# IRIS: A Portable Runtime System Exploiting Multiple Heterogeneous Programming Systems

*Jungwon Kim*, Seyong Lee, Beau Johnston, and Jeffrey S. Vetter
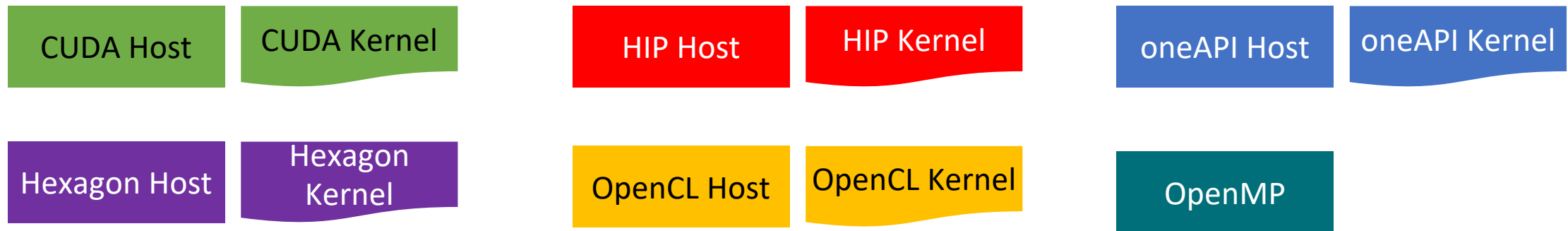
Oak Ridge National Laboratory

20 September 2021 @ IEEE HPEC '21

# No De Facto Standard for Heterogeneous Programming

- ORNL Experimental Computing Laboratory (ExCL) systems*

| Systems | Snapdragon | Jetson | Zynq | DGX | | | Oswald | | | Summit | Frontier | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **CPU** | ARM | ARM | ARM | I | I | I | I | I | I | IBM | AMD | |
| **GPU** | Qualcomm | NVIDIA | | NVIDIA | | | NV | | NV | NVIDIA | AMD | AMD |
| **FPGA** | | | Xilinx | | | | Intel | | Intel | | | |
| **DSP** | Qualcomm | | | | | | | | | | | |

CUDA Host    CUDA Kernel    HIP Host    HIP Kernel    oneAPI Host    oneAPI Kernel

Hexagon Host    Hexagon Kernel    OpenCL Host    OpenCL Kernel    OpenMP

* ORNL ExCL: https://excl.ornl.gov/

OAK RIDGE
National Laboratory

# We Need Portability in Heterogeneous Programming

- Not portable program across different HW configurations

| Systems | Snapdragon | Jetson | Zynq | DGX | Oswald | Summit | Frontier |
|---------|------------|--------|------|-----|--------|--------|----------|
| CPU | ARM | ARM | ARM | I  I  I | I  I  I | IBM | AMD |
| GPU | Qualcomm | NVIDIA | | NVIDIA | NV  NV | NVIDIA | AMD  AMD |
| FPGA | | | Xilinx | | Intel  Intel | | |
| DSP | Qualcomm | | | | | | |

