# Striving for SDR Performance Portability in the Era of Heterogeneous SoCs

Jeffrey S. Vetter

Seyong Lee

Mehmet Belviranli

Jungwon Kim

Richard Glassbrook

Abdel-Kareem Moadi

Seth Hitefield

FOSDEM
Brussels
2 Feb 2020

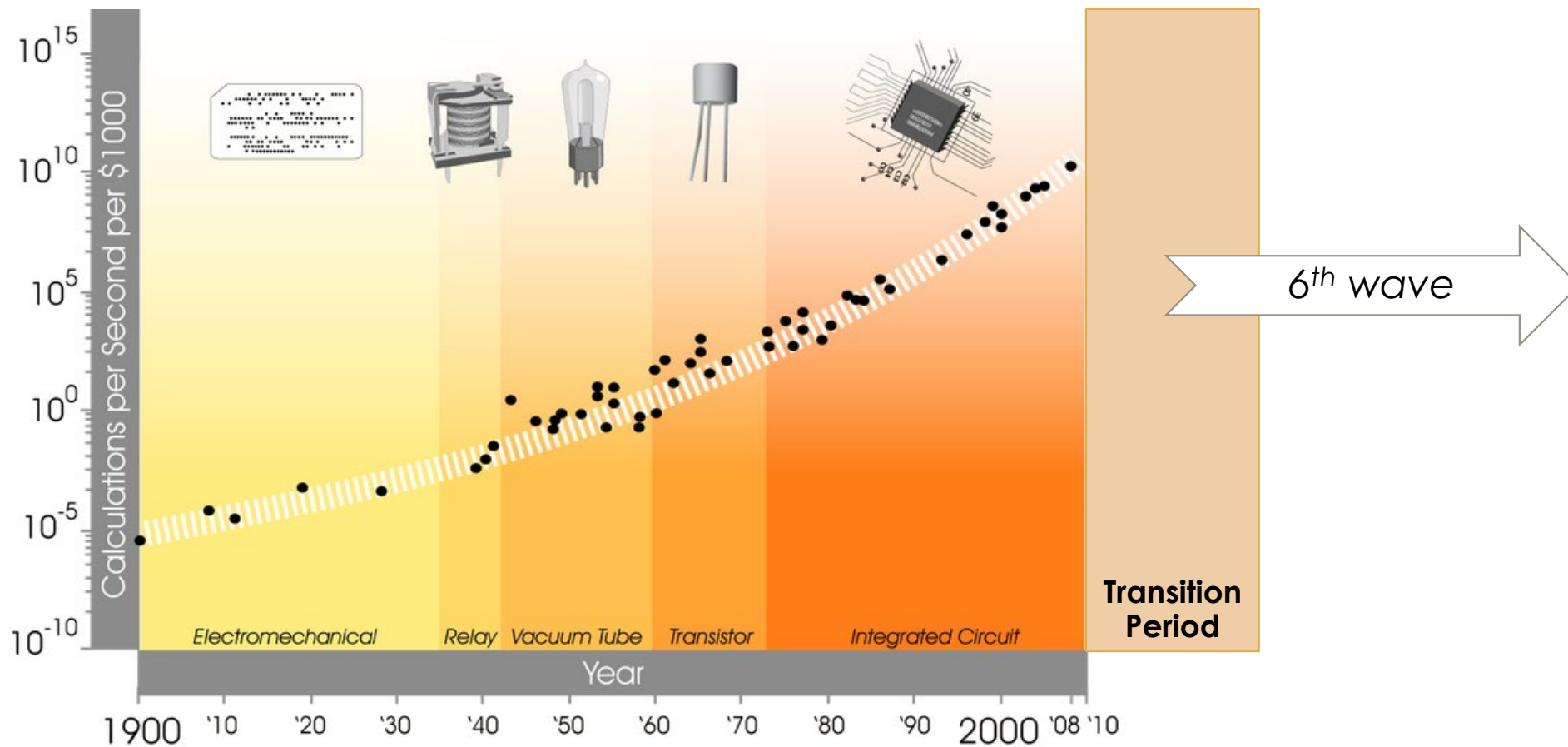ORNL is managed by UT-Battelle, LLC for the US Department of Energy

**U.S. DEPARTMENT OF ENERGY**

http://ft.ornl.gov     vetter@computer.org

# Highlights

- Architectural specialization

- Performance portability of applications and software

- DSSoC ORNL project investigating on performance portability of SDR
  - Understand applications and target architectures
  - Use open programming models (e.g., OpenMP, OpenACC, OpenCL)
  - Develop intelligent runtime systems

- Goal: scale applications from Qualcomm Snapdragon to DoE Summit Supercomputer with minimal programmer effort

OAK RIDGE
National Laboratory

# Sixth Wave of Computing



http://www.kurzweilai.net/exponential-growth-of-computing

OAK RIDGE
National Laboratory

37

# Predictions for Transition Period

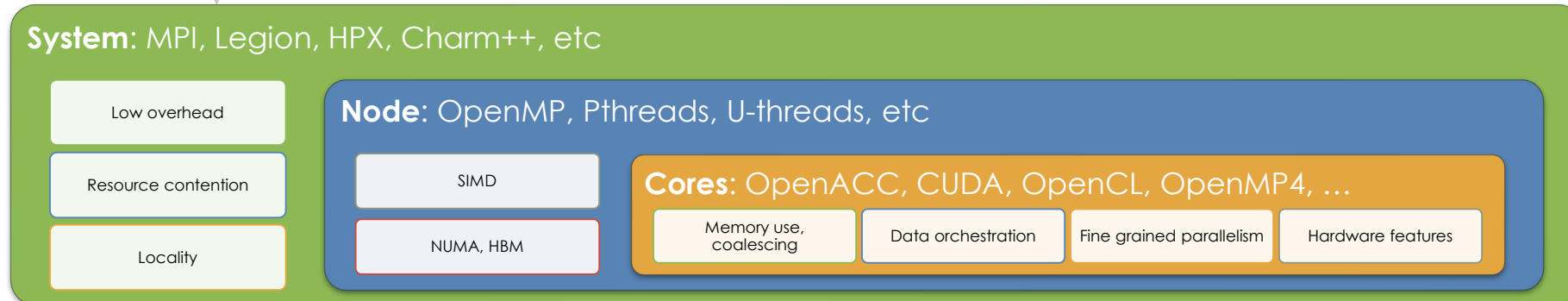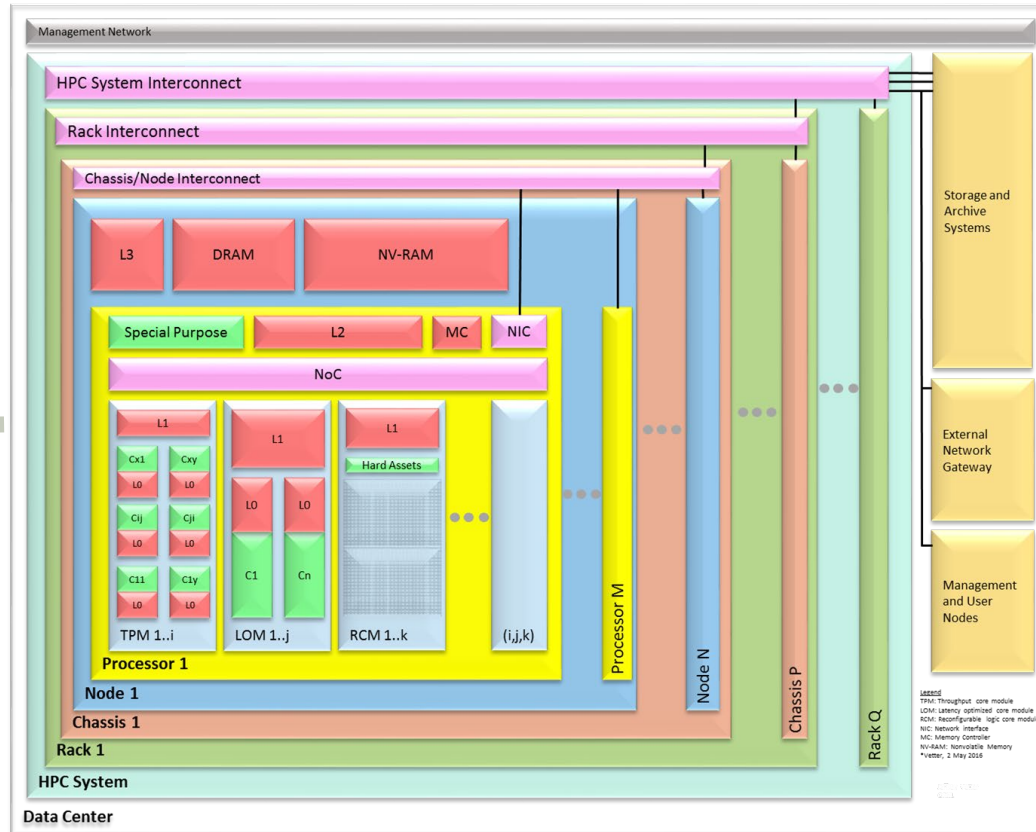| Optimize Software and Expose New Hierarchical Parallelism | Architectural Specialization and Integration | Emerging Technologies |
|---|---|---|
| • Redesign software to boost performance on upcoming architectures<br><br>• Exploit new levels of parallelism and efficient data movement | • Use CMOS more effectively for specific workloads<br><br>• Integrate components to boost performance and eliminate inefficiencies<br><br>• Workload specific memory+storage system design | • Investigate new computational paradigms<br>  • Quantum<br>  • Neuromorphic<br>  • Advanced Digital<br>  • Emerging Memory Devices |

OAK RIDGE
National Laboratory

# Complex architectures yield...



Complex Programming Models

# During this Sixth Wave transition, Complexity is our major challenge!

**Design**: How do we design future systems so that they are better than current systems on mission applications?

- Entirely possible that the new system will be slower than the old system!
- Expect 'disaster' procurements

