# Towards Generalized On-Chip Communication for Programmable Accelerators in Heterogeneous Architectures
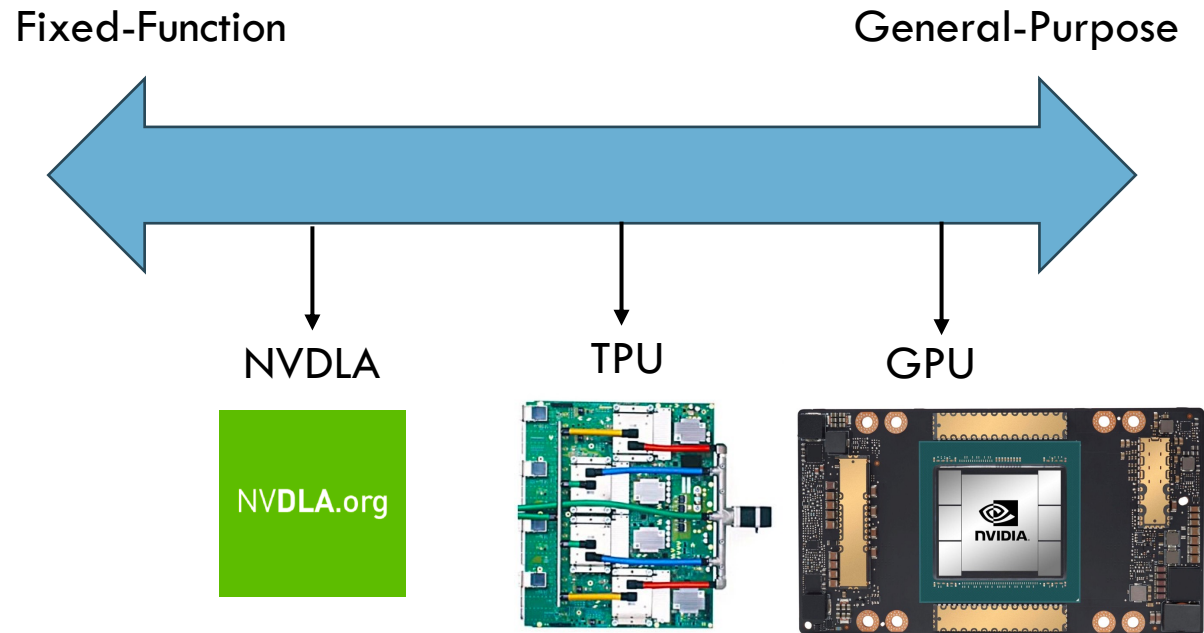
**Joseph Zuckerman**[1]

John-David Wellman[2], Ajay Vanamali[1], Manish Shankar[1], Gabriele Tombesi[1],

Karthik Swaminathan[2], Kevin Lee[1], Mohit Kapur[2], Robert Philhower[2],

Pradip Bose[2], Luca P. Carloni[1]

