# V1-26: Performance, Power, and Precision, of Heterogeneous Systems

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

**SHREC Annual Workshop (SAW25-26)**

University of Pittsburgh

BYU BRIGHAM YOUNG UNIVERSITY

VIRGINIA TECH

UF UNIVERSITY of FLORIDA

January 13-14, 2026

Faculty:  Wu Feng and Krish Sundararajah

Students and Post-Students:

Nabayan Chaudhury, Atharva Gondhalekar, Poorna Gunathilaka, Pratheek Prakash Shetty, Ritvik Prabhu, Eric Rippey, Paul Sathre
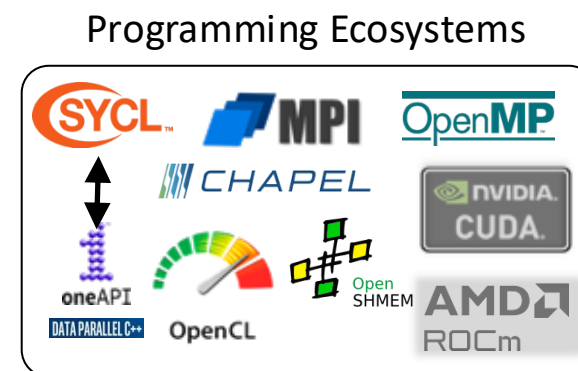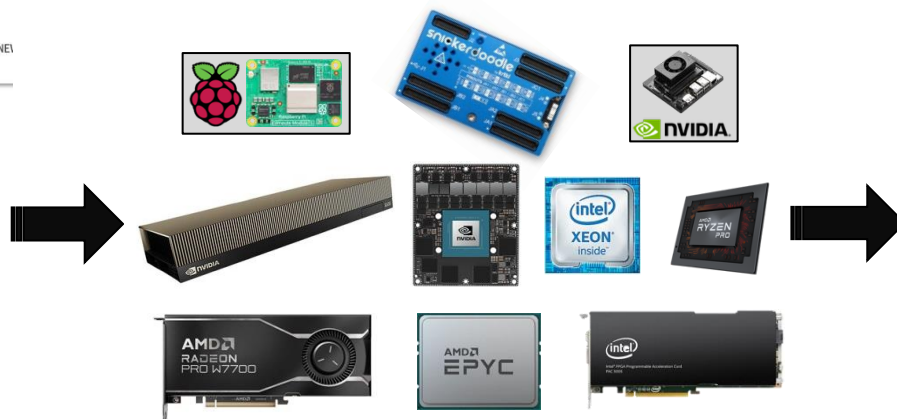
Number of requested memberships ≥ 6

# Goal: Productivity → Performance, Power, Precision

- Tuning & optimization of performance, power, and precision (manual → automatic) for productivity in **heterogeneous** computing systems:  CPU + {CPU, GPU, TPU, ...}
    - Akin to DARPA HPCS program for <u>homogeneous</u> systems (e.g., Chapel, Fortress, X10) but for <u>heterogeneous</u> systems (e.g., Chapel, oneAPI → SYCL, OpenSHMEM)
    - Preferred Vehicle: Modern, Open Standard Languages & Runtimes → write once, run anywhere
    - Metrics of Evaluation: Performance, Power / Energy Efficiency, and Precision (e.g., int vs. SP vs. DP)



Programming Ecosystems

# Motivation & Background

### Tools & Environments



### Programming Ecosystems



### Benchmarks



Community Detection

## Performance (still matters but …)



logarithmic axis

ILP + DLP + TLP + RPP

## Power (becoming an issue …)



Evolution of Power Consumption & Dissipation per Rack (2000-2030)

logarithmic axis

9 kW/rack = 126 kW

Highest-ranked commodity supercomputer in the USA on the Green500

## Precision (can address both)



Single-Chip Inference Performance - 1000X in 10 years

RPP: reduced-precision parallelism

logarithmic axis

ILP: instruction-level parallelism
DLP: data-level parallelism
TLP: thread-level parallelism

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

University of Pittsburgh · BYU · Virginia Tech · UF
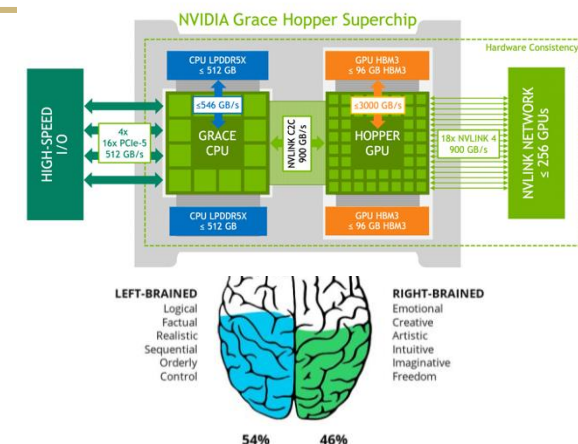
# Background: Performance & Power



- **Ranking of fastest supercomputers** (Nov. 2025) based on
  - High-performance LINPACK (HPL) → **REGULAR** workloads
  - High-performance conjugate gradient (HPCG) → **IRREGULAR** workloads
  - GPU: 60%-85% of peak for **regular** workloads; 0%-5% for **irregular** workloads



HPL vs HPCG Efficiency on Top Supercomputers (Nov 2025)

| Supercomputer | TOP500 Rank | HPCG Rank | Accelerator | Peak Performance $R_{peak}$(PFlop/s) | Power (MW) |
|---|---|---|---|---|---|
| El Capitan (LLNL, USA) | 1 | 1 | AMD Instinct MI300A GPUs | 2821.1 | 29.7 |
| Frontier (ORNL, USA) | 2 | 3 | AMD Instinct MI250X GPUs | 2055.0 | 24.6 |
| Aurora (ANL, USA) | 3 | 4 | Intel Data Center Max GPUs | 1980.0 | 38.7 |
| Fugaku (RIKEN, Japan) | 7 | 2 | Fujitsu A64FX CPUs (no GPUs) | 537.2 | 29.9 |
| LUMI (CSC, Finland) | 9 | 5 | AMD Instinct MI250X GPUs | 531.5 | 7.1 |

- **Optimizations for performance and power**
  - Mixed-precision, tuning voltage/frequency
- **Simultaneous co-scheduling for heterogeneity**
  - CPU + GPU co-scheduling, as appropriate