**Programmability**: How do we design applications with some level of performance portability?

- Software lasts much longer than transient hardware platforms
- Adapt or die

OAK RIDGE
National Laboratory

# DARPA Domain-Specific System on a Chip (DSSoC) Program
# Getting the best out of specialization when we need programmability
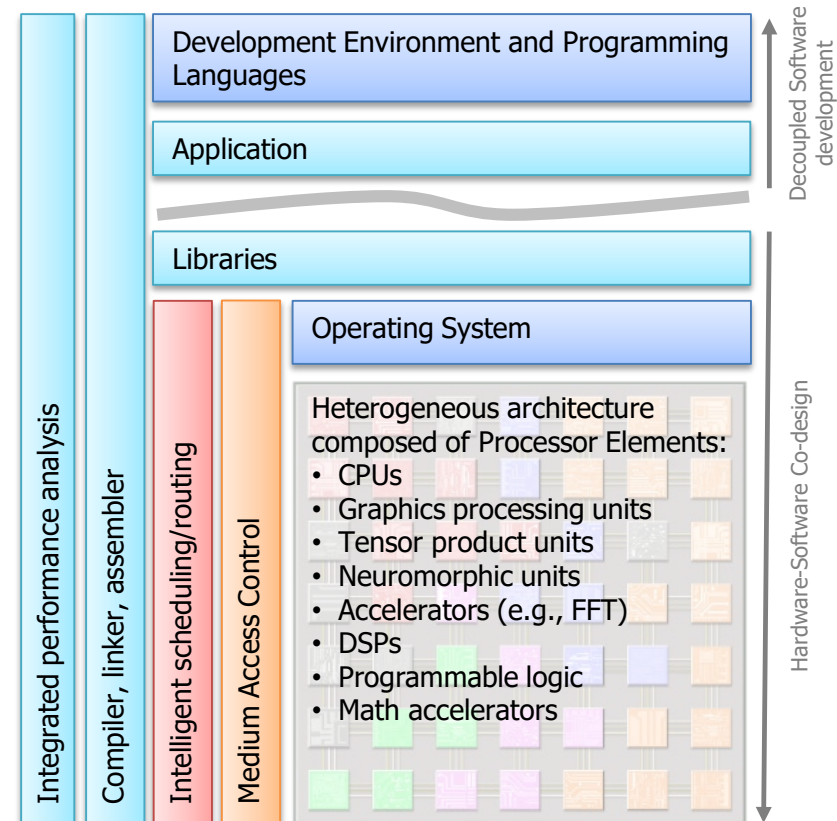
DARPA ERI DSSoC Program: Dr. Tom Rondeau

## Three Optimization Areas

1. Design time
2. Run time
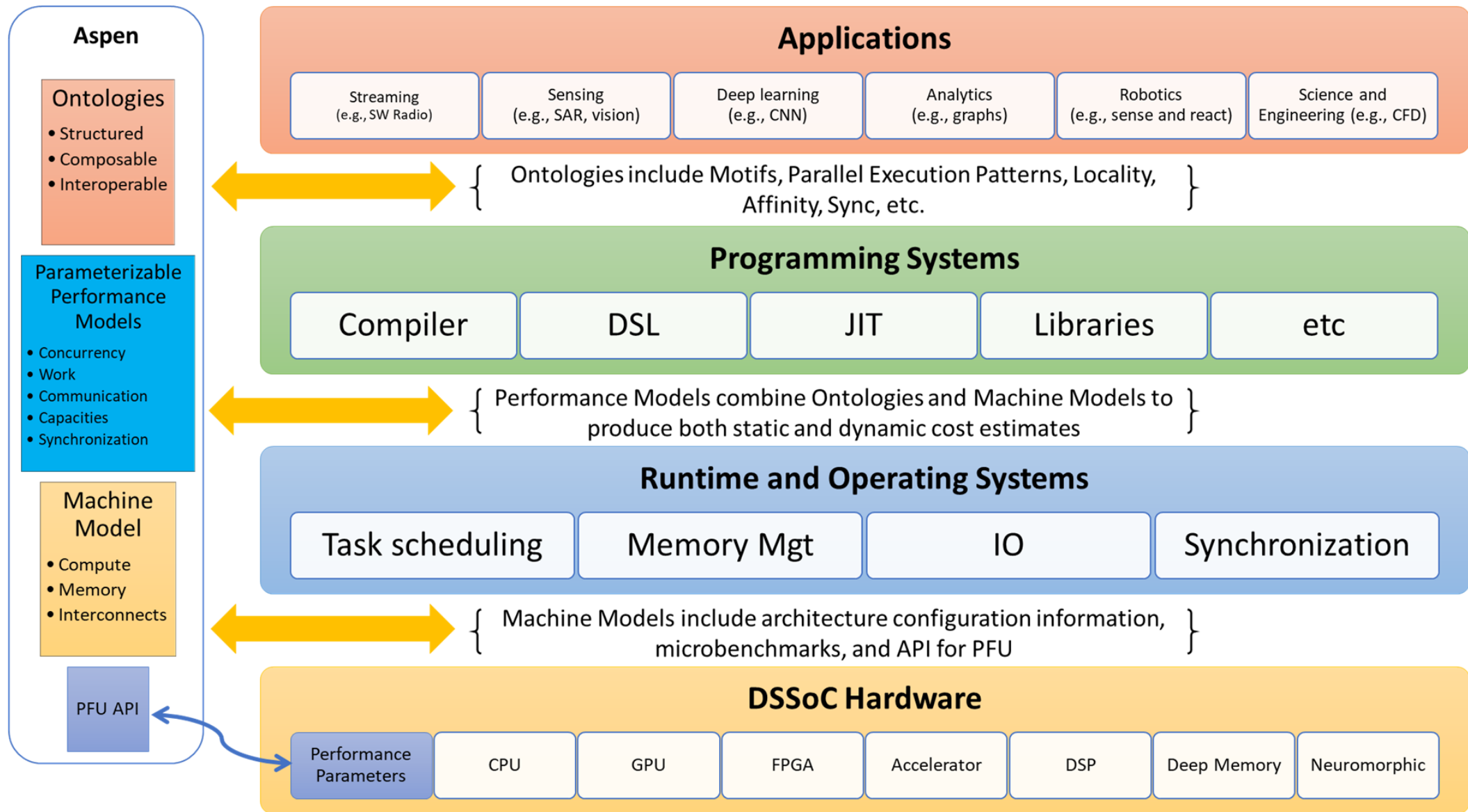3. Compile time

## Addressed via five program areas

1. Intelligent scheduling
2. Domain representations
3. Software
4. Medium access control (MAC)
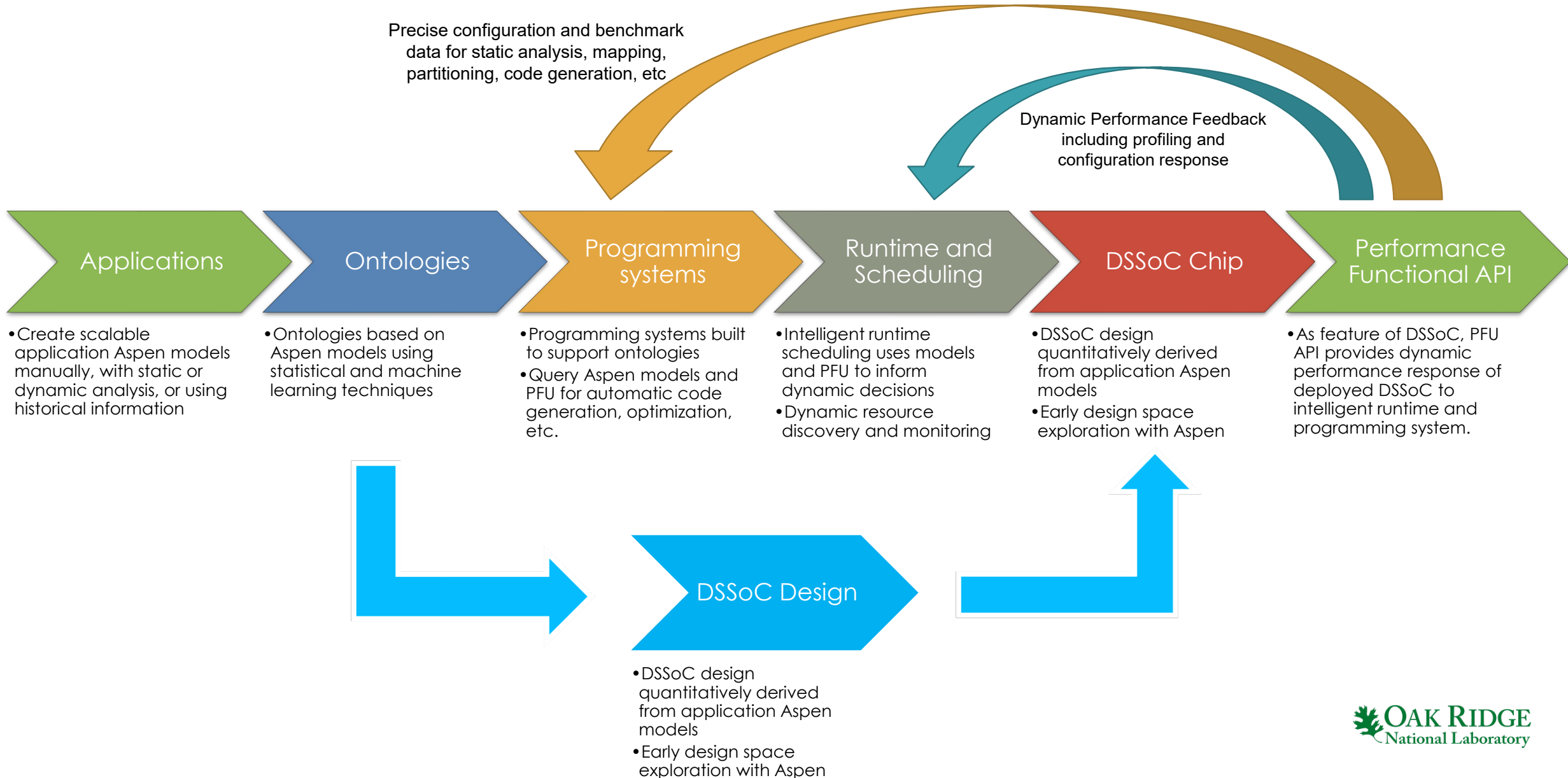5. Hardware integration

### DSSoC's Full-Stack Integration

Integrated performance analysis

Compiler, linker, assembler

Intelligent scheduling/routing

Medium Access Control

Development Environment and Programming Languages

Application

Libraries

Operating System

Heterogeneous architecture composed of Processor Elements:
- CPUs
- Graphics processing units
- Tensor product units
- Neuromorphic units
- Accelerators (e.g., FFT)
- DSPs
- Programmable logic
- Math accelerators

Decoupled Software development

Hardware-Software Co-design

*Looking at how Hardware/Software co-design is an enabler for efficient use of processing power*