**Snapdragon Host**

**OpenMP + OpenCL + Hexagon**

OpenMP Kernel

OpenCL Kernel

Hexagon Kernel

**Frontier Host**

**OpenMP + HIP**
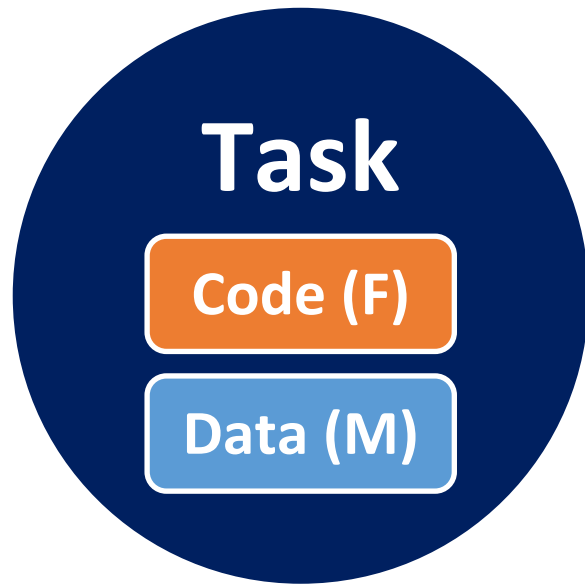
OpenMP Kernel

HIP Kernel

# Orchestrating Multiple Programming Systems
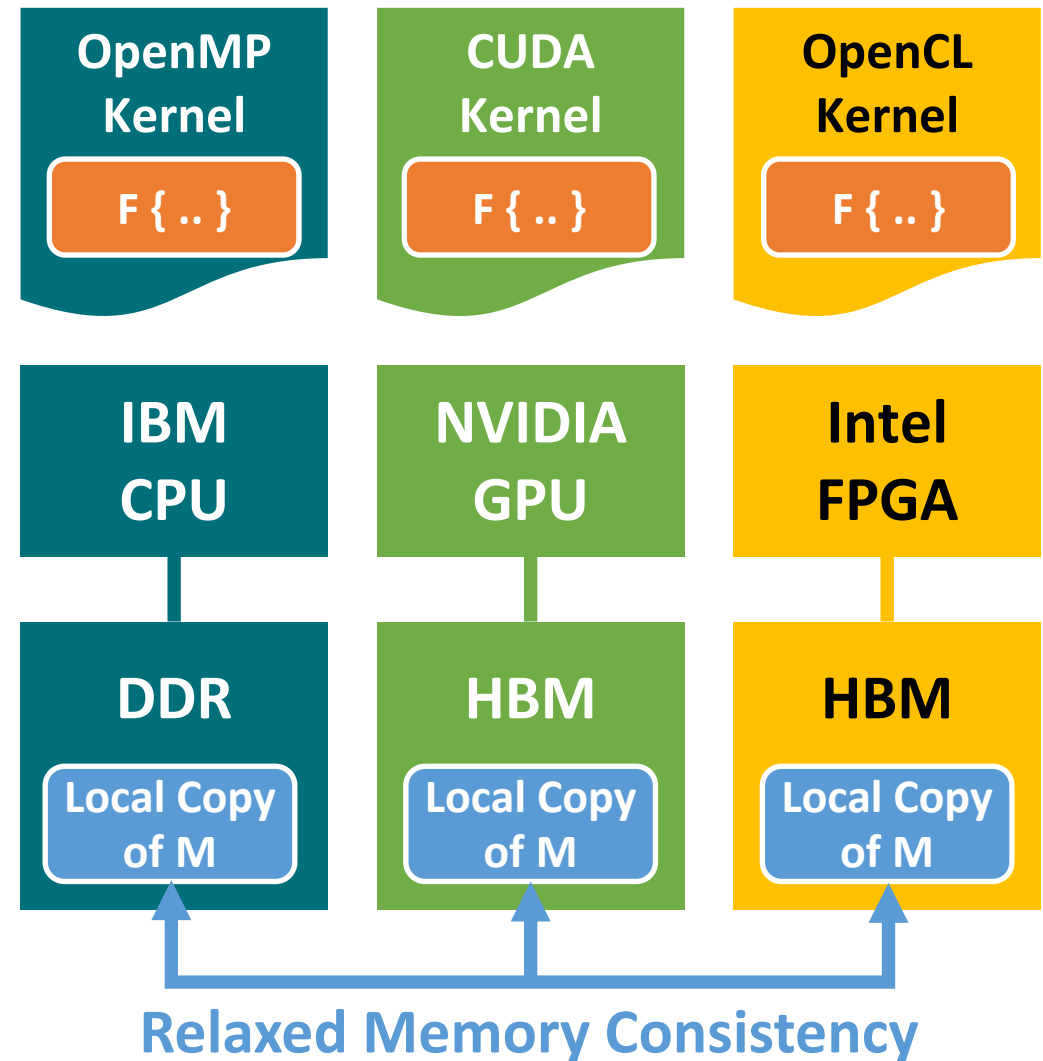
- ## The IRIS Architecture



- Compiler
  - High level application → IRIS unified host code + native kernels

- Dynamic Platform Loader
  - Automatically discover all available accelerators and their programming systems

- Task Scheduler
  - Task: memory copy + kernel launch
  - DAG-style tasks graph across multiple devices
  - Device Selection Policies

- Shared Virtual Device Memory
  - An Illusion of single logical device memory across all physical device memories
  - Multiple local copies on multiple device memories (relaxed consistency model)

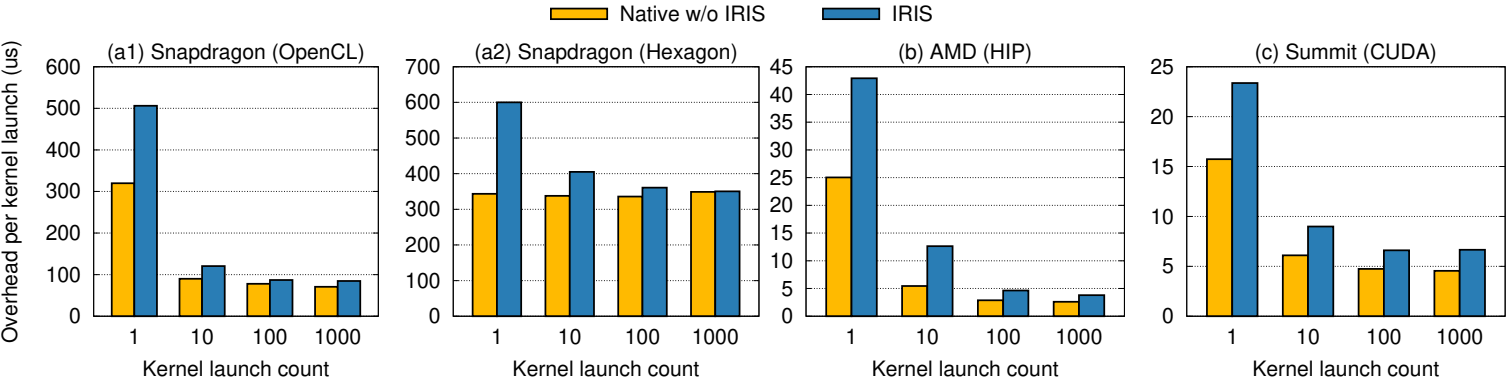# Multiple Native Kernels + SVDM = *Portable Tasks & Flexible Scheduling*

