipher underlying logic of domainspecific CUDA code, enabling more accurate translation.



- Existing Source-to-Source Translators
  - OpenMP → CUDA, OpenCL, OpenACC, ISPC, MPI
  - OpenACC  $\rightarrow$  OpenMP
  - CUDA → OpenCL, OpenMP, OpenACC, HIP
  - OpenCL  $\rightarrow$  CUDA
- Task 6a, 6b, 6c
  - a) Update CU2CL source-to-source translator from V1-12
  - b) Train NLP to auto-generate CUDA-to-OpenCL translator
  - c) Evaluate accuracy of original CU2CL to autogenerated CU2CL on pre-existing suite of application codes
  - d) Assess efficacy of auto-generator for other languages