# DSSoC ORNL Project Overview
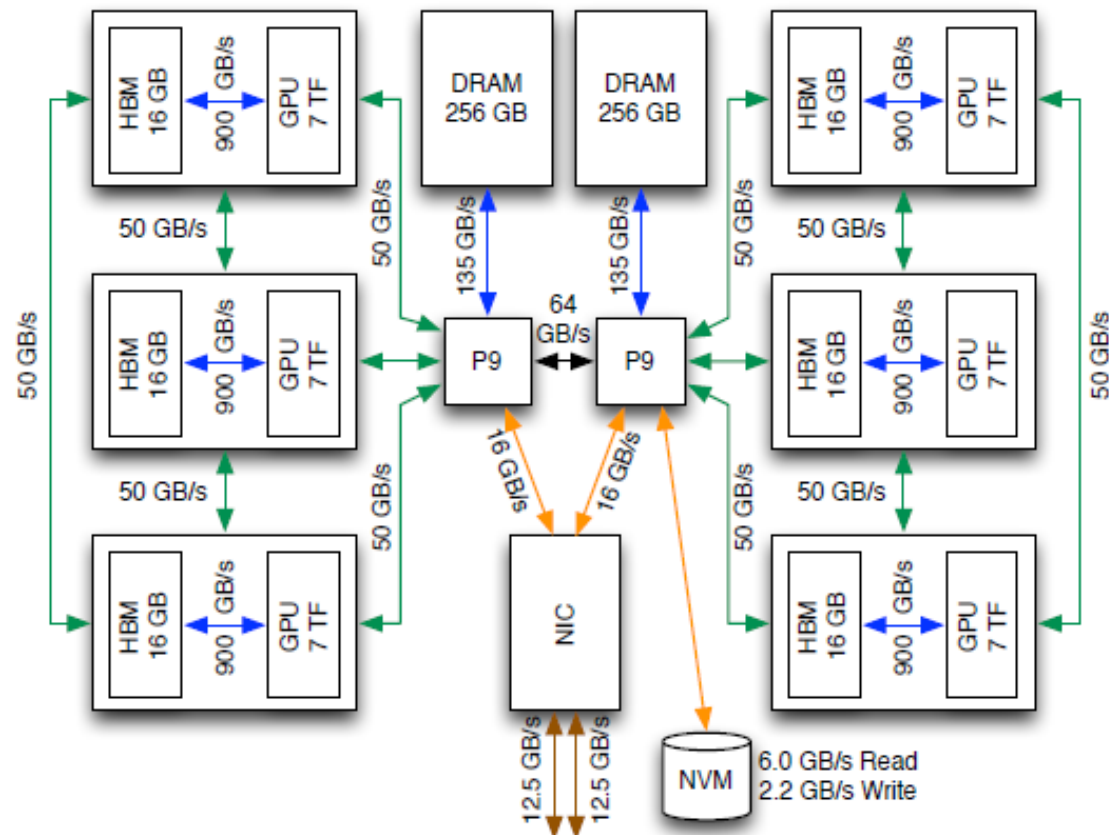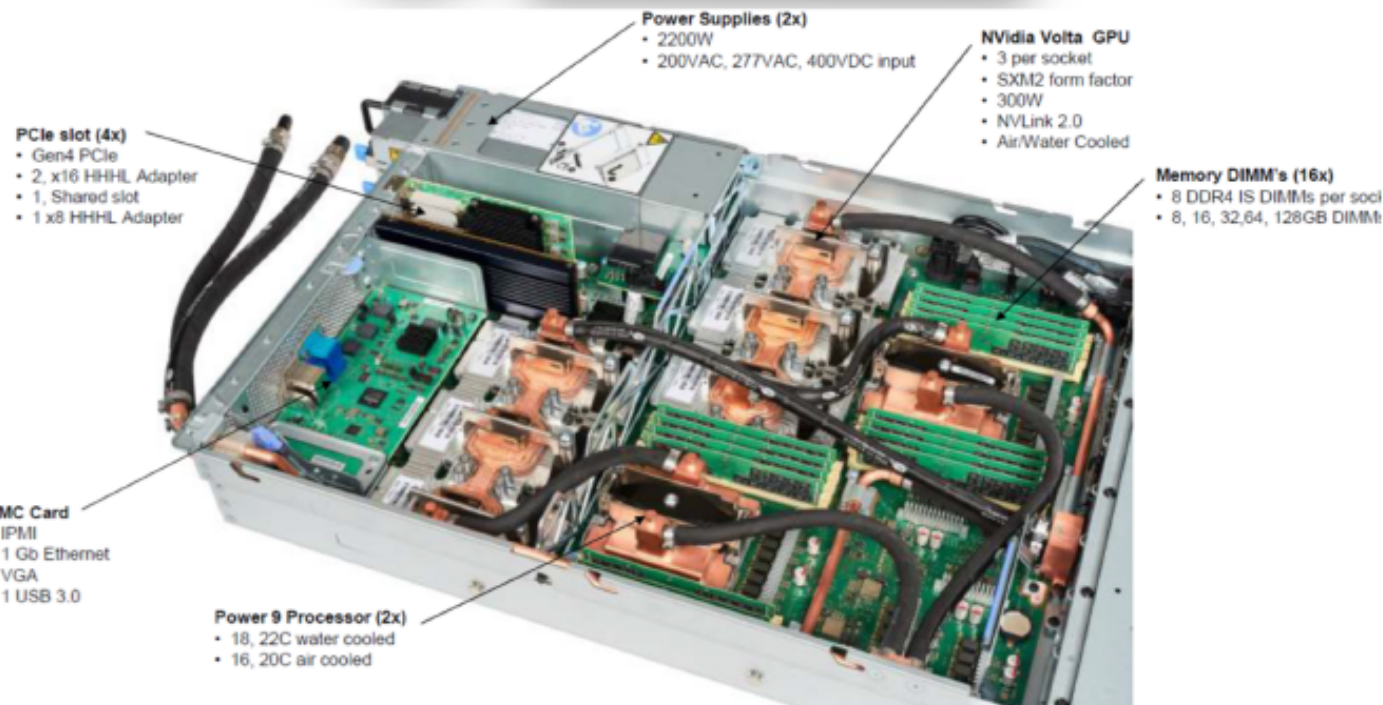
# Development Lifecycle



Precise configuration and benchmark data for static analysis, mapping, partitioning, code generation, etc

Dynamic Performance Feedback including profiling and configuration response

**Applications**
- Create scalable application Aspen models manually, with static or dynamic analysis, or using historical information

**Ontologies**
- Ontologies based on Aspen models using statistical and machine learning techniques

**Programming systems**
- Programming systems built to support ontologies
- Query Aspen models and PFU for automatic code generation, optimization, etc.

**Runtime and Scheduling**
- Intelligent runtime scheduling uses models and PFU to inform dynamic decisions
- Dynamic resource discovery and monitoring

**DSSoC Chip**
- DSSoC design quantitatively derived from application Aspen models
- Early design space exploration with Aspen

**Performance Functional API**
- As feature of DSSoC, PFU API provides dynamic performance response of deployed DSSoC to intelligent runtime and programming system.

**DSSoC Design**
- DSSoC design quantitatively derived from application Aspen models
- Early design space exploration with Aspen

OAK RIDGE
National Laboratory

# Architectures

# Summit Node Overview

| Application Performance | 200 PF |
|---|---|
| Number of Nodes | 4,608 |
| Node performance | 42 TF |
| Memory per Node | 512 GB DDR4 + 96 GB HBM2 |
| NV memory per Node | 1600 GB |
| Total System Memory | >10 PB DDR4 + HBM2 + Non-volatile |
| Processors | 2 IBM POWER9™ 9,216 CPUs<br>6 NVIDIA Volta™ 27,648 GPUs |
| File System | 250 PB, 2.5 TB/s, GPFS™ |
| Power Consumption | 13 MW |
| Interconnect | Mellanox EDR 100G InfiniBand |
| Operating System | Red Hat Enterprise Linux (RHEL) version 7.4 |



**Power Supplies (2x)**
- 2200W
- 200VAC, 277VAC, 400VDC input

**NVidia Volta GPU**
- 3 per socket
- SXM2 form factor
- 300W
- NVLink 2.0
- Air/Water Cooled

**PCIe slot (4x)**
- Gen4 PCIe
- 2, x16 HHHL Adapter
- 1, Shared slot
- 1 x8 HHHL Adapter

**Memory DIMM's (16x)**
- 8 DDR4 IS DIMMs per soci
- 8, 16, 32,64, 128GB DIMM

**BMC Card**
- IPMI
- 1 Gb Ethernet
- VGA
- 1 USB 3.0

**Power 9 Processor (2x)**
- 18, 22C water cooled
- 16, 20C air cooled



| TF | 42 TF (6x7 TF) | ← → | HBM/DRAM Bus (aggregate B/W) |
| HBM | 96 GB (6x16 GB) | ← → | NVLINK |
| DRAM | 512 GB (2x16x16 GB) | ← → | X-Bus (SMP) |
| NET | 25 GB/s (2x12.5 GB/ | | |
| MMsg/s | 83 | | |

HBM & DRAM speeds are aggregate (Read+Write).
All other speeds (X-Bus, NVLink, PCIe, IB) are bi-directional.

OAK RIDGE National Laboratory | 75 YEARS

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

# Intel Stratix 10 FPGA

**Experimental Computing Lab (ExCL) managed by the ORNL Future Technologies Group**

- Intel Stratix 10 FPGA and four banks of DDR4 external memory
  - Board configuration: Nallatech 520 Network Acceleration Card

- Up to 10 TFLOPS of peak single precision performance

- 25MBytes of L1 cache @ up to 94 TBytes/s peak bandwidth

- 2X Core performance gains over Arria® 10

- Quartus and OpenCL software   (Intel SDK v18.1) for using FPGA

- Provide researcher access to advanced FPGA/SOC environment

For more information or to apply for an account, visit https://excl.ornl.gov/

Mar 2019

**OAK RIDGE**
National Laboratory

# NVIDIA Jetson AGX Xavier SoC

**Experimental Computing Lab (ExCL) managed by the ORNL Future Technologies Group**

- NVIDIA Jetson AGX Xavier:

- High-performance system on a chip for autonomous machines

- Heterogeneous SoC contains:
  - Eight-core 64-bit ARMv8.2 CPU cluster (Carmel)
  - 1.4 CUDA TFLOPS (FP32) GPU with additional inference optimizations (Volta)
  - 11.4 DL TOPS (INT8) Deep learning accelerator (NVDLA)
  - 1.7 CV TOPS (INT8) 7-slot VLIW dual-processor Vision accelerator (PVA)
  - A set of multimedia accelerators (stereo, LDC, optical flow)

- Provides researchers access to advanced high-performance SOC environment

For more information or to apply for an account, visit https://excl.ornl.gov/

Mar 2019

# Qualcomm 855 SoC (SM8510P) Snapdragon™

**Experimental Computing Lab (ExCL) managed by the ORNL Future Technologies Group**



Adreno 640 | Snapdragon X24 modem
Hexagon 690 | Wi-Fi/BT/Location
Spectra 360
Kyro 485 | Security
© Qualcomm Inc.



Snapdragon 855 Mobile Platform
7nm TSMC
Wi-Fi
855
5G
RF | RFFE
Audio
Qualcomm snapdragon
Finger Print
PMIC | Quick Charge
Not to scale; chipset enlarged for visibility.
© Qualcomm Inc.

**Qualcomm Development Board connected to (mcmurdo) HPZ820**



## Kyro 485 (8-ARM Prime+BigLittle Cores)

| @2.84G | @ 2.42G | | | @ 1.8G | |
|---|---|---|---|---|---|
| Prime Core A76 | A76 | A76 | A76 | A55 128 KB | A55 128 KB |
| | | | | A55 128 KB | A55 128 KB |
| 512 KB | 256 KB | 256 KB | 256 KB | | |
| DSU | | 2048KB | | | |

## Connectivity (5G)

- Snapdragon X24 LTE (855 built-in) modem LTE Category 20
- Snapdragon X50 5G (external) modem (for 5G devices)
- Qualcomm Wi-Fi 6-ready mobile platform: (802.11ax-ready, 802.11ac Wave 2, 802.11ay, 802.11ad)
- Qualcomm 60 GHz Wi-Fi mobile platform: (802.11ay, 802.11ad)
- Bluetooth Version: 5.0
- Bluetooth Speed: 2 Mbps
- High accuracy location with dual-frequency GNSS.