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

# Approach

- Tune and optimize the *performance* of a heterogeneous system
- Characterize the *power* and *energy* of a heterogeneous system
  - GPU device power via vendor tools (e.g., nvidia-smi)
  - Total system power via power meters and software tools (e.g., RAPL)
- Characterize the *performance-vs-power* tradeoff
  - Performance per watt or energy-delay product
  - Power vs. runtime → energy
- Evolve the diversity of *app benchmarks* to evaluate the above
  - Regular vs irregular. Double vs. single precision. CPU- vs memory-intensive
- *Identify the best platform(s)* and associated ecosystem(s) for performance, power, and/or precision (across many apps)
- *Enable further performance, power, and precision-aware research*:  automated co-scheduling at run-/compile-time, performance vs. power vs. precision tradeoff
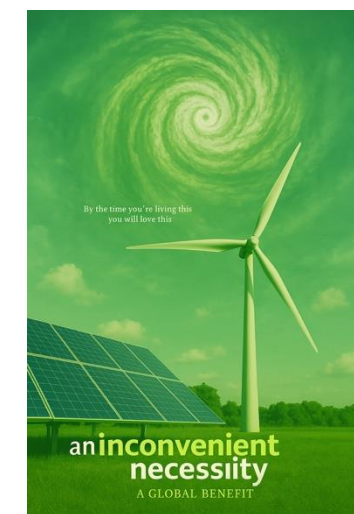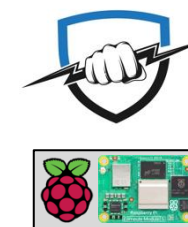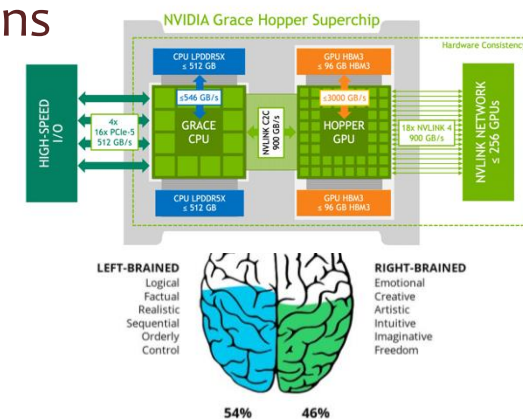
Open Source                                    Closed Source

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

# Proposed Tasks for V1-26

- Task 1: Performance, Power/Energy, & Precision for *Parallel* Hetero Computing (**2**+5)
  - Task 1a: Energy-Efficient/Energy-Dominant Computing for Irregular Applications
  - Task 1b: @Runtime: Simultaneous Co-scheduling on Heterogeneous Devices
  - Task 1c: @Compiler: Simultaneous Co-scheduling on Heterogeneous Devices
  - Task 1d: Portable Runtimes for Heterogeneous Task Graphs
  - Task 1e: Concurrent Data Structures for the GPU

- Task 2: High-Performance *Distributed* Computing with GPUs (**2**+2)
  - Heterogeneous PGAS vs MPI+X for Large-Scale Compute

- Task 3: Performance & Power/Energy for *Edge* Computing (**1**+2)
  - Task 3a: Analysis of Portable Kernel Pipelines for Edge Devices
  - Task 3b: Modeling of Power/Energy Draw via Generative AI

Courtesy: ChatGPT

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

University of Pittsburgh

BYU BRIGHAM YOUNG UNIVERSITY

VIRGINIA TECH

UF UNIVERSITY of FLORIDA

# Task 1a: Energy-Efficient/Energy-Dominant Computing

## Motivation

- Power & energy are now *first-order constraints*
  - Hyperscale data center guzzles 20 MW – 50 MW on avg. (with energy consumption ~ 32 TWh)
  - Modern supercomputer uses 10 MW – 40 MW



**The power of data** [2]

US, data-centre energy consumption, TWh By type
- Hyperscalers (red)
- Other (grey)

*Low  †High  ‡Estimate  §Forecast

Source: Lawrence Berkeley National Laboratory



**AI's Power Drive is Fuelling a Green Boom**

Hyperscale data centres are guzzling electricity and water at an alarming rate. For one, Google's sustainability report showed that its electricity use went up 27% in 2024 to 32 Terawatt-hours. Barclays says that hyperscalers Google, Microsoft and Meta are on track for their 7th consecutive year of 25% growth in power demand. Their Scope 3 emissions have also risen 23%-100% from the 2019 baseline. But this has also translated into a clean energy drive. **Himanshi Lohchab** curates facts and figures linked to AI's impact on power use.

HYPERSCALER GLOBAL ELECTRICITY USE
- Google, Microsoft, Meta

*Meta is yet to publish 2024 figures
Source: Compnay reports & Barclays Research

AI is projected to double the electricity use of data centres by 2030, reaching 945 terawatt hours – more than Japan's consumption today

- In 2024, data centres made up just 1.5% of global electricity consumption, expected to grow to 3% by 2030
- In US, power consumption by data centres will account for almost half of overall growth in demand by 2030
- In US, data processing will consume more energy than manufacturing of goods like aluminium, steel, cement and chemicals
- India has 2GW of installed data centre capacity consuming electricity equivalent to 6.5 million Indian households

INDIA'S TOTAL DATA CENTRE CAPACITY AND ELECTRICITY GENERATION MIX, 2020-2030
- Total installed data centre capacity
- Coal, Solar PV, Hydro, Wind, Natural gas, Nuclear, Bioenergy

Total installed capacity of date centres in India is set to double by 2030; While coal dominates the elecitity mix in India, the share of renewables increases to 35% by then

THE CLEAN ENERGY BOOM
- **Google marked** a record year procuring 8GW of clean energy in 2024 including solar and wind projects in Mumbai and Delhi
- **Clean energy** procurement marked a record 68GW of deals in 2024 with data centres accounting for 17GW
- **Amazon, Google,** Meta and Microsoft now have combined clean energy portfolio of 84 GW across 29 markets
- In May 2024, Microsoft announced a $10-billion deal with Brookfield to develop green energy capacity of 10.5GW, the largest ever
- **Asia-Pacific secured** a total of 27 GW of renewable capacity, up 60% year on year
- **India accounted** for more than 69% of this capacity, surpassing the US to become the largest market globally at 18 GW

Sources: INTERNATIONAL ENERGY AGENCY, GOOGLE SUSTAINABILITY REPORT, S&P GLOBAL COMMODITY INDEX, GOLDMAN SACHS RESEARCH

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

University of Pittsburgh  BYU BRIGHAM YOUNG UNIVERSITY  VIRGINIA TECH  UF UNIVERSITY of FLORIDA

# Task 1a: Energy-Efficient/Energy-Dominant Computing

## Motivation

- Power & energy are now *first-order constraints*
  - Hyperscale data center guzzles 20 MW – 50 MW on avg. (with energy consumption ~ 32 TWh)
  - Modern supercomputer uses ~ 10 MW – 40 MW
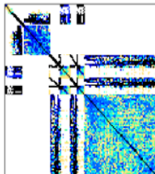- **Power- and energy-aware approaches needed** to align application behaviour with system-level power budgets

| Supercomputer | TOP500 Rank | HPCG Rank | Accelerator | Peak Performance $R_{peak}$(PFlop/s) | Power (MW) |
|---|---|---|---|---|---|
| El Capitan (LLNL, USA) | 1 | 1 | AMD Instinct MI300A GPUs | 2821.1 | 29.7 |
| Frontier (ORNL, USA) | 2 | 3 | AMD Instinct MI250X GPUs | 2055.0 | 24.6 |
| Aurora (ANL, USA) | 3 | 4 | Intel Data Center Max GPUs | 1980.0 | 38.7 |
| Fugaku (RIKEN, Japan) | 7 | 2 | Fujitsu A64FX CPUs (no GPUs) | 537.2 | 29.9 |
| LUMI (CSC, Finland) | 9 | 5 | AMD Instinct MI250X GPUs | 531.5 | 7.1 |

## Approach: Target Apps, Platforms, and Optimizations

- Apps: Jaccard similarity (JS), conjugate gradient (CG), triangle counting (TC), **[ your workload here ]** *(see appendix)*
- Platforms: CPU/GPU/APU from AMD, Intel, or NVIDIA
- App-level optimizations
  - (1) mixed- or reduced-precision computing and (2) fine-grained domain decomposition

## Milestones

1. Application suite for tuning & optimization (i.e., power and energy efficiency)
2. Software-based scripting framework for power and energy measurement
3. Profiling database of power and energy data via hardware & software meters
   - Power → vendor tools (e.g., nvidia-smi, rocm-smi)
   - Energy → integration of power over runtime

$$\frac{1}{h^2} \begin{pmatrix} -2 & 1 & 0 & 0 & \dots & 0 \\ 1 & -2 & 1 & 0 & \dots & 0 \\ 0 & 1 & -2 & 1 & \dots & 0 \\ \vdots & & & & & \vdots \\ 0 & \dots & & 1 & -2 & 1 \\ 0 & \dots & & & 1 & -2 \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ \vdots \\ u_{N-3} \\ u_{N-2} \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ f_3 \\ \vdots \\ f_{N-3} \\ f_{N-2} \end{pmatrix}$$
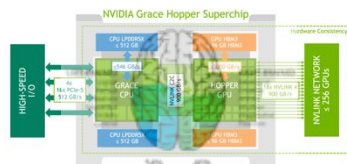
Domain decomposition

Validation of software meters w.r.t. hardware meters?