**Columbia University**[1], **IBM Research**[2]

**DOSSA Workshop '24**
**6/30/24**

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

ESP

CS @CU COMPUTER SCIENCE

# Accelerator Specialization

Fixed-Function                General-Purpose

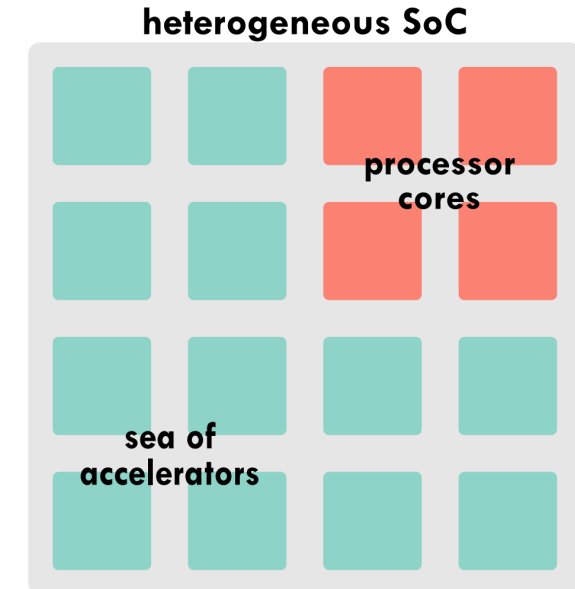NVDLA          TPU        GPU



- *Programmable accelerators* are optimized for an application domain, but execute instructions
  - Generated by a compiler, another tool, or hand written
- Balance specialization with programmability
  - Have made them widely popular in industry and academia
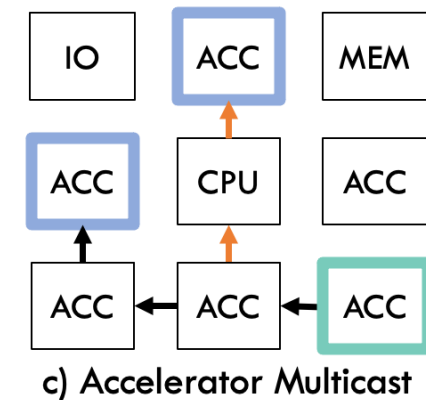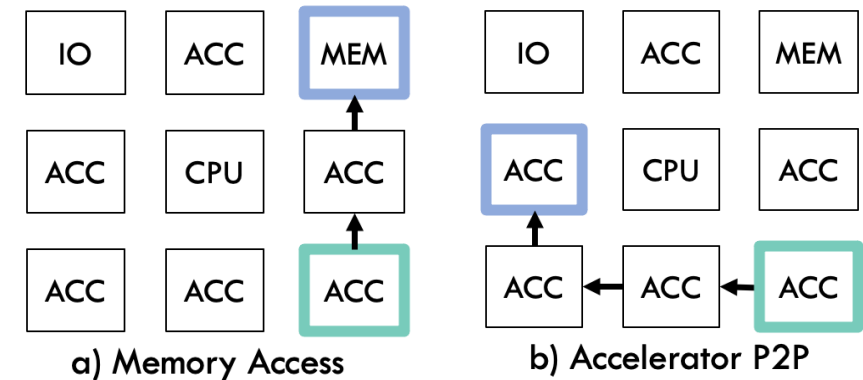- Mostly studied as monolithic entities

# The Age of Heterogeneous Computing

- SoCs are growing increasingly heterogeneous across various computing domains
  - CPUs, GPUs, accelerators, I/O peripherals, sensors…
  - Possibly multiple instances of complex, programmable accelerators
- Heterogeneity increases engineering effort
  - Capabilities of new generations limited by team size
  - Biggest challenges are in system-integration

**heterogeneous SoC**

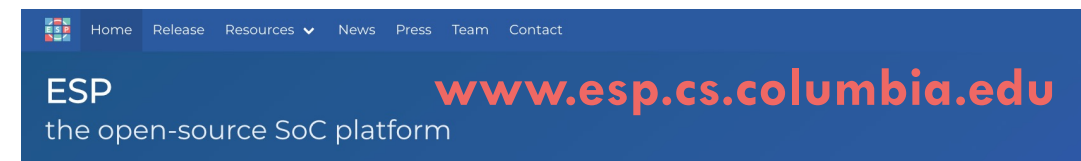**processor cores**

**sea of accelerators**

# Challenges and Opportunities for SoC Design

- Integration of programmable accelerators in heterogeneous architectures is not well-studied

- Programmability at the accelerator level could benefit from flexibility at the system level

- Complex workloads might feature multiple data-dependency patterns and require synchronization

- Can we develop *platform-level services* that provide these capabilities transparently to the design of the accelerator?
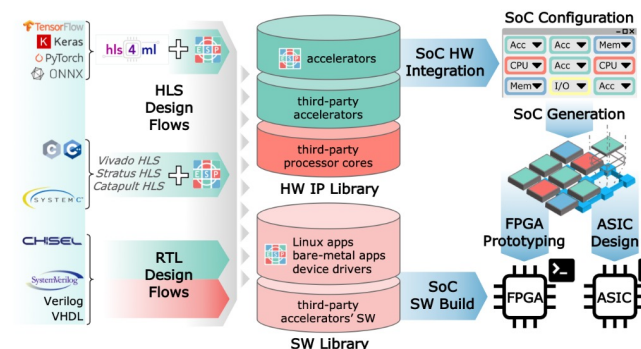


a) Memory Access

b) Accelerator P2P

c) Accelerator Multicast

# **ESP** : An Open-Source Platform for SoC Research

- We propose a set of enhancements to an SoC architecture to support the integration of programmable accelerators
  - Flexible point-to-point accelerator
  - Multicast communication supported by the NoC
  - Intra-accelerator communication
  - A simple accelerator interface to manage communication modes
  - Accelerator ISA extensions for on-chip communication
- Implement these enhancements in **ESP**, but the principles can apply to other architectures
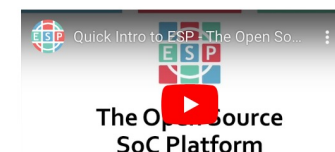


www.esp.cs.columbia.edu

# ESP Background

# ESP Architecture

- Multi-Plane NoC
- Many-Accelerator
- RISC-V Cores
- Distributed Memory

The ESP architecture implements a **distributed** system, which is **scalable**, **modular** and **heterogeneous**, giving processors and accelerators similar weight in the SoC [3]

memory tile

ESP accelerator tile

processor tile

third-party accelerator tile

I/O tile

memory tile

multi-plane NoC

[3] Mantovani, ICCAD '20

# ESP Accelerator Socket

# ESP Accelerator Flow

Developers focus on the **high-level specification**, **decoupled** from memory access, system communication, hardware/software interface