## Hexagon 690 (DSP + AI)

- Quad threaded Scalar Core
- DSP + 4 Hexagon Vector Xccelerators
- New Tensor Xccelerator for AI
- Apps: AI, Voice Assistance, AV codecs

## Adreno 640

- Vulkan, OpenCL, OpenGL ES 3.1
- Apps: HDR10+, HEVC, Dolby, etc
- Enables 8k-360° VR video playback
- 20% faster compared to Adreno 630

## Spectra 360 ISP

- New dedicated Image Signal Processor (ISP)
- Dual 14-bit CV-ISPs; 48MP @ 30fps single camera
- Hardware CV for object detection, tracking, streo depth process
- 6DoF XR Body tracking, H265, 4K60 HDR video capture, etc.

- Connected Qualcomm board to HPZ820 through USB
- Development Environment: Android SDK/NDK
- Login to mcmurdo machine
    - $ ssh –Y mcmurdo
- Setup Android platform tools and development environment
    - $ source /home/nqx/setup_android.source
- Run Hello-world on ARM cores
    - $ git clone https://code.ornl.gov/nqx/helloworld-android
    - $ make compile push run
- Run OpenCL example on GPU
    - $ git clone https://code.ornl.gov/nqx/opencl-img-processing
    - Run Sobel edge detection
        - $ make compile push run fetch
- Login to Qualcomm development board shell
    - $ adb shell
    - $ cd /data/local/tmp

For more information or to apply for an account, visit https://excl.ornl.gov/

OAK RIDGE National Laboratory

Created by Narasinga Rao Miniskar, Steve Moulton

Xavier SoC #1

GR IEEE-802.11 Transmit (TX)

Xavier SoC #2

IEEE-802.11 Receive (RX)

Antenna

UDP

OpenCV

OpenCL

- Signal processing: An open-source implementation of IEEE-802.11 WIFI a/b/g with GR OOT modules.
- Input / Output file support via Socket PDU (UDP server) blocks
- Image/Video transcoding with OpenCL/OpenCV

141

- GR-Tools

  - First tools are released

    - Block-level Ontologies [ontologyAnalysis]

      - Following properties are extracted from a batch of block definition files: Descriptions and IDs, source and sink ports (whether input/output is scalar, vector or multi-port), allowed data types, and additional algorithm-specific parameters

    - Flowgraph Characterization [workflowAnalysis]

      - Characterization of GR workloads at the flowgraph level.

      - Scripts automatically run for for 30 seconds and reports a breakdown of high-level library module calls

    - Design-space Exploration [designSpaceCL]

      - Script to run 13 blocks included in gr-clenabled
        - Both on a GPU and on a single CPU core
        - By using input sizes varying between 24 and 227 elements.

  - Two prototype tools have been added recently

    - cgran-scraper

    - GRC-analyzer



libgnuradio CPU-time Breakdown

- libgnuradio-analog
- libgnuradio-blocks
- libgnuradio-channels
- libgnuradio-digital
- libgnuradio-dtv

4% 13% 6% 3% 4% 1% 28% 10% 1% 22% 8% 0%

https://github.com/cosmic-sdr

- **Preliminary SDR Application Profiling:**

  - **Created fully automated GRC profiling toolkit**

  - Ran each of the 89 flowgraph for 30 seconds

  - Profiled with performance counters

  - Major overheads:

    - Python glue code (libpython), O/S threading & profiling (kernel.kallsysms, libpthread), libc, ld, Qt

  - Runtime overhead:

    - Will require significant consideration when run on SoC

    - Cannot be executed in parallel

    - Hardware assisted scheduling is essential

| Library | Percentage |
|---|---|
| [kernel.kallsyms] | 27.8547 |
| libpython | 18.6281 |
| libgnuradio | 11.7548 |
| libc | 7.7503 |
| ld | 3.8839 |
| libvolk | 3.7963 |
| libperl | 3.7837 |
| [unknown] | 3.6465 |
| libQt5 | 2.9866 |
| libpthread | 2.1449 |

libgnuradio CPU-time Breakdown



- libgnuradio-analog
- libgnuradio-blocks
- libgnuradio-channels
- libgnuradio-digital
- libgnuradio-dtv
- libgnuradio-fec
- libgnuradio-fft

## Block proximity analysis

- Creates a graph:
  - <u>Nodes</u>: Unique block types
  - <u>Edges</u>: Blocks used in the same GRC file.
  - Every co-occurrence increases edge weight by 1.
- This example was run
  - With `--mode proximityGraph`
  - On randomly selected sub-set of GRC files

# Programming Solution for DSSoC

# New OpenACC GR Block Mapping Strategy for Heterogeneous Architectures

```
//Constructor
accLog_impl::accLog_impl(...
  int contextType, int deviceId, int copy_in, int copy_out)
  : gr::sync_block("accLog", ...), ...  GRACCBase(contextType, deviceId) {
    accLog_init(deviceType, deviceId, threadID);
    ...
}

//Reference CPU implementation
int accLog_impl::testCPU(...) {
    ...
    for (int i=0;i<noi;i++) {
        out[i] = n_val * log10(in1[i]) + k_val;
    }
    ...
}

//OpenACC implementation
int accLog_impl::testOpenACC(...) {
    ...
    if( acc_init_done == 0 ) {
        gracc_pcopyin(...); //Create and copy input data to device memory.
        gracc_pcreate(...); //Create device buffer for output data.
        acc_init_done = 1;
    } else if( gracc_copy_in == 1 ) {
        gracc_update_device(...); //Copy input data to device memory.
    }
    accLog_kernel(...); //Execute an OpenACC kernel.
    if( gracc_copy_out == 1 ) {
        gracc_update_self(...); //Copy output data to host memory.
    }
    ...
}

int accLog_impl::work(...) {
    ...
    if( contextType == ACCTYPE_CPU ) {
        retVal = testCPU(...); //Execute reference CPU version.
    } else {
        retVal = testOpenACC(...); //Execute OpenACC version.
    }
    ...
}
```

## Constructor
- OpenACC GR block class inherits GRACCBase class as a base class.
- GRACCBase constructor assigns a unique thread ID per OpenACC GR block instantiation, which is internally used for thread safety.
- OpenACC backend runtime is also initialized.

## Reference CPU Implementation
- Contains the same code as that in the original GR block, which may have already been vectorized using Volk library.

## OpenACC Implementation
- Contains the OpenACC version of the reference CPU implementation.
- Performs the following tasks:
  - Copy input data to device memory.
  - Execute the OpenACC kernel.
  - Copy output data back to host memory.
- OpenARC will translate the OpenACC kernel to multiple different output programming models (e.g., CUDA, OpenCL, OpenMP, HIP, etc.)

## Main Entry Function
- Main entry function executed whenever GR scheduler invokes the OpenACC GR block.
- The GR block argument, contextType decides which to execute between the reference CPU version and OpenACC version.
  - OpenACC backend runtime may choose CPU as an offloading target (e.g., offloading OpenMP3 kernel to CPU).

Input OpenACC code

Output host code

Output CUDA kernel code

- In the basic memory management scheme, each invocation of an OpenACC GR block performs the following three tasks:
    1) Copy input data to device memory.
    2) Run a kernel on device.
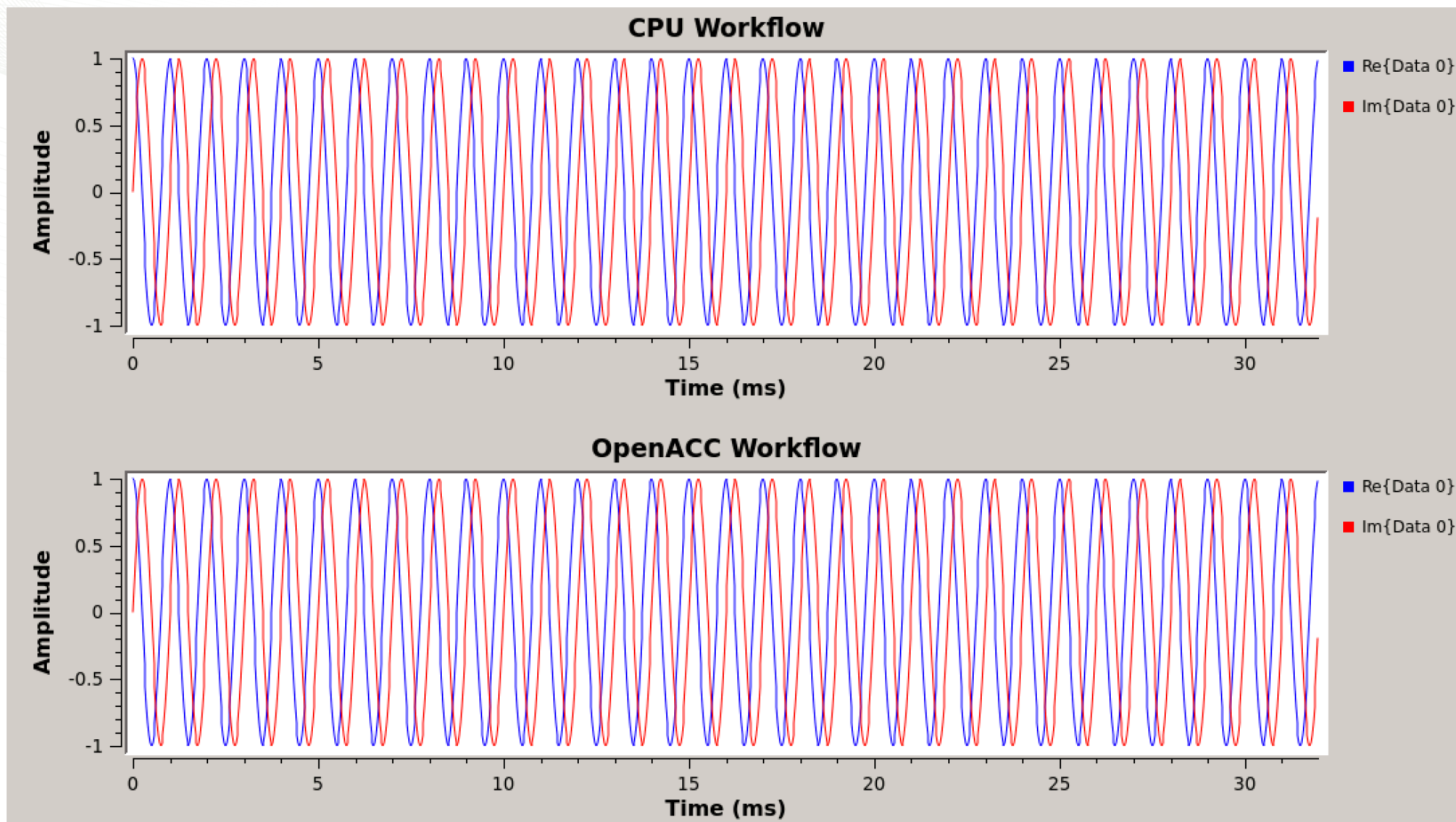    3) Copy output data back to host memory.