Tasks: Baseline & Optional
( 1 + 0 )

## Motivation

- **Supercomputing nodes lean towards GPUs** (3:1 GPU-to-CPU) → **presumption: more GPUs will benefit every application**
  - Alas, **NOT** the case for **irregular applications** → non-coalesced memory accesses, branch divergence, heavy data movement
- Physiologically, we use two brains simultaneously – left and right What about "in silico"? CPU **and** GPU simultaneously?
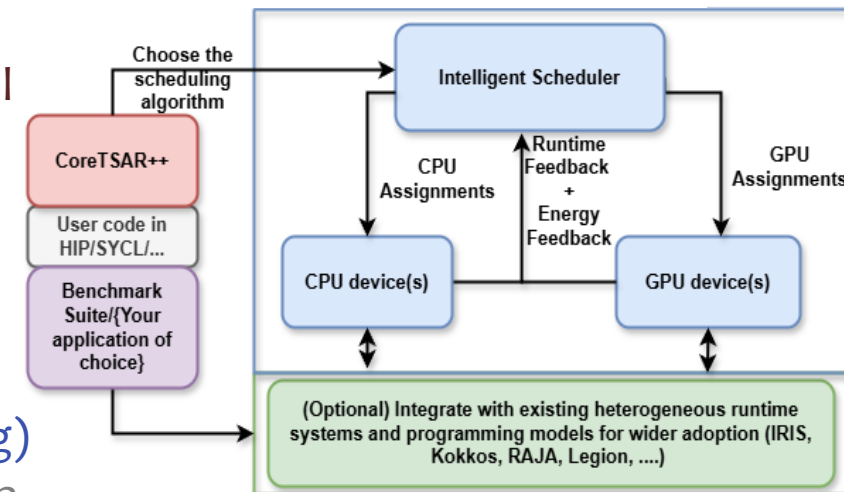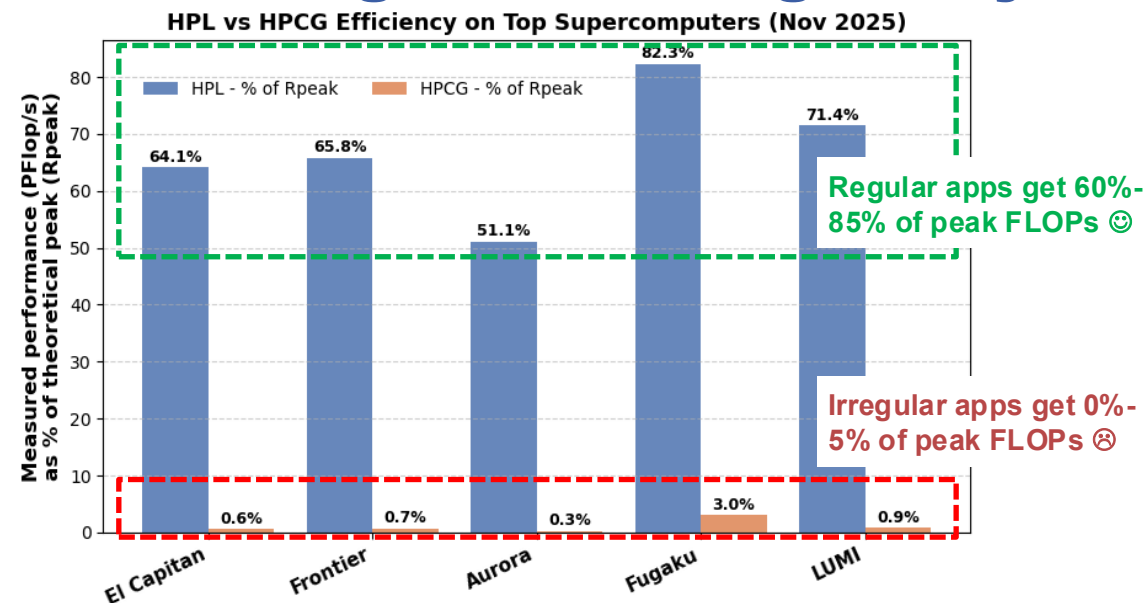


NVIDIA Grace Hopper Superchip

## Approach

- Efficacy of (manual) simultaneous co-scheduling of apps
  - ✓ Broader evaluation across a larger collection of benchmarks (see appendix)
- Build an automated scheduler that fits an accelerator programming model (e.g., OpenMP or [backend of interest] on any xPU) w/ minimal hassle
- Enable automated scheduling for runtime performance & energy cost

## Milestones      0. Manual "oracle" co-scheduling of resources

1. CoreTSAR++ automated scheduler with accelerator programming model (e.g., OpenMP on *any* CPU+GPU; optionally, AMD HIP, oneAPI/SYCL Intel)
2. Power measurement methodology (prelude to energy-efficient scheduling)
3. EnergyTSAR++ → New scheduler to minimize runtime energy consumption



**HPL vs HPCG Efficiency on Top Supercomputers (Nov 2025)**

Regular apps get 60%-85% of peak FLOPs ☺

Irregular apps get 0%-5% of peak FLOPs ☹



Tasks:  Baseline & Optional
( **1** + 1 )

# Task 1c: @Compile Simultaneous Co-scheduling for Heterogeneity

## Data-Layout Optimizations for Irregular Apps



- Motivation
  - Data layout often influences simultaneous co-scheduling of an app
    - Changing data layout → a schedule favoring different execution targets
  - Challenge: Realizing different schedules by changing data layout

- Approach
  - Template-based embedded DSL (in C++) to abstract data layout from description of computation
  - Compiler plugin to manipulate data layout of program (e.g., MLIR dialect) and generate code accordingly
  - Compiler support to guide co-scheduling of CoreTSAR++ runtime



- Milestones

  1. Bidirectional compiler support for CoreTSAR++ runtime to provide data layout-aware scheduling for hetero execution targets.