# ESP SoC Flow
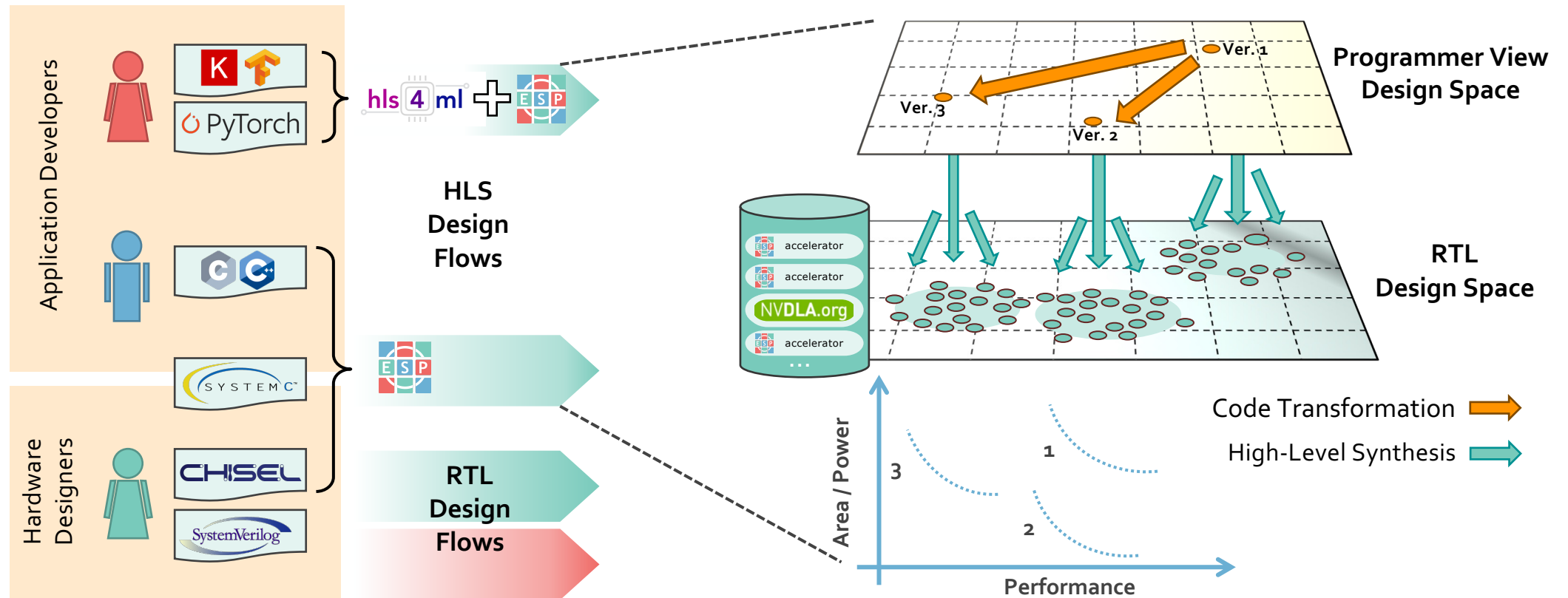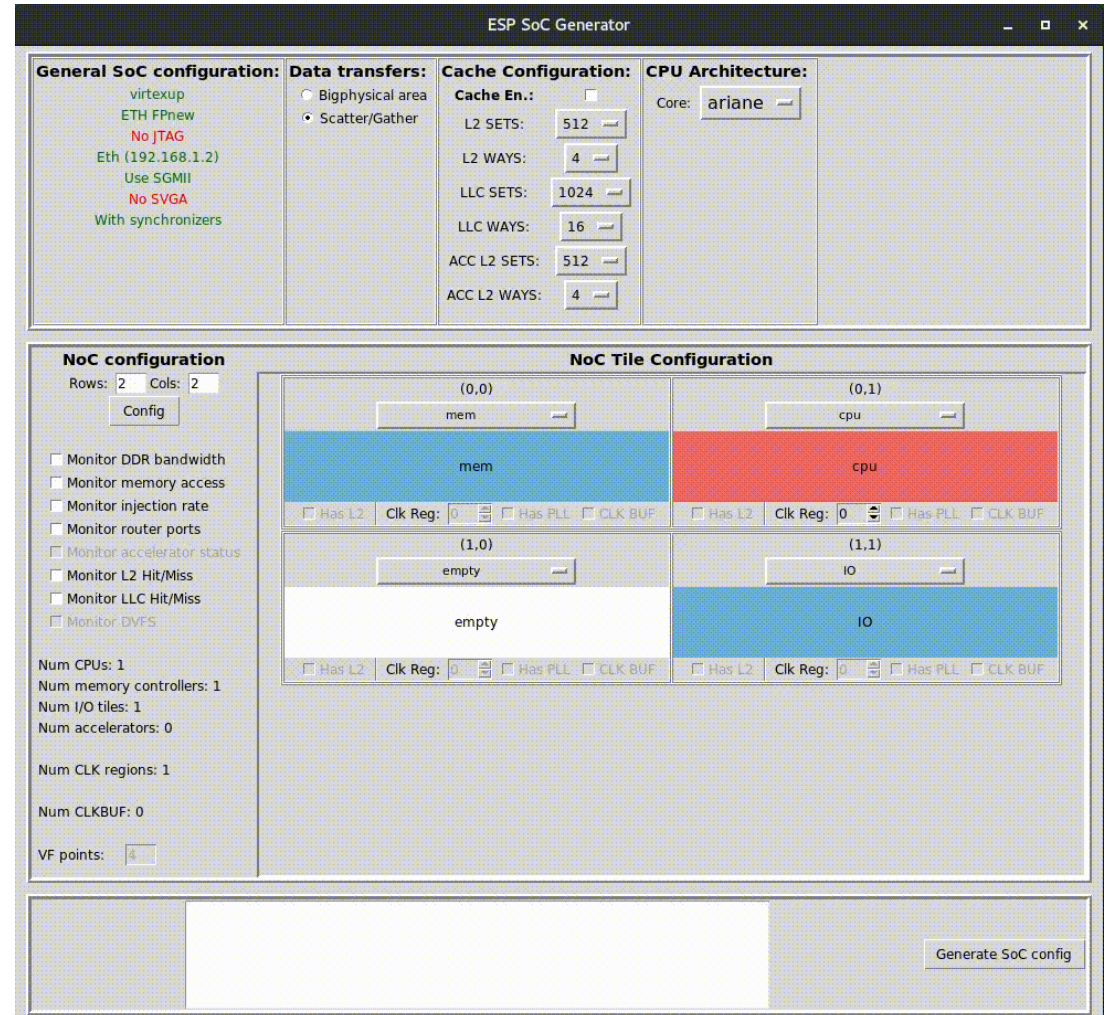