# Optimized Memory Management for OpenACC-Enabled GR Workflow



- In the optimized memory management scheme, some blocks can bypass unnecessary memory transfers between host and device and directly communicate each other using device memory if both producer and consumer blocks are running on the same device.
- Notice that device buffer needs extra padding to handle the overwriting feature in the host circular buffer.

# SDR Workflow Profiling Using a Built-in GR Performance Monitoring Tool

- CPU versions of OpenACC blocks are algorithmically equivalent to those in the original GR blocks.



Some OpenACC blocks (B, D) use a simple register caching optimization, which causes them to perform better than the original GR blocks.

- OpenACC blocks are automatically translated to OpenMP3 versions and run on Xavier CPU.



OpenACC Blocks on Xavier CPU via OpenMP

Original GR Blocks on Xavier CPU

Some of original GR blocks (A, C) were already vectorized with Volk library.

Some of original GR blocks (B, C) performed better than OpenACC blocks (B, C).

- OpenACC blocks are automatically translated to CUDA versions and run on Xavier GPU.
- Each invocation of an OpenACC block executes three tasks: 1) copy input data to device memory, 2) run a kernel on device, and 3) copy output data back to host memory



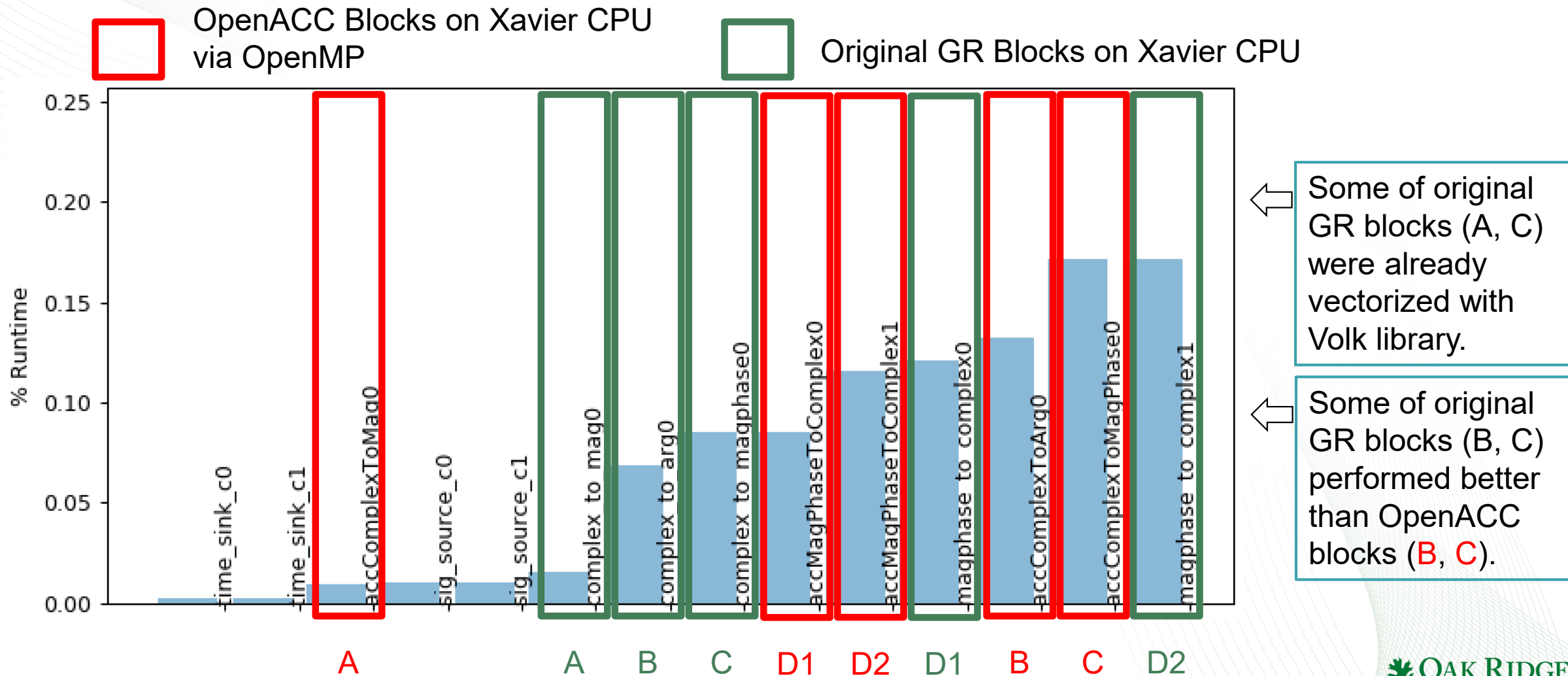OpenACC Blocks on Xavier GPU            Original GR Blocks on Xavier CPU

Due to extra memory transfer overheads, most OpenACC blocks perform worse than original GR blocks, except for the OpenACC block D1 and D2.

- OpenACC blocks are automatically translated to CUDA versions and run on Xavier GPU.
- Optimized OpenACC blocks bypass memory transfers between host and device and directly communicate each other using device memory if both producer and consumer blocks are running on the same device.

Opt. OpenACC Blocks on Xavier GPU

Original GR Blocks on Xavier CPU



Most of the OpenACC blocks perform better than original GR blocks, except for the block A; the original GR block A is vectorized with Volk library, which performs better than the OpenACC block A.

This example offloads more OpenACC blocks to Xavier GPU than the previous example.

OpenACC-enabled workflow using **gr-openacc** blocks

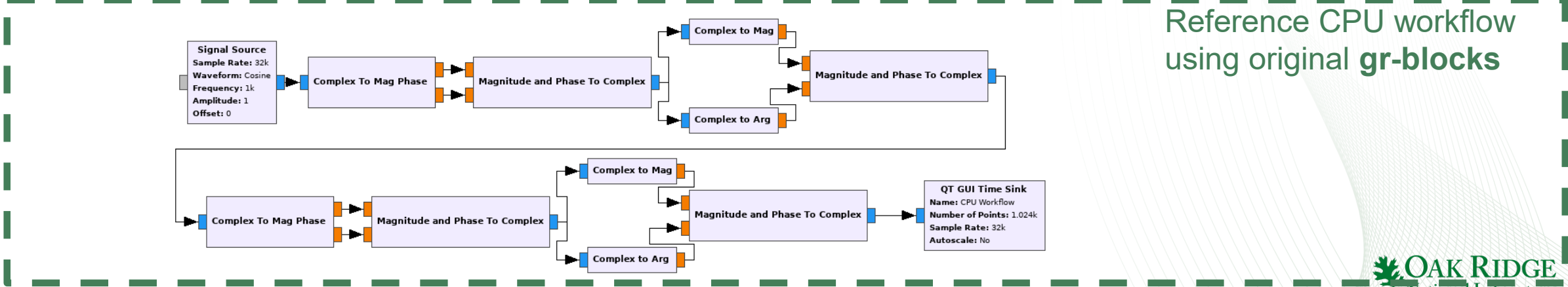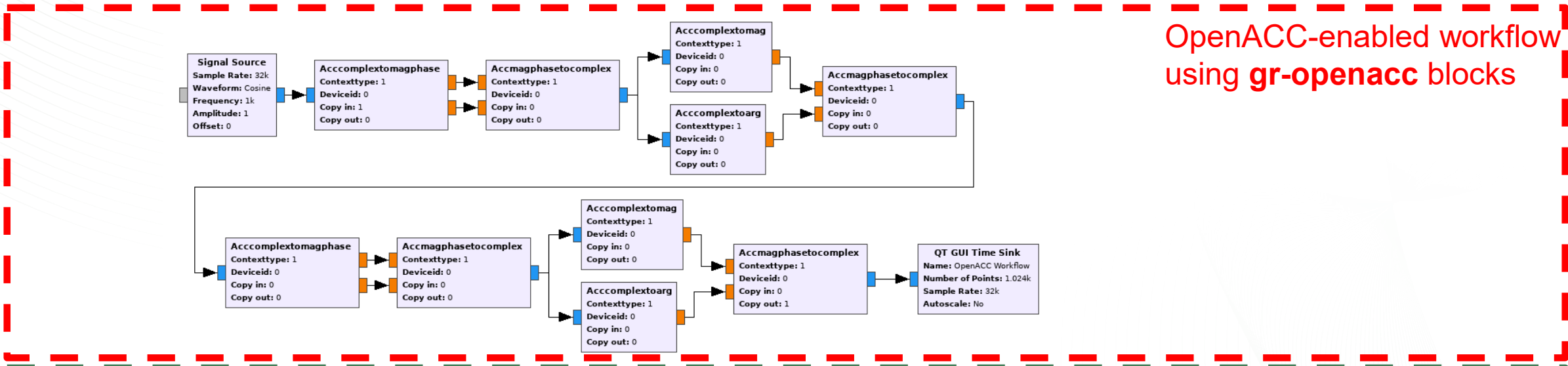Reference CPU workflow using original **gr-blocks**

- OpenACC blocks are automatically translated to CUDA versions and run on Xavier GPU.
- Optimized OpenACC blocks bypass memory transfers between host and device and directly communicate each other using device memory if both producer and consumer blocks are running on the same device.