  2. Implementation & evaluation of irregular applications from HeCBench in this compilation pipeline.



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

Tasks:  Baseline & Optional
( **0** + 1 )

10

# Task 1d: Portable Runtimes for Heterogeneous Task Graphs
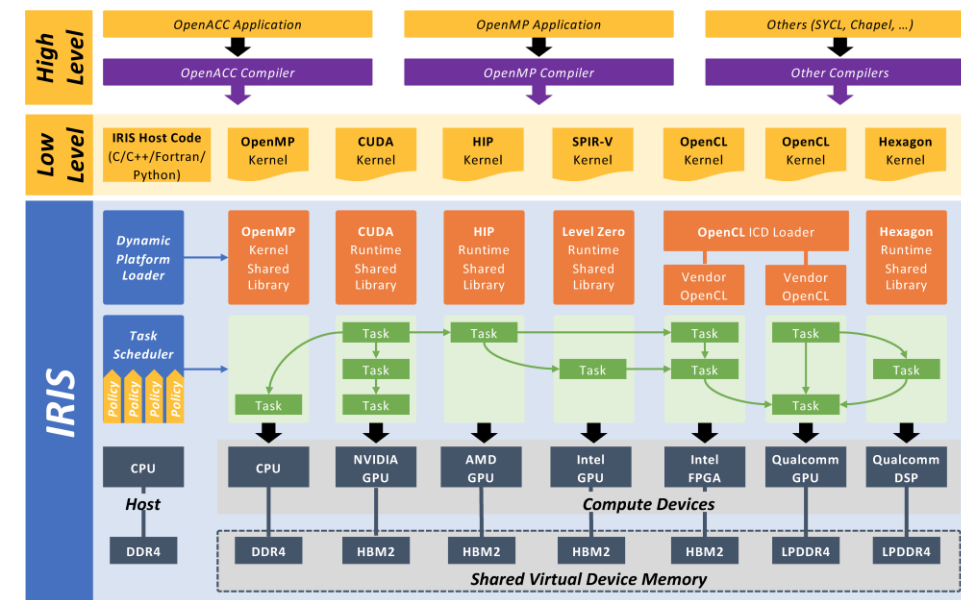
## Motivation

- Modern HPC requires device- and system-aware mapping of kernels, communication, and I/O to hardware
- Hardware *migration* (translate, remap, retune) is a significant *cost* [time, human, $$] which slows mission progress
- Portable languages help reduces translation component
- Remapping and re-tuning for new hardware still takes effort!
  - *Intelligent heterogeneous tasking systems* can help!
  - Given a portable representation, model and predict tradeoffs in mapping kernels to different hardware in the system



## Approach

- Implement SHREC-related applications using either in OpenMP / OpenARC, or emerging UniSYCL compiler
- Leverage and evaluate the *IRIS portable heterogeneous tasking system*'s ability to achieve high performance

## Milestones

1. Identify and migrate/implement a SHREC workload in the IRIS runtime, analyze perf./prod. ($\Pi$) (0.5)
2. Evaluate perf./prod. ($\Pi$) on traditional heterogeneous HPC (CPU+GPU, **homogeneous across** nodes) (1)
3. Evaluate perf./prod. ($\Pi$) on *multiply*-heterogenous HPC (CPU+**X, where X differs between nodes**) (1)
4. Evaluate perf./prod. ($\Pi$) w/ edge+centralized hybrid workloads w/ heterogeneous platforms (2)
   (i.e. data collection/reduction at the low-power edge, tightly coupled to high-power centralized analysis)

Tasks: Baseline & Optional
( **0** + 2 )

# Task 1e: Concurrent Data Structures for the GPU

## Motivation

- CPU concurrent data structures? MATURE (e.g., see Michael & Scott)
- GPU concurrent data structures? NASCENT to NON-EXISTENT
- GPU concurrency bottlenecks NOT addressed by current abstractions



**Obstruction-Free** *progress only if isolated*   **Lock-Free** *system progresses individual may starve*   **Wait-Free** *all thread completes in ≤ k steps* ≤ k steps (bounded)

## Approach

**Analysis of Concurrent Queues for CPU**

- Evaluate CPU queue designs

  ... with variations in progress guarantees (lock-free / wait-free)

| Michael & Scott Queue | Yang & Mellor-Crummey Queue |

**Synthesis of Concurrent Queues for GPU**

- Refactor for GPU / CPU+GPU architecture
  - ✓ Discrete: AMD MI200, NVIDIA A100/H200
  - ✓ Fused: AMD MI300A, NVIDIA GH200

SIMT + atomics

| Refactored M&S Queue | Refactored Y&M Queue |

## Milestones

1. Adaptation of CPU concurrent queues for GPU
   - Wait-Free-Queue (WFQ), Fetch-and-add based Queue (F&AQ)
2. Synthesis of GPU concurrent queues, e.g., bounded memory, cache-aware
3. Evaluation via microbenchmarks (BFS) and application (ray tracing)

**Breadth-First Search**        **Ray-Tracing**

# Task 2: Heterogeneous PGAS vs MPI+X for Large-Scale Compute
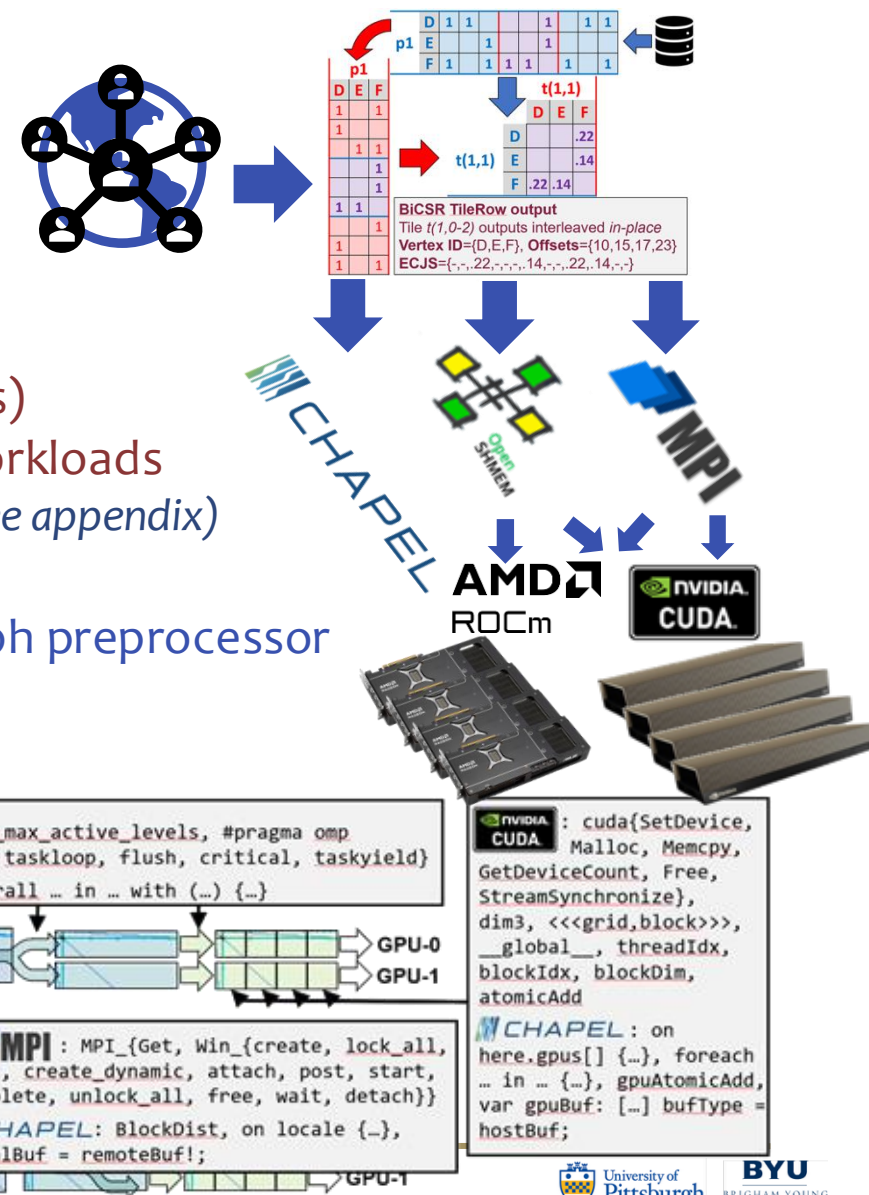
## Motivation

- Proliferation of programming models for *distributed GPUs*
  - MPI+CUDA/HIP, Chapel (PGAS), SHMEM (Open-, NV-/ROCm-)
- Scale-out evaluation needs new inputs too large for 1 GPU (> 100GB)
  - Small (<~1B-edge) graphs starve the GPUs of work when partitioned