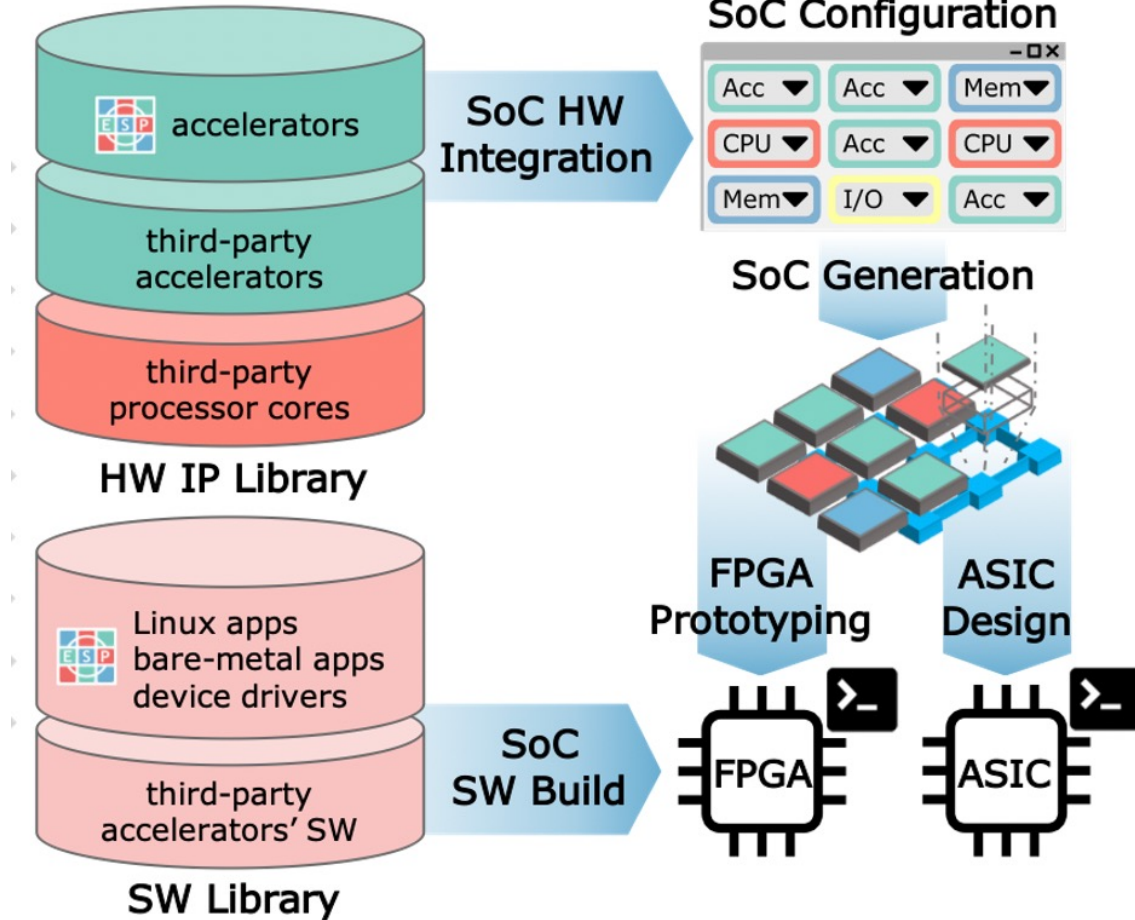
# Proposed Architecture