**Task**

Code (F)

Data (M)

A task can be scheduled and run on any device.

An IRIS application is portable across all heterogeneous systems.

OpenMP Kernel

F { .. }

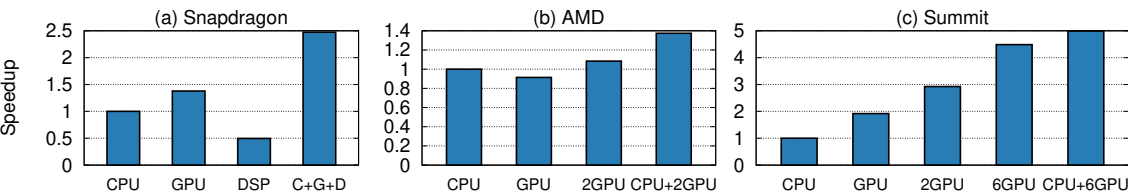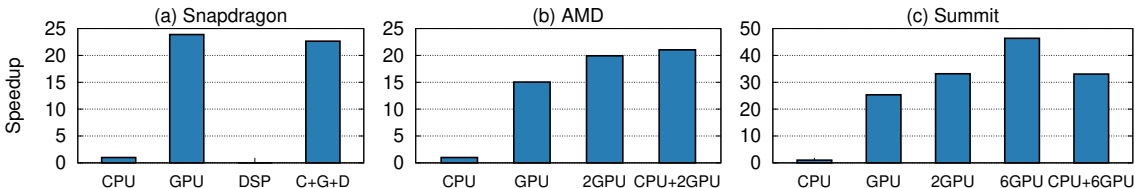CUDA Kernel

F { .. }

OpenCL Kernel

F { .. }

IBM CPU

NVIDIA GPU

Intel FPGA

DDR

Local Copy of M

HBM

Local Copy of M

HBM

Local Copy of M

**Relaxed Memory Consistency**

OAK RIDGE
National Laboratory

# Evaluation: Negligible Runtime Overhead



***Kernel Launch Overhead***

Legend: Native w/o IRIS (yellow), IRIS (blue)

(a1) Snapdragon (OpenCL)
(a2) Snapdragon (Hexagon)
(b) AMD (HIP)
(c) Summit (CUDA)

Overhead per kernel launch (us) vs Kernel launch count

***SAXPY***

(a) Snapdragon — CPU, GPU, DSP, C+G+D
(b) AMD — CPU, GPU, 2GPU, CPU+2GPU
(c) Summit — CPU, GPU, 2GPU, 6GPU, CPU+6GPU

***DGEMM***

(a) Snapdragon — CPU, GPU, DSP, C+G+D
(b) AMD — CPU, GPU, 2GPU, CPU+2GPU
(c) Summit — CPU, GPU, 2GPU, 6GPU, CPU+6GPU

***LULESH***

(a) Snapdragon — OpenCL, IRIS
(b) AMD — HIP, IRIS
(c) Summit — CUDA, IRIS

| Systems | Snapdragon | AMD | Summit |
|---|---|---|---|
| **CPU** | Qualcomm OpenMP | AMD OpenMP | IBM OpenMP |
| **GPU** | Qualcomm OpenCL | AMD HIP | NVIDIA CUDA |
| **DSP** | Qualcomm Hexagon | | |

OAK RIDGE
National Laboratory

# Recap

**Situation**    No de facto standard for heterogeneous programming

**Task**    Achieving portability in heterogeneous programming

**Activity**    We designed and implemented a new portable runtime system, ***IRIS***

- Orchestrating multiple programming systems (CUDA, Hexagon, HIP, Level Zero, OpenCL, OpenMP)
- Portable Tasks & Flexible Scheduling from Multiple Native Kernels + Shared Virtual Device Memory

**Result**    IRIS achieves portability, programmability, and performance

IRIS is freely available at

## https://iris-programming.com

**OAK RIDGE**
National Laboratory

# Acknowledgments

**OAK RIDGE**
National Laboratory