## Approach

- Implement new preprocessor for *web-scale graphs* (~10B-100B edges)
- Refine and rigorously compare & contrast distributed GPU graph workloads
  - Jaccard similarity (MPI+X, Chapel, OpenSHMEM), … **[ your app here ]** (*see appendix*)

## Milestones

1. LWA (web-scale) → *bidirectional* compressed sparse row (CSR) graph preprocessor
2. Comparative language analysis
   - Intra- and inter-node performance, productivity, power
   - HIP for AMD GPU (MPI+X and OpenSHMEM)
   - CPU- vs. GPU-driven communication models
3. Integration of intra-node hybrid JS approaches
   - Coarse/fine, CPU/GPU, 2d kernels, co-scheduled
4. Additional graph workloads: triangle count, k-truss, etc.
5. Explore reduced-width quantization of intersection kernels



Tasks:  Baseline & Optional

( 2 + 2 )

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

University of Pittsburgh

BYU BRIGHAM YOUNG UNIVERSITY

VIRGINIA TECH

UF UNIVERSITY of FLORIDA

# Task 2: Heterogeneous PGAS vs MPI+X vs SHMEM

- Capturing high-level language tradeoffs for parallel & distributed computing (work in progress)

| Feature | Chapel | SHMEM | | MPI+X |
| --- | --- | --- | --- | --- |
| | | Open- | NV- | |
| GPU-driven communication | NO | NO | YES | Depends on implementation |
| Consistent GPU/node API | YES | NO | Depends | NO |
| Vendor-neutral API | YES | YES | NO | Depends on kernel language |
| Independently-sized per-node GETable allocations | YES | NO | NO | Dynamic-only |
| Exposed thread blocking | NO | YES | ? | YES |
| Node-specific partial data | YES | If same size | If same size | YES |

"Can I transfer GPU pointers without manually staging on the CPU?"

"Do I copy between nodes and between the CPU and GPU with the same API?"

"Can my code run on another vendor's hardware without rewrites?"

"Can I allocate different sizes of communicable data on different ranks?"

"Can I force a user-facing thread to spin in the foreground?"

"Can I allocate *non-overlapping* subsets of data across ranks, and still communicate them?"

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

Tasks: Baseline & Optional
( 2 + 2 )

University of Pittsburgh    BYU BRIGHAM YOUNG UNIVERSITY
VIRGINIA TECH    UNIVERSITY of FLORIDA

# Task 3a: Analysis of Portable Kernel Pipelines for Edge Devices
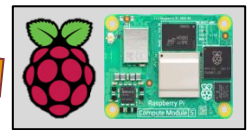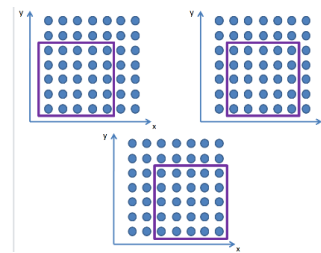
## Motivation

Datacenter Servers

- Proliferation of edge devices creates a data bottleneck with centralized processing
- Moving pre-filtering and other compute to the edge reduces aggregate bandwidth and storage
  - Issue: Wide range of edge devices that require *different* programming approaches
  - Challenge: Is there a robust path from portable HPC languages to low-power edge devices?
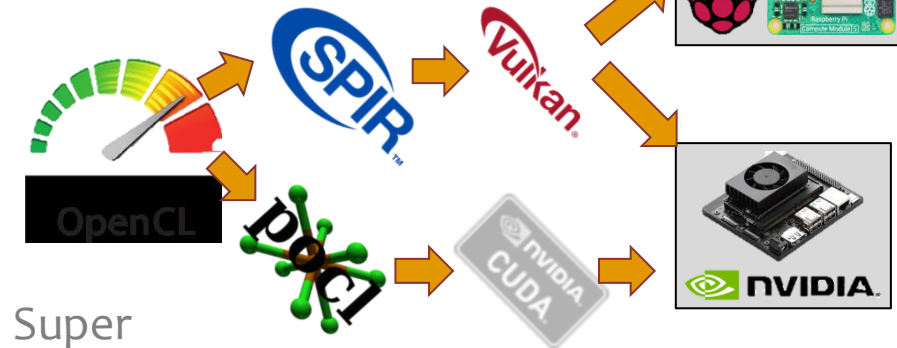
Edge Devices

## Approach

- Leverage portable, open HPC standards, and open-source toolchains to compute on edge GPU(s)
  - Examples: OpenCL, SYCL, SPIR, Vulkan
- Metrics: performance, power, productivity, and performance/power
- Platforms: Raspberry Pi Compute Module (CM) 5 and Nvidia Jetson Orin Nano Super
- Workloads: **Estimation of signal parameters via rotational invariant techniques (ESPRIT), FFT convolution,** [ **your workload here** ], multiple signal classification (MUSIC)

## Milestones

1. Raspberry Pi CM5 via OpenCL C → SPIR-V → Vulkan (0.5)
   - FFT convolution, ESPRIT → Perf./Power (Π) analysis

2. Nvidia Jetson Orin Nano via OpenCL C → SPIR-V → Vulkan (0.5)
3. MUSIC algorithm on Raspberry Pi CM5 and Nvidia Jetson Orin Nano Super

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

Tasks: Baseline & Optional
( 1 + 1 )

University of Pittsburgh   BYU BRIGHAM YOUNG UNIVERSITY
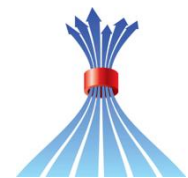VIRGINIA TECH   UNIVERSITY of FLORIDA

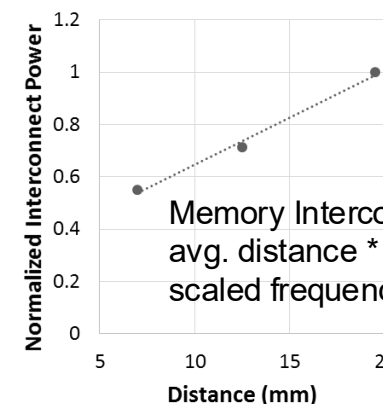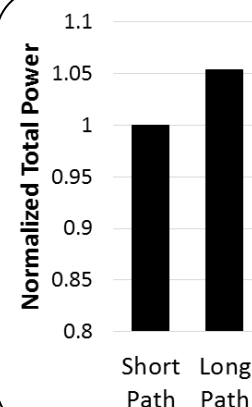# Task 3b: Modeling of Power/Energy Draw via Generative AI

## Motivation

- **Environment: Tactical computing on the edge**
  (i.e., move from datacenter to the edge)

- **Challenge: Edge devices energy-constrained**
  - How to deploy generative AI models that are constrained by the "AI Memory Wall," where energy consumption is dominated by data movement