Opt. OpenACC Blocks on Xavier GPU
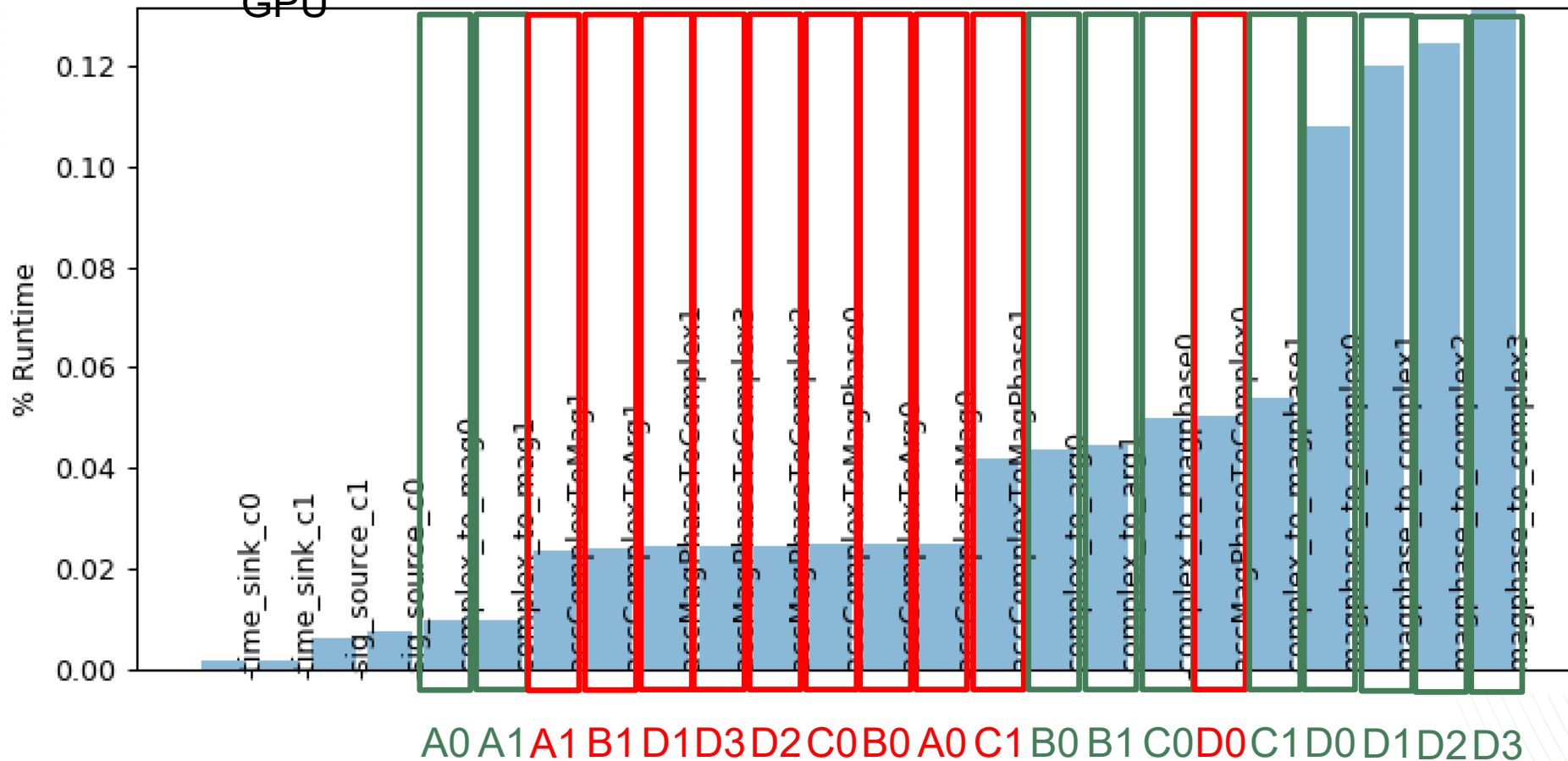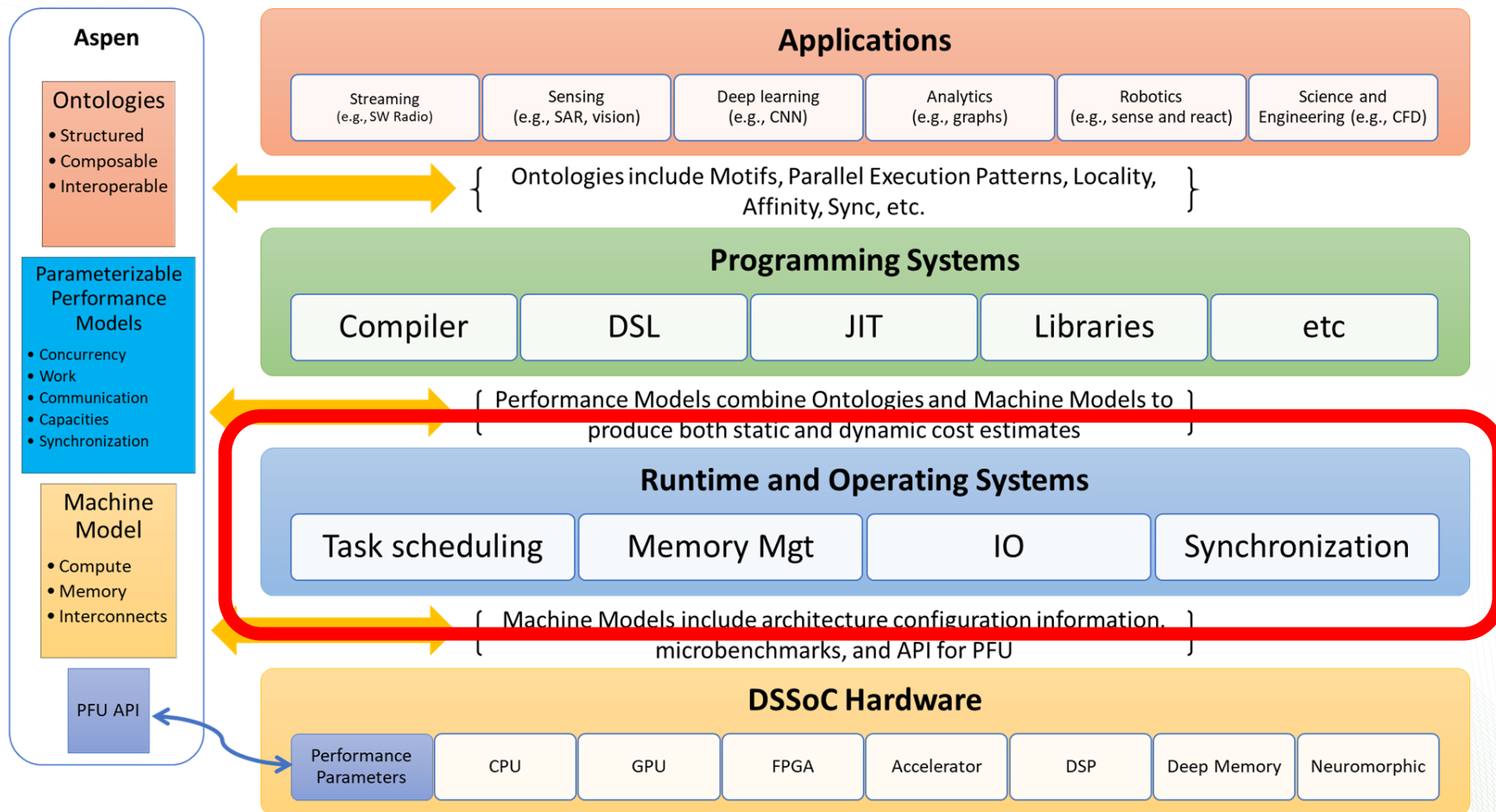
Original GR Blocks on Xavier CPU

This example shows similar performance behaviors as the previous example.

# Programming Systems Update Summary and Next Steps

- Updated the programming system to use our new heterogeneous runtime system, called IRIS, as the common backend runtime.

    - IRIS allows intermixing of multiple different output programming models (e.g., OpenMP3, OpenMP4, OpenACC, CUDA, HIP, etc.) and runs them on heterogeneous devices concurrently.

- Developed a host-device memory transfer optimization scheme, which allows OpenACC GR blocks to bypass memory transfers between host and device and directly communicate each other if both producer and consumer blocks are running on the same device.

- Performed preliminary evaluation of the new programming system by creating synthetic SDR workflow using the OpenACC GR blocks.

- Next Steps

    - Port more complex GR blocks to OpenACC and evaluate more complex SDR workflow.

    - Continue to improve and fix bugs in the programming system.

# IRIS: An Intelligent Runtime System for Extremely Heterogeneous Architectures

- Provide programmers a unified programming environment to write portable code across heterogeneous architectures (and preferred programming systems)

- Orchestrate diverse programming systems (OpenCL, CUDA, HIP, OpenMP for CPU) in a single application
  - OpenCL
    - NVIDIA GPU, AMD GPU, ARM GPU, Qualcomm GPU, Intel CPU, Intel Xeon Phi, Intel FPGA, Xilinx FPGA
  - CUDA
    - NVIDIA GPU
  - HIP
    - AMD GPU
  - OpenMP for CPU
    - Intel CPU, AMD CPU, PowerPC CPU, ARM CPU, Qualcomm CPU



https://github.com/swiftcurrent2018

**OAK RIDGE** National Laboratory

# The IRIS Architecture

- Platform Model
  - A single-node system equipped with host CPUs and multiple compute devices (GPUs, FPGAs, Xeon Phis, and multicore CPUs)

- Memory Model
  - Host memory + shared device memory
  - All compute devices share the device memory

- Execution Model
  - DAG-style task parallel execution across all available compute devices

- Programming Model
  - High-level OpenACC, OpenMP4, SYCL* (* planned)
  - Low-level C/Fortran/Python IRIS host-side runtime API + OpenCL/CUDA/HIP/OpenMP kernels (w/o compiler support)

# Supported Architectures and Programming Systems by IRIS

| ExCL* Systems | Oswald | Summit-node | Radeon | Xavier | Snapdragon |
|---|---|---|---|---|---|
| CPU | Intel Xeon | IBM Power9 | Intel Xeon | ARMv8 | Qualcomm Kryo |
| Programming Systems | • Intel OpenMP<br>• Intel OpenCL | • IBM XL OpenMP | • Intel OpenMP<br>• Intel OpenCL | • GNU GOMP | • Android NDK OpenMP |
| GPU | NVIDIA P100 | NVIDIA V100 | AMD Radeon VII | NVIDIA Volta | Qualcomm Adreno 640 |
| Programming Systems | • NVIDIA CUDA<br>• NVIDIA OpenCL | • NVIDIA CUDA | • AMD HIP<br>• AMD OpenCL | • NVIDIA CUDA | • Qualcomm OpenCL |
| FPGA | Intel/Altera Stratix 10 | | | | |
| Programming Systems | Intel OpenCL | | | | |

*ORNL Experimental Computing Laboratory (ExCL) https://excl.ornl.gov/
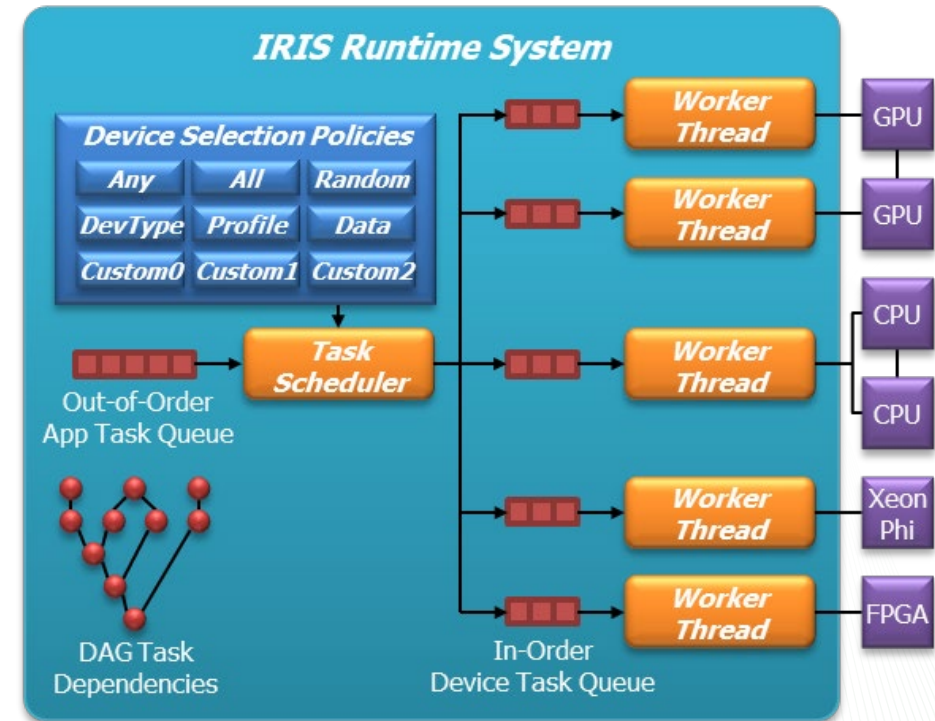
OAK RIDGE
National Laboratory

