# FPGA or CGRA?
# Reconfigurable Architectures Suitable for High-Performance Computing

**Kentaro Sano**

**RIKEN Center for Computational Science (R-CCS)**

# Summary of this Talk

- **Reconfigurable data-flow computing as promising architecture for HPC**

- **FPGA**
  - ✓ **ESSPER** : Elastic and scalable FPGA-cluster system for high-performance reconfigurable computing, as Prototype FPGA cluster for HPC
  - ✓ Lessons we learned

- **CGRA (Coarse-grained reconfigurable array)**
  - ✓ **RIKEN CGRA research**
  - ✓ What we have studied so far.
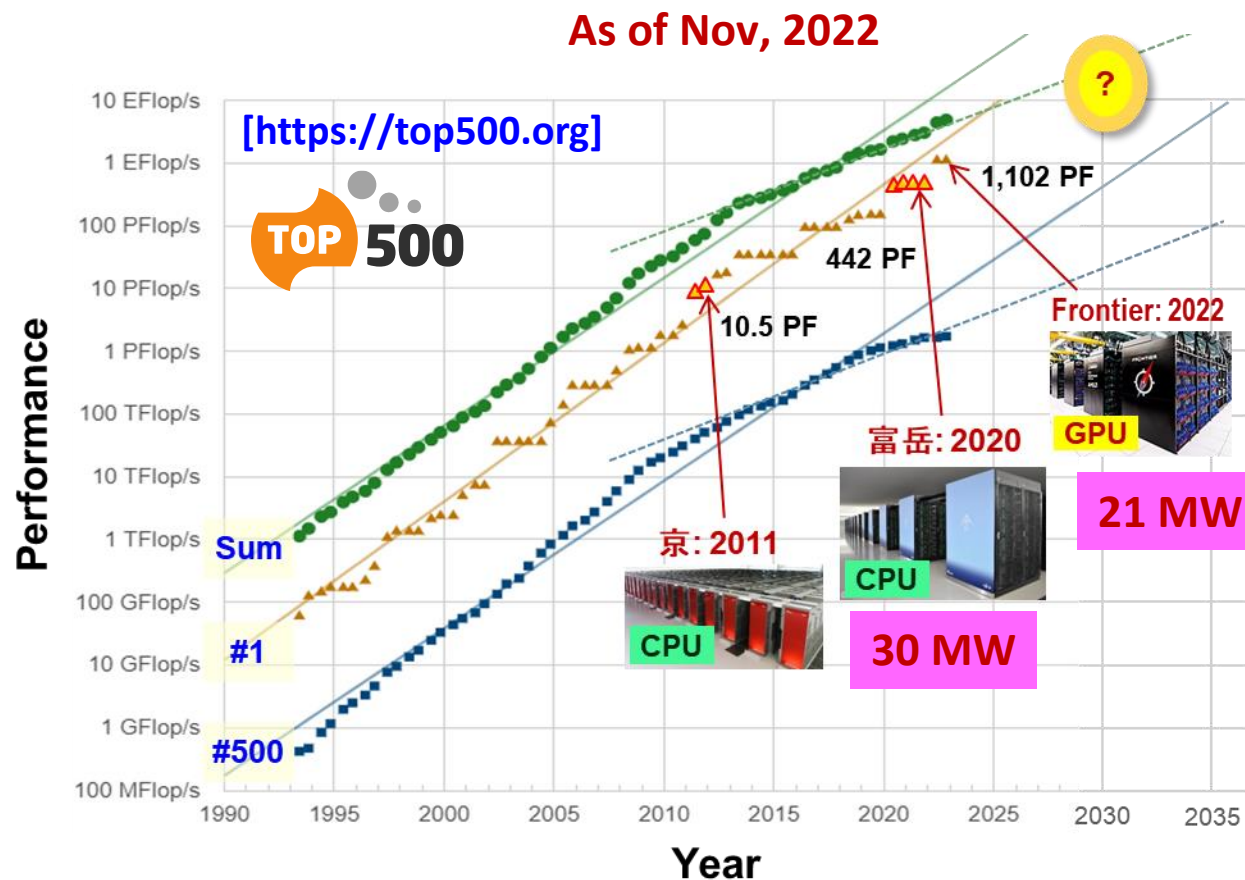
# Introduction

- **World ranking of supercomputers**
  - ✓ TOP500: Ranking of HPL performance
  - ✓ CPU-based vs. GPU/Acc-based
  - ✓ Perf improvement slowed down around 2015.
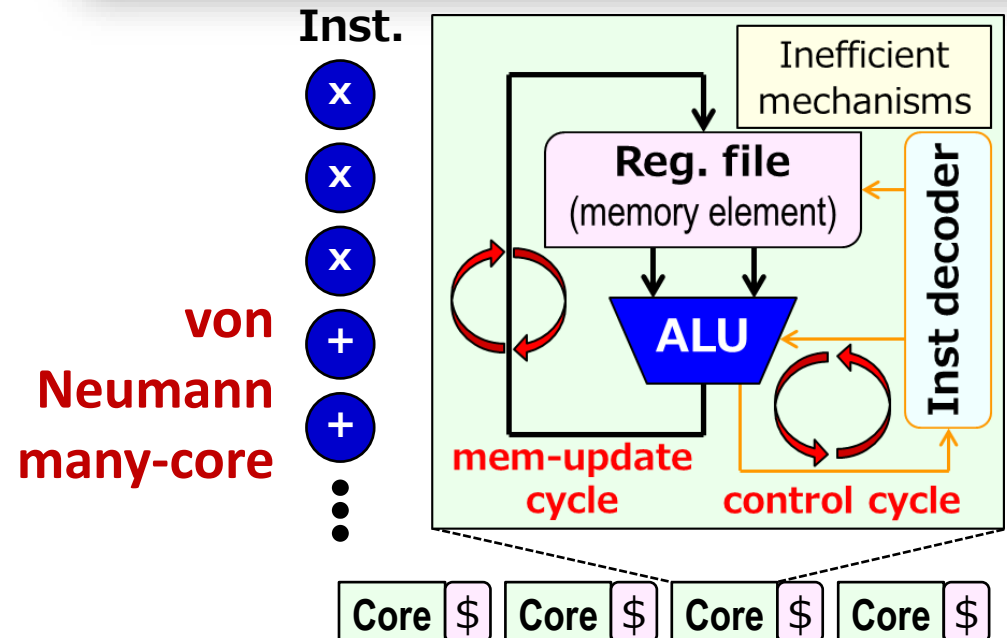