- **State of the Art**
  - Reliance on *oversimplified* linear models that do NOT capture complex, non-linear dynamics of today's GPU

## Approach

- Memory-centric energy-modeling framework that integrates actual measurements and intelligent optimization to enable accurate prediction and systematic minimization of energy use in transformer-based models
  - Memory access can be ***100-200 times more energy-intensive*** than computation
  - Re-orientation of the energy optimization problem around data movement, creating new pathways to deploy powerful foundational models on edge & tactical hardware

Datacenter Servers

Edge Devices
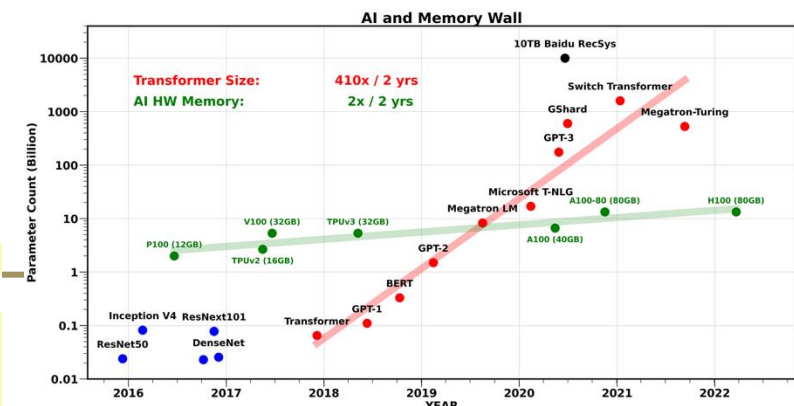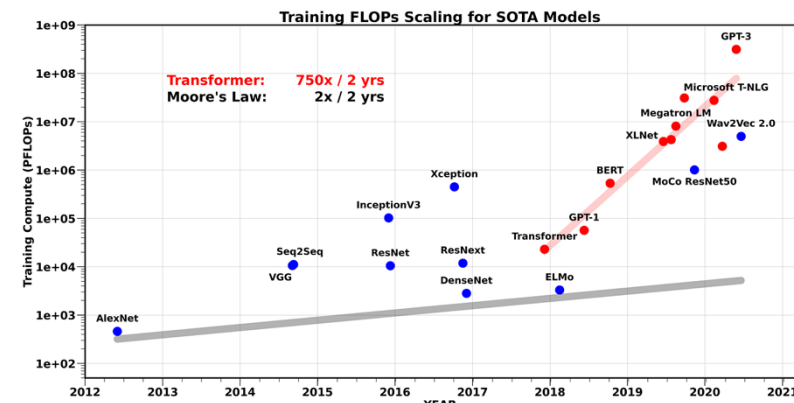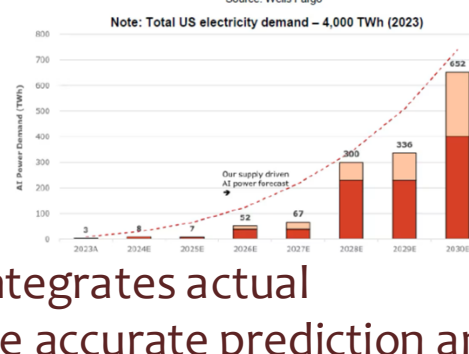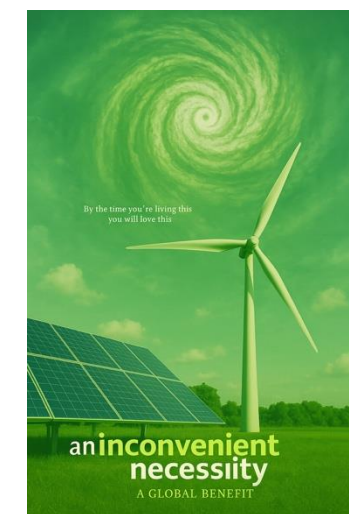


Source: Adhinarayanan, Paul, Greathouse, Huang, Pattnaik, Feng. *IISWC*, 2016. **Best Paper Award**

Memory Interconnect Power = Energy/bit/mm * avg. distance * avg. bits/sec * scaled voltage$^2$ * scaled frequency * avg. toggle rate

**Summary of GenAI demand forecast**
Source: Wells Fargo
Note: Total US electricity demand – 4,000 TWh (2023)

**Training FLOPs Scaling for SOTA Models**

Transformer: 750x / 2 yrs
Moore's Law: 2x / 2 yrs

**AI and Memory Wall**

Transformer Size: 410x / 2 yrs
AI HW Memory: 2x / 2 yrs

All generative AI queries, all users per year
**5,100,000,000,000**

Electricity required
**15 TWh**

All Gen AI → 🌐 → 365 = Equivalent to the annual output of two typical nuclear reactors

Tasks: Baseline & Optional
( **0** + 1 )

# Proposed Tasks for V1-26

- Task 1: Performance, Power/Energy, & Precision for *Parallel* Hetero Computing (**2**+5)
  - Task 1a:  Energy-Efficient/Energy-Dominant Computing for Irregular Applications
  - Task 1b:  @Runtime: Simultaneous Co-scheduling on Heterogeneous Devices
  - Task 1c:  @Compiler: Simultaneous Co-scheduling on Heterogeneous Devices
  - Task 1d:  Portable Runtimes for Heterogeneous Task Graphs
  - Task 1e:  Concurrent Data Structures for the GPU

- Task 2: High-Performance *Distributed* Computing with GPUs (**2**+2)
  - Heterogeneous PGAS vs MPI+X for Large-Scale Compute

- Task 3: Performance & Power/Energy for *Edge* Computing (**1**+2)
  - Task 3a:  Analysis of Portable Kernel Pipelines for Edge Devices
  - Task 3b:  Modeling of Power/Energy Draw via Generative AI

Courtesy: ChatGPT

# Appendix

# Background & Motivation

- Extend our R&D to create and analyze an ecosystem of *tools, environments, and benchmarks for heterogeneous computing*
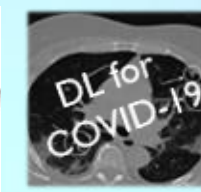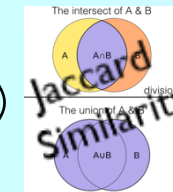
Past CHREC- and SHREC-funded R&D



Tools & Environments

CoreTSAR  FLOCL  CU2CL  MetaCL

Benchmarks

OpenDwarfs  Jaccard Similarity  DL for COVID-19

Devices

Programming Ecosystems

OpenMP  CHAPEL  OpenACC  oneAPI

OpenCL  MPI  SYCL  nvidia CUDA  AMD ROCm  Open SHMEM

- **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?
    - Ensure the system's power and precision/accuracy constraints are met?
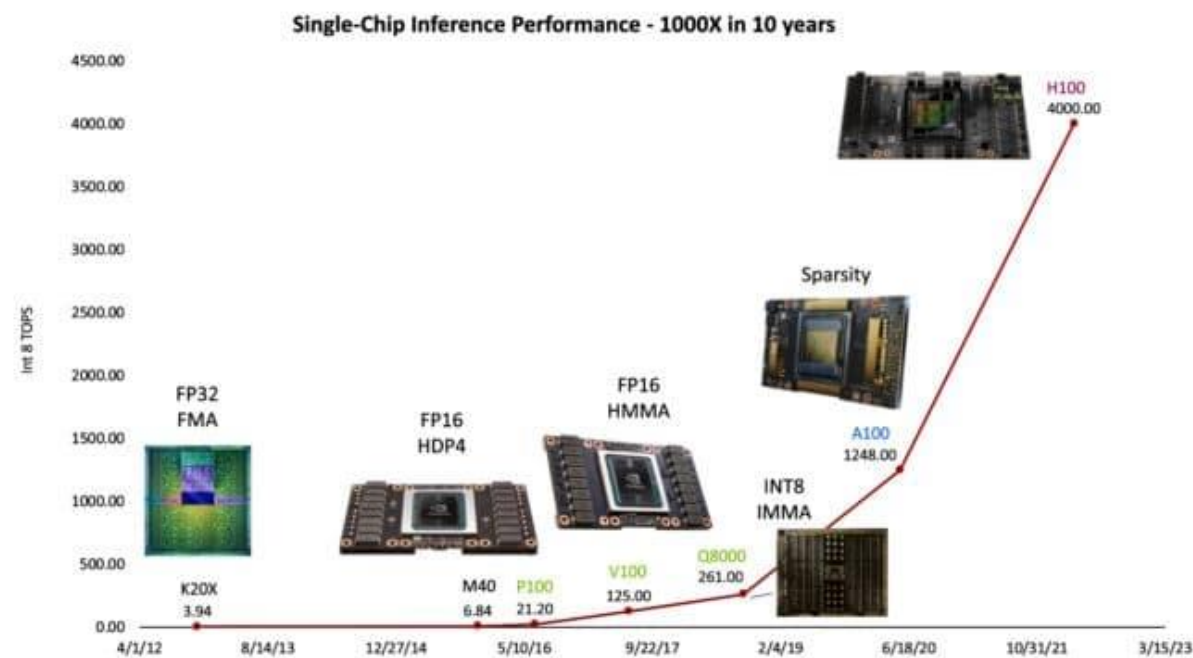
Application and platform -dependent

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

University of Pittsburgh  BYU BRIGHAM YOUNG UNIVERSITY  VIRGINIA TECH  UF UNIVERSITY of FLORIDA

# Motivation

## How BIG tech plans to feed AI's voracious appetite for power
### As data centers get more energy-hungry, the hyperscalers get more creative

Gains from

Number representation
FP32, FP16, Int8
(TF32, BF16)

Complex instructions
DP4, HMMA, IMMA

Process
28nm, 16nm, 7nm, 5nm



Source: Nvidia Blog

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

University of Pittsburgh

BYU
BRIGHAM YOUNG UNIVERSITY

VIRGINIA TECH

UF
UNIVERSITY of FLORIDA

# Tasks 1a & 2: Context for Jaccard Similarity

*You are here*

**Unpartitionable** Coarse-/Fine-grained work scheduling

Hybrid CPU+GPU

V1-25
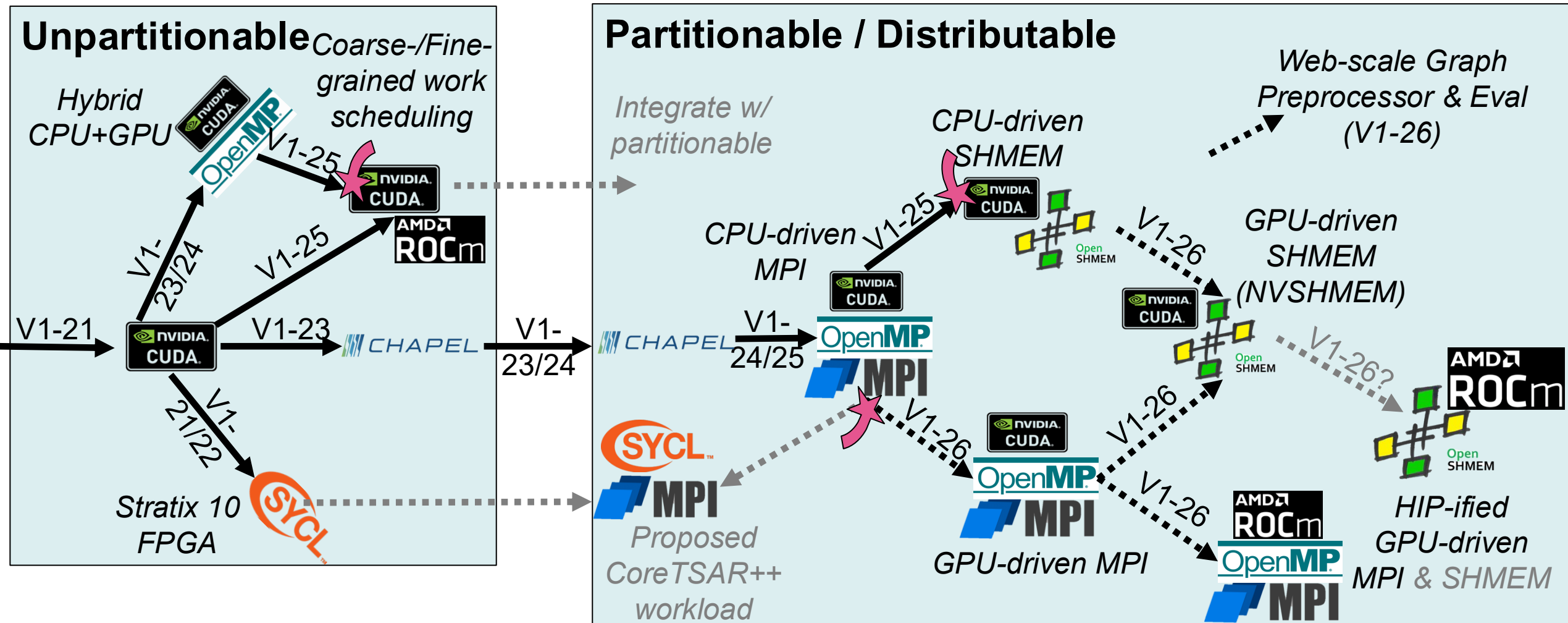
V1-23/24

V1-25

V1-21

V1-23

V1-23/24

V1-21/22

Stratix 10 FPGA

**Partitionable / Distributable**

*Integrate w/ partitionable*

Web-scale Graph Preprocessor & Eval (V1-26)

CPU-driven SHMEM

CPU-driven MPI

V1-25

GPU-driven SHMEM (NVSHMEM)

V1-26

V1-24/25

V1-26?