# IRIS Booting on Various Platforms

eck@xavier:~/work/brisbane-rts/apps/saxpy-py$ ./saxpy.py

```
[I] xavier [Platform.cpp:98:Init] Brisbane architectures[openmp:cuda:hip:opencl]
[T] xavier [Platform.cpp:238:InitOpenMP] OpenMP platform[0] ndevs[1]
[I] xavier [DeviceOpenMP.cpp:27:DeviceOpenMP] device[0] platform[0] device[ARMv8 Processor rev 0 (v8l)] type[64]
[T] xavier [Platform.cpp:180:InitCUDA] CUDA platform[1] ndevs[1]
[I] xavier [DeviceCUDA.cpp:44:DeviceCUDA] device[1] platform[1] vendor[NVIDIA Corporation] device[Xavier] type[128] ve
 10000] max_compute_units[8] max_work_group_size_[1024] max_work_item_sizes[2199023254528,67107840,4194240] max_block
]
[T] xavier [Loader.cpp:19:Load] libhip_hcc.so: cannot open shared object file: No such file or directory
[T] xavier [Platform.cpp:203:InitHIP] skipping HIP architecture
[T] xavier [Platform.cpp:261:InitOpenCL] OpenCL nplatforms[0]
[T] xavier [Loader.cpp:19:Load] brisbane.poly.so: cannot open shared object file: No such file or directory
[I] xavier [Platform.cpp:133:Init] nplatforms[2] ndevs[2] hub[0] polyhedral[0] profile[0]
[T] xavier [DeviceCUDA.cpp:64:Init] dev[1][Xavier] kernels[kernel.ptx]
X [ 0.  1.  2.  3.  4.  5.  6.  7.]
Y [ 0.  1.  2.  3.  4.  5.  6.  7.]
[T] xavier [DeviceCUDA.cpp:92:MemH2D] mem[4] off[0] size[32] host[0x5572cd2bf0]
[T] xavier [DeviceCUDA.cpp:142:KernelLaunch] kernel[saxpy0] dim[1] grid[8,1,1] block[1,1,1] blockOff_x[0]
[I] xavier [Device.cpp:44:Execute] task[8][saxpy0] complete dev[1][Xavier] time[0.000877]
[T] xavier [DeviceCUDA.cpp:100:MemD2H] mem[6] off[0] size[32] host[0x7f68001110]
[T] xavier [Consistency.cpp:104:ResolveWithoutPolymem] kernel[saxpy1] memcpy[6] [Xavier] -> [ARMv8 Processor rev 0 (v8
[I] xavier [Device.cpp:44:Execute] task[10][saxpy1] complete dev[0][ARMv8 Processor rev 0 (v8l)] time[0.001016]
S = 10.0 * X + Y [  0.  11.  22.  33.  44.  55.  66.  77.]
[I] xavier [Platform.cpp:595:ShowKernelHistory] kernel[saxpy1] k[0.000997][1] h2d[0.000018][2] d2h[0.000001][1]
[I] xavier [Platform.cpp:595:ShowKernelHistory] kernel[saxpy0] k[0.000189][1] h2d[0.000688][1] d2h[0.000000][0]
[I] xavier [Platform.cpp:595:ShowKernelHistory] kernel[brisbane_null] k[0.000000][0] h2d[0.000000][0] d2h[0.000155][1
[I] xavier [Platform.cpp:600:ShowKernelHistory] total kernel[0.001186] h2d[0.000706] d2h[0.000156]
[I] xavier [Platform.cpp:624:Finalize] total execution time:[0.106565] sec. initialize:[0.098530] sec. t-i:[0.008035]
```

eck@xavier:~/work/brisbane-rts/apps/saxpy-py$

# Task Scheduling in IRIS

- A task
  - A scheduling unit
  - Contains multiple in-order commands
    - Kernel launch command
    - Memory copy command (device-to-host, host-to-device)
  - May have DAG-style dependencies with other tasks
  - Enqueued to the application task queue with a device selection policy
    - Available device selection policies
      - Specific Device (compute device #)
      - Device Type (CPU, GPU, FPGA, XeonPhi)
      - Profile-based
      - Locality-aware
      - Ontology-base
      - Performance models (Aspen)
      - Any, All, Random, 3rd-party users' custom policies

- The task scheduler dispatches the tasks in the application task queue to available compute devices
  - Select the optimal target compute device according to task's device selection policy

# SAXPY Example on Xavier

- Computation
  - S[] = A * X[] + Y[]

- Two tasks
  - S[] = A * X[] on NVIDIA GPU (CUDA)
  - S[] += Y[] on ARM CPU (OpenMP)
    - S[] is shared between two tasks
    - Read-after-write (RAW), true dependency

- Low-level Python IRIS host code + CUDA/OpenMP kernels
  - saxpy.py
  - kernel.cu
  - kernel.openmp.h

OAK RIDGE
National Laboratory

# SAXPY: Python host code & CUDA kernel code

## saxpy.py (1/2)

```python
#!/usr/bin/env python

import iris
import numpy as np
import sys

iris.init()

SIZE = 1024
A = 10.0

x = np.arange(SIZE,
dtype=np.float32)
y = np.arange(SIZE,
dtype=np.float32)
s = np.arange(SIZE,
dtype=np.float32)

print 'X', x
print 'Y', y

mem_x = iris.mem(x.nbytes)
mem_y = iris.mem(y.nbytes)
mem_s = iris.mem(s.nbytes)
```

## saxpy.py (2/2)

```python
kernel0 = iris.kernel("saxpy0")
kernel0.setmem(0, mem_s, iris.iris_w)
kernel0.setint(1, A)
kernel0.setmem(2, mem_x, iris.iris_r)

off = [ 0 ]
ndr = [ SIZE ]

task0 = iris.task()
task0.h2d_full(mem_x, x)
task0.kernel(kernel0, 1, off, ndr)
task0.submit(iris.iris_gpu)

kernel1 = iris.kernel("saxpy1")
kernel1.setmem(0, mem_s, iris.iris_rw)
kernel1.setmem(1, mem_y, iris.iris_r)

task1 = iris.task()
task1.h2d_full(mem_y, y)
task1.kernel(kernel1, 1, off, ndr)
task1.d2h_full(mem_s, s)
task1.submit(iris.iris_cpu)

print 'S =', A, '* X + Y', s

iris.finalize()
```

## kernel.cu (CUDA)

```cuda
extern "C" __global__ void saxpy0(float*
S, float A, float* X) {
  int id = blockIdx.x * blockDim.x +
threadIdx.x;
  S[id] = A * X[id];
}

extern "C" __global__ void saxpy1(float*
S, float* Y) {
  int id = blockIdx.x * blockDim.x +
threadIdx.x;
  S[id] += Y[id];
}
```

OAK RIDGE
National Laboratory

# SAXPY: Python host code & OpenMP kernel code

## saxpy.py (1/2)

```python
#!/usr/bin/env python

import iris
import numpy as np
import sys

iris.init()

SIZE = 1024
A = 10.0

x = np.arange(SIZE,
dtype=np.float32)
y = np.arange(SIZE,
dtype=np.float32)
s = np.arange(SIZE,
dtype=np.float32)

print 'X', x
print 'Y', y

mem_x = iris.mem(x.nbytes)
mem_y = iris.mem(y.nbytes)
mem_s = iris.mem(s.nbytes)
```

## saxpy.py (2/2)

```python
kernel0 = iris.kernel("saxpy0")
kernel0.setmem(0, mem_s, iris.iris_w)
kernel0.setint(1, A)
kernel0.setmem(2, mem_x, iris.iris_r)

off = [ 0 ]
ndr = [ SIZE ]

task0 = iris.task()
task0.h2d_full(mem_x, x)
task0.kernel(kernel0, 1, off, ndr)
task0.submit(iris.iris_gpu)

kernel1 = iris.kernel("saxpy1")
kernel1.setmem(0, mem_s, iris.iris_rw)
kernel1.setmem(1, mem_y, iris.iris_r)

task1 = iris.task()
task1.h2d_full(mem_y, y)
task1.kernel(kernel1, 1, off, ndr)
task1.d2h_full(mem_s, s)
task1.submit(iris.iris_cpu)

print 'S =', A, '* X + Y', s

iris.finalize()
```

## kernel.openmp.h (OpenMP)

```c
#include <iris/iris_openmp.h>

static void saxpy0(float* S, float A, float*
X, IRIS_OPENMP_KERNEL_ARGS) {
  int id;
#pragma omp parallel for shared(S, A, X)
private(id)
  IRIS_OPENMP_KERNEL_BEGIN
  S[id] = A * X[id];
  IRIS_OPENMP_KERNEL_END
}

static void saxpy1(float* S, float* Y,
IRIS_OPENMP_KERNEL_ARGS) {
  int id;
#pragma omp parallel for shared(S, Y)
private(id)
  IRIS_OPENMP_KERNEL_BEGIN
  S[id] += Y[id];
  IRIS_OPENMP_KERNEL_END
}
```

OAK RIDGE
National Laboratory

# Memory Consistency Management

## saxpy.py (1/2)

```python
#!/usr/bin/env python

import iris
import numpy as np
import sys

iris.init()

SIZE = 1024
A = 10.0

x = np.arange(SIZE,
dtype=np.float32)
y = np.arange(SIZE,
dtype=np.float32)
s = np.arange(SIZE,
dtype=np.float32)

print 'X', x
print 'Y', y

mem_x = iris.mem(x.nbytes)
mem_y = iris.mem(y.nbytes)
mem_s = iris.mem(s.nbytes)
```

## saxpy.py (2/2)

```python
kernel0 = iris.kernel("saxpy0")
kernel0.setmem(0, mem_s, iris.iris_w)
kernel0.setint(1, A)
kernel0.setmem(2, mem_x, iris.iris_r)

off = [ 0 ]
ndr = [ SIZE ]

task0 = iris.task()
task0.h2d_full(mem_x, x)
task0.kernel(kernel0, 1, off, ndr)
task0.submit(iris.iris_gpu)

kernel1 = iris.kernel("saxpy1")
kernel1.setmem(0, mem_s, iris.iris_rw)
kernel1.setmem(1, mem_y, iris.iris_r)

task1 = iris.task()
task1.h2d_full(mem_y, y)
task1.kernel(kernel1, 1, off, ndr)
task1.d2h_full(mem_s, s)
task1.submit(iris.iris_cpu)

print 'S =', A, '* X + Y', s

iris.finalize()
```