# Programmable Accelerator Model

- Loosely-coupled from the host core
  - Won't boot an OS, etc.

- Executes coarse-grained tasks on large workloads
  - Needs a DMA interface to bring in large amounts of data from the SoC

- Highly-customized datapath and private local memory

- Can rely on a lightweight core for instruction fetch, decode, commit, etc.
  - E.g. RISC-V Rocket Core + RoCC interface
  - ISA extensions for leveraging the datapath



Control Core
(e.g. RISC-V Rocket)

Control Interface
(e.g. RoCC)

Compute Substrate
(e.g. CGRA)

L1 Cache

Private Local Memory

Programmable Accelerator

Memory Interface (e.g. AXI, TileLink, ESP)

DMA Controller

# Flexible Point-to-Point Accelerator Communication

- Point-to-point (P2P) accelerator communication improves performance of applications with data dependencies
  - Avoids round-trip to memory
  - Pipelining with fine-grained synchronization
- Implemented with a *pull-based model*
  - Consumer initiates each exchange to satisfy the *consumption assumption*
- Previously, data access patterns must be exactly the same
  - Relaxed this by adding a *length*
  - Total amount of data must match



13

# Multicast NoC

- Support multiple-consumer dependencies with a single NoC message

- With minimal modifications to the ESP NoC

- Header *flit* modified to carry an array of destinations
  - Max # of consumers limited by NoC bitwidth

- Lookahead routing logic replicated for each destination

- Integrated with ESP's P2P capabilities
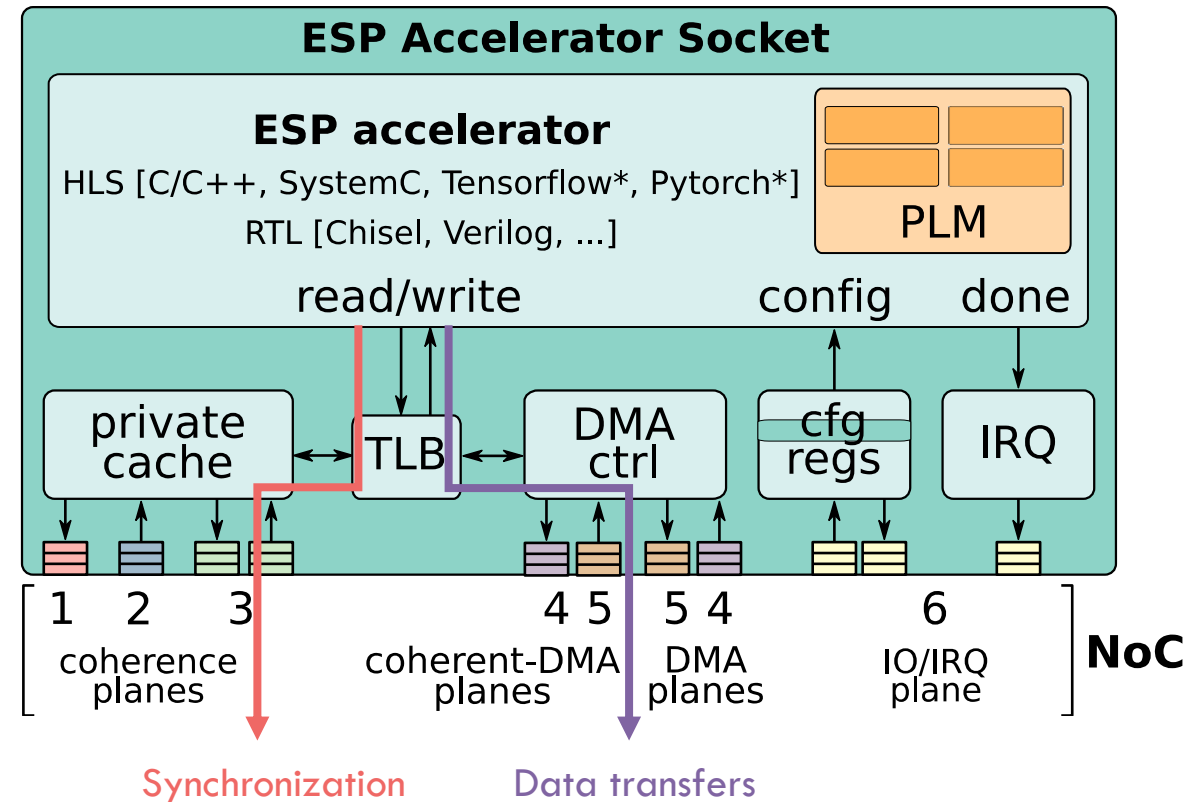
Unicast Header Flit

| PREAMBLE | SRC | DST | MSG TYPE | RESERVED | ROUTING |
|---|---|---|---|---|---|

Multicast Header Flit

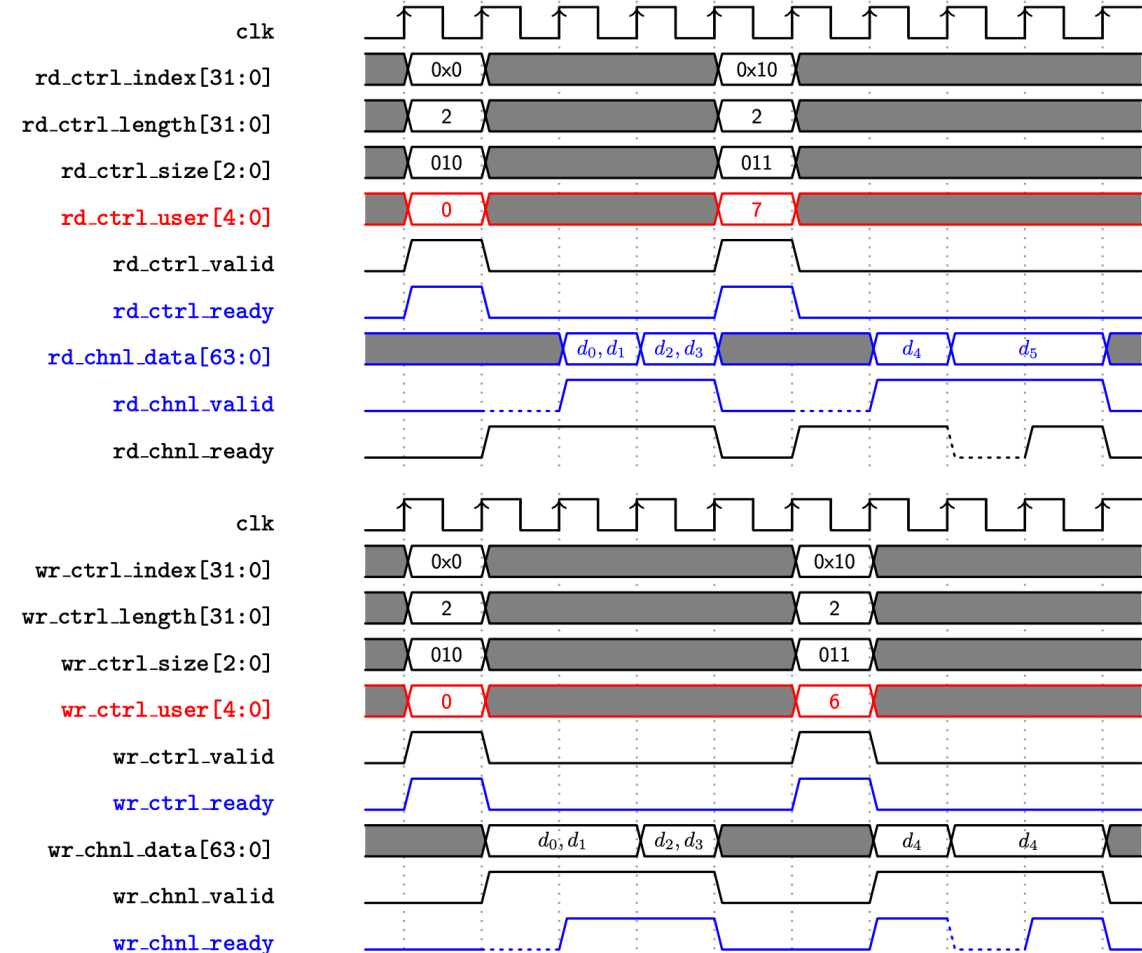| PREAMBLE | SRC | DST | MSG TYPE | DST [ ] | VAL [ ] | RESERVED | ROUTING |
|---|---|---|---|---|---|---|---|

# Accelerator Synchronization

- Programmable accelerators might need to *synchronize* with each other throughout their execution
  - E.g. ensuring they have both arrived to a pre-defined point in the workload before exchanging data
- Requires *coherence* between accelerators
  - ESP supports coherent accelerators, but coherence must be used for all data transfers → inefficient
- Enable *dynamic switching* between DMA for data transfers and coherence for synchronization

**ESP Accelerator Socket**

**ESP accelerator**

HLS [C/C++, SystemC, Tensorflow*, Pytorch*]

RTL [Chisel, Verilog, …]

PLM

read/write          config        done

private cache     TLB     DMA ctrl     cfg regs     IRQ

1   2   3          4  5    5  4          6          **NoC**

coherence planes     coherent-DMA planes     DMA planes     IO/IRQ plane

Synchronization          Data transfers

15

# Accelerator Interface + ISA Extensions

- Give control of the mode of each transaction to the accelerator itself
  - Just one additional signal on the read and write control interfaces
- For reads, the *user* field specifies the source of the data
  - 0 for memory
  - 1 to N-1 for other accelerators
- For writes, the *user* field specifes the # of consumers of the data
  - 0 for memory
  - 1 for regular P2P
  - 2 to N-1 for a multicast P2P
- ISA Extensions
  - Initiate DMA (IDMA) – encodes the index, length, size, and user as expected by ESP
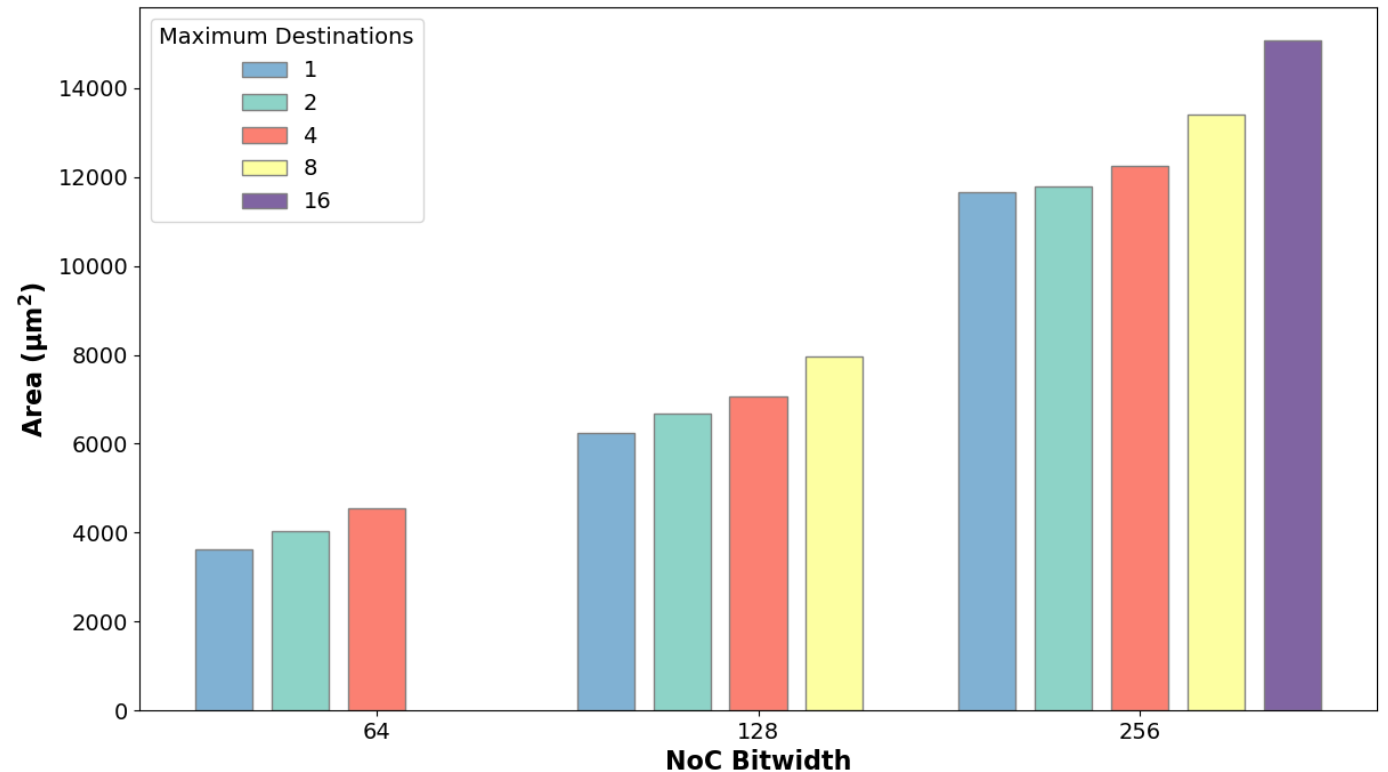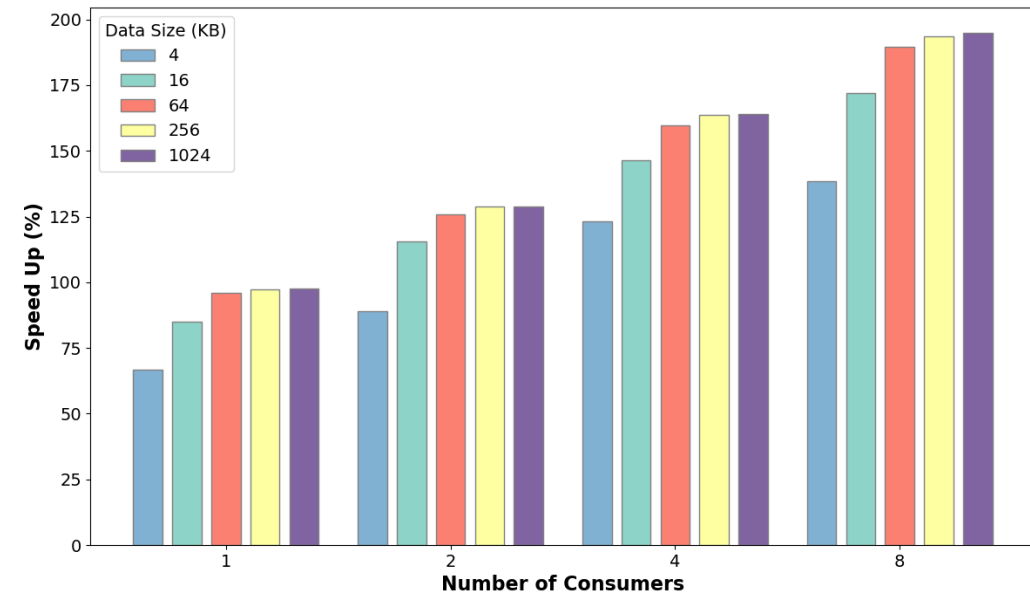  - Check DMA (CDMA) – query if DMA has finished
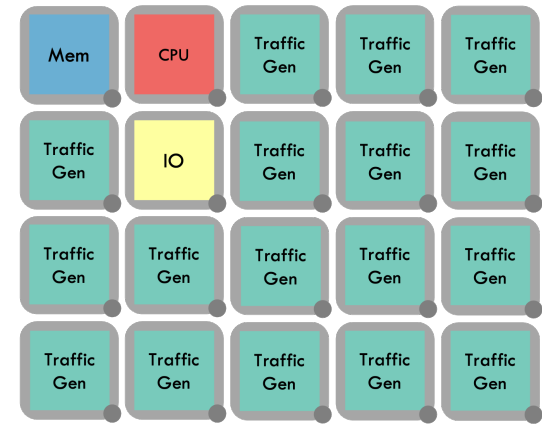
# Results

# Results: Multicast Router Area

- Synthesized in Cadence Genus targeting a 12nm technology

- Varying the bitwidth of the NoC and the maximum number of supported multicast destinations

# Results: Multicast Performance



- Developed a 5x4 SoC prototype with 17 traffic generator accelerators

- 256-bit NoC allows us to test multicast to 16 consumers

- Deployed on a Xilinx VCU128 FPGA at 78MHz

- Baremetal "workload" that creates a 1 to N data dependency

- Compared multicast performance to communication through shared memory

# Thank you from the ESP team!

**SLD** sld.cs.columbia.edu    ColumbiaSld    esp.cs.columbia.edu    sld-columbia/esp    c/ESP-platform

## Towards Generalized On-Chip Communication for Programmable Accelerators in Heterogeneous Architectures

**Joseph Zuckerman[1]**

John-David Wellman[2], Ajay Vanamali[1], Manish Shankar[1], Gabriele Tombesi[1],

Karthik Swaminathan[2], Kevin Lee[1], Mohit Kapur[2], Robert Philhower[2],

Pradip Bose[2], Luca P. Carloni[1]

**Columbia University[1]**, **IBM Research[2]**

DOSSA 2024

COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

CS@CU COMPUTER SCIENCE