V1-26

V1-26

*Proposed CoreTSAR++ workload*

V1-26

V1-26

GPU-driven MPI

HIP-ified GPU-driven MPI & SHMEM

1/8/2026

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

University of Pittsburgh

BYU BRIGHAM YOUNG UNIVERSITY

VIRGINIA TECH

UF UNIVERSITY of FLORIDA

Optimal graph clustering is NP-hard → GraphChallenge MIT amazon web services *

**Enable fast and accurate graph clustering in large graphs by accelerating SBP**



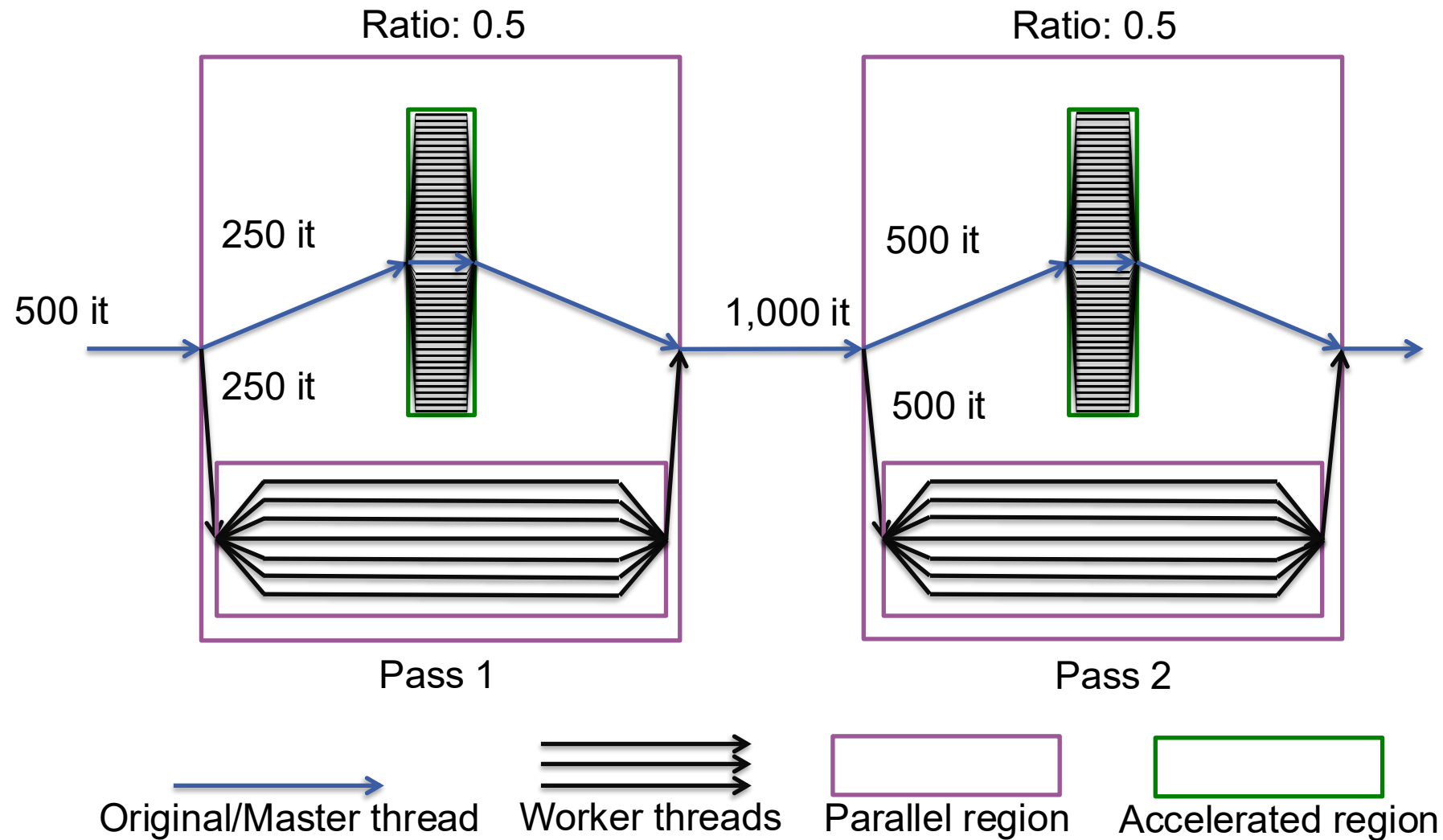Wanye, Gleyzer, Kao, and Feng.
MIT Graph Challenge Champion, 2023

**Port SBP to GPU**

Apply to applications like network intrusion detection

Speed

Quality

* E. Kao, V. Gadepally, M. Hurley, M. Jones, J. Kepner, S. Mohindra, P. Monticciolo, A. Reuther, S. Samsi, W. Wong, D. Staheli, S. Smith, "Streaming Graph Challenge: Stochastic Block Partition," *Proc. IEEE HPEC*,

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

# Task 1b: CoreTSAR: Core Task-Size Adapting Runtime

Ratio: 0.5

Ratio: 0.5

250 it

250 it

500 it

1,000 it

500 it

500 it

Pass 1

Pass 2

Original/Master thread    Worker threads    Parallel region    Accelerated region

# Task 1b: CoreTSAR: Core Task-Size Adapting Runtime

**Static vs. Dynamic Scheduler**

**Static**

250 it

500 it

250 it

Pass 1

1,000 it

500 it

500 it

Pass 2

Argument:
Ratio=0.5

**Dynamic**

250 it

500 it

250 it

Pass 1

1,000 it

900 it

100 it

Pass 2

Original/Master thread   Worker threads   Parallel region   Accelerated region

# Task 1b: CoreTSAR: Core Task-Size Adapting Runtime

**Dynamic vs. Split Scheduler**



Argument: Ratio=0.5

**Dynamic**

500 it — 250 it / 250 it — 1,000 it — 900 it / 100 it

Pass 1 — Pass 2

**Split**

500 it — 63 it / 113 it / 113 it / 113 it — 62 it / 12 it / 12 it / 12 it — 1,000 it — 113 it / 113 it / 113 it / 113 it — 12 it / 12 it / 12 it / 12 it

**Reschedule** — **Original/Master thread** — **Worker threads** — **Parallel region** — **Accelerated region**

**Split vs. Quick Scheduler**



Argument: Ratio=0.5

**Dynamic**

500 it  250 it  250 it  1,000 it  900 it  100 it

Pass 1  Pass 2

**Quick**

500 it  63 it  338 it  62 it  37 it  1,000 it  900 it  100 it

Reschedule  Original/Master thread  Worker threads  Parallel region  Accelerated region

# "Task 2b": Modernization of OpenDwarfs

**Goal:** *"Write once, run anywhere"*

- ■ Motivation:
  - ■ OpenDwarfs was a CHREC project to show how to map 13 *parallel computational idioms* to *GPUs* via OpenCL
    - • Part learning tool, part benchmark suite
    - • Eventually extended to Intel/Altera FPGAs
  - ■ Now *many more paths* to portable, heterogeneous computing
  - ■ To bridge *programming gap* between high-level, *library-driven* heterogeneity, need examples of how to write *novel* kernels
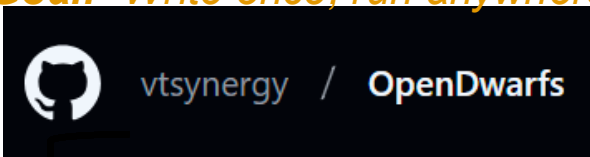- ■ Approach:
  - ■ Showcase idiomatic parallel codes using modern portable langs.
  - ■ Modernize for new classes of devices, and compute modalities
    - • {unified memory, PGAS, tensor cores, HBM, hybrid co-scheduling, DSPs, edge GPUs, … }
- ■ Milestones
  1. Update existing OpenCL Dwarfs for modern devices → characterize perf. shifts (0.25)
  2. Implement Dwarfs in new lang(s)., analyze perf./prod. (Π) vs. OpenCL (0.5 per lang.)
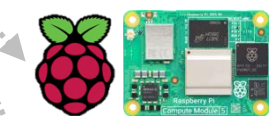  3. Design *partitionable/distributable* variants of existing dwarfs (1+)

vtsynergy / **OpenDwarfs**

OpenCL
OpenMP 4+
SYCL
Chapel
HIP
Vulkan Compute
`std::parallel`

?

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

**Tasks:** Optional (1-2 memberships)

27

University of Pittsburgh
BYU BRIGHAM YOUNG UNIVERSITY
VIRGINIA TECH
UF UNIVERSITY of FLORIDA