**mem_s is shared between GPU and CPU**



```
cuMemcpyHtoD(dst, src)
cuMemcpyDtoH(dst, src)
memcpy(dst, src)
```

**task0**

cuMemcpy HtoD (X_{gpu}, x)

saxpy0 kernel.cu

$S_{gpu}[] = A*X_{gpu}[]$

cuMemcpy DtoH ($S_{cpu}$, $S_{gpu}$)

An inter-device memory copy command to guarantee memory consistency

Injected from the IRIS runtime and hidden from the application

**task1**

memcpy ($Y_{cpu}$, y)

saxpy1 kernel.openmp.h

$S_{cpu}[]+= Y_{cpu}[]$

memcpy (s, $S_{cpu}$)

NVIDIA GPU    *Xavier*    ARM CPU

185

OAK RIDGE National Laboratory

# Locality-aware Device Selection Policy

## saxpy.py (1/2)

```python
#!/usr/bin/env python

import iris
import numpy as np
import sys

iris.init()

SIZE = 1024
A = 10.0

x = np.arange(SIZE,
dtype=np.float32)
y = np.arange(SIZE,
dtype=np.float32)
s = np.arange(SIZE,
dtype=np.float32)

print 'X', x
print 'Y', y

mem_x = iris.mem(x.nbytes)
mem_y = iris.mem(y.nbytes)
mem_s = iris.mem(s.nbytes)
```

## saxpy.py (2/2)

```python
kernel0 = iris.kernel("saxpy0")
kernel0.setmem(0, mem_s, iris.iris_w)
kernel0.setint(1, A)
kernel0.setmem(2, mem_x, iris.iris_r)

off = [ 0 ]
ndr = [ SIZE ]

task0 = iris.task()
task0.h2d_full(mem_x, x)
task0.kernel(kernel0, 1, off, ndr)
task0.submit(iris.iris_gpu)

kernel1 = iris.kernel("saxpy1")
kernel1.setmem(0, mem_s, iris.iris_rw)
kernel1.setmem(1, mem_y, iris.iris_r)

task1 = iris.task()
task1.h2d_full(mem_y, y)
task1.kernel(kernel1, 1, off, ndr)
task1.d2h_full(mem_s, s)
task1.submit(iris.iris_data)

print 'S =', A, '* X + Y', s

iris.finalize()
```
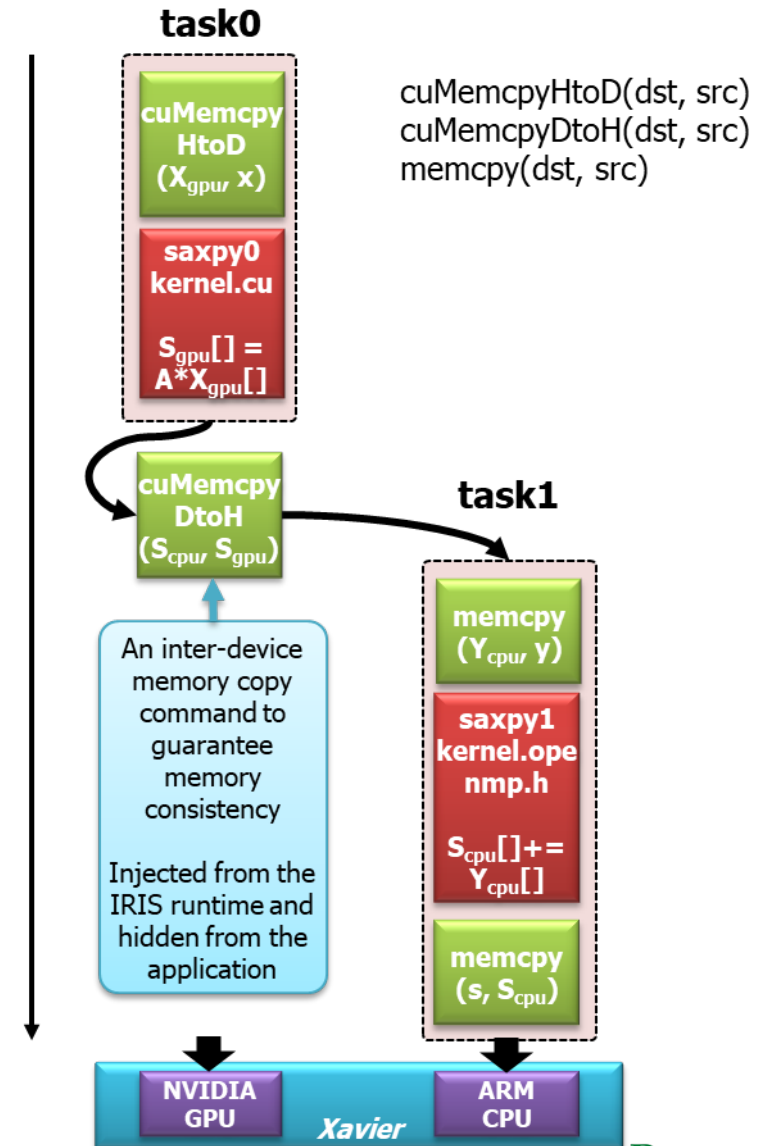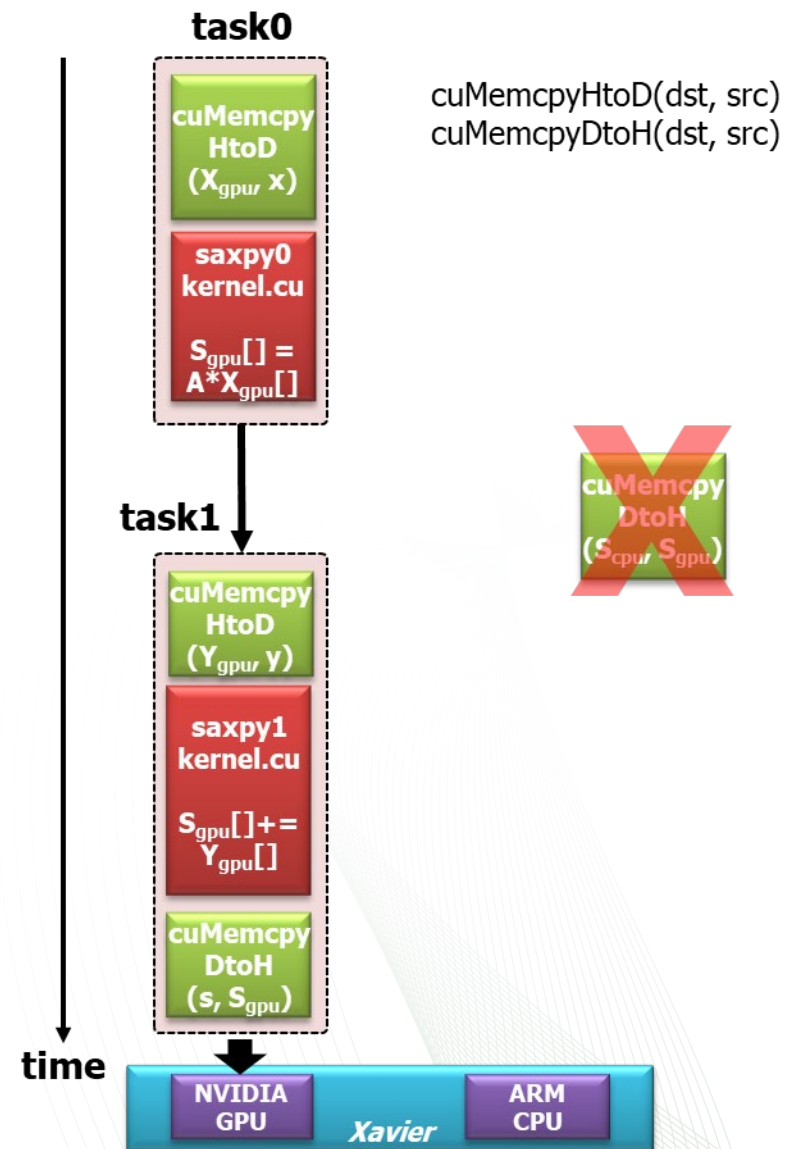
iris_data selects the device that requires minimum data transfer to execute the task

cuMemcpyHtoD(dst, src)
cuMemcpyDtoH(dst, src)

**task0**
cuMemcpy HtoD ($X_{gpu}$, x)
saxpy0 kernel.cu
$S_{gpu}[] = A*X_{gpu}[]$

**task1**
cuMemcpy HtoD ($Y_{gpu}$, y)
saxpy1 kernel.cu
$S_{gpu}[] += Y_{gpu}[]$
cuMemcpy DtoH (s, $S_{gpu}$)

cuMemcpy DtoH ($S_{cpu}$, $S_{gpu}$)

**time**

NVIDIA GPU    Xavier    ARM CPU

OAK RIDGE
National Laboratory

# IRIS: Task Scheduling Overhead – Running One Million (Empty) Tasks

## ntasks.py

```
#!/usr/bin/env python

import iris

iris.init()

NTASKS = 1000000

t0 = iris.timer_now()

for i in range(NTASKS):
  task = iris.task()
  task.submit(iris.iris_random, False)

iris.synchronize()

t1 = iris.timer_now()
print 'Time:', t1 - t0

iris.finalize()
```

**CPU or GPU randomly**

**asynchronous task submission**

**concurrent tasks execution on multiple devices**

user@xavier:~/work$ ./ntasks.py
Time: 11.46s

| Throughput | Latency |
|---|---|
| 87,268 tasks/sec | 11.4 µs/task |

# Closing



| | | |
|---|---|---|
| Jeffrey Vetter, PI | Seyong Lee, Programming Systems | Jungwon Kim, Runtime Systems |
| Mehmet Belviranli, Apps, Modeling, Ontologies | Richard Glassbrook, Project Manager | Steve Moulton, Systems Engineer |
| Abdel-Kareem Moadi, Software/Hardware Engineer | Seth Hitfield SDR | Blaise Tine, Intern, Georgia Tech |
| Mohammad Monil, Intern, Oregon | Austin Latham SYCL | |

## Summary

- Architectural specialization

- Performance portability of applications and software

- DSSoC ORNL project investigating on performance portability of SDR
  - Understand applications and target architectures
  - Use open programming models: OpenACC, OpenCL, OpenMP
  - Developing intelligent runtime systems: IRIS

- Goal: scale applications from Qualcomm Snapdragon to DoE Summit Supercomputer with minimal programmer effort

- Work continues…

## Acknowledgements

- Thanks to staff and students for the work!

- Thanks to DARPA, DOE for funding our work!

**OAK RIDGE** National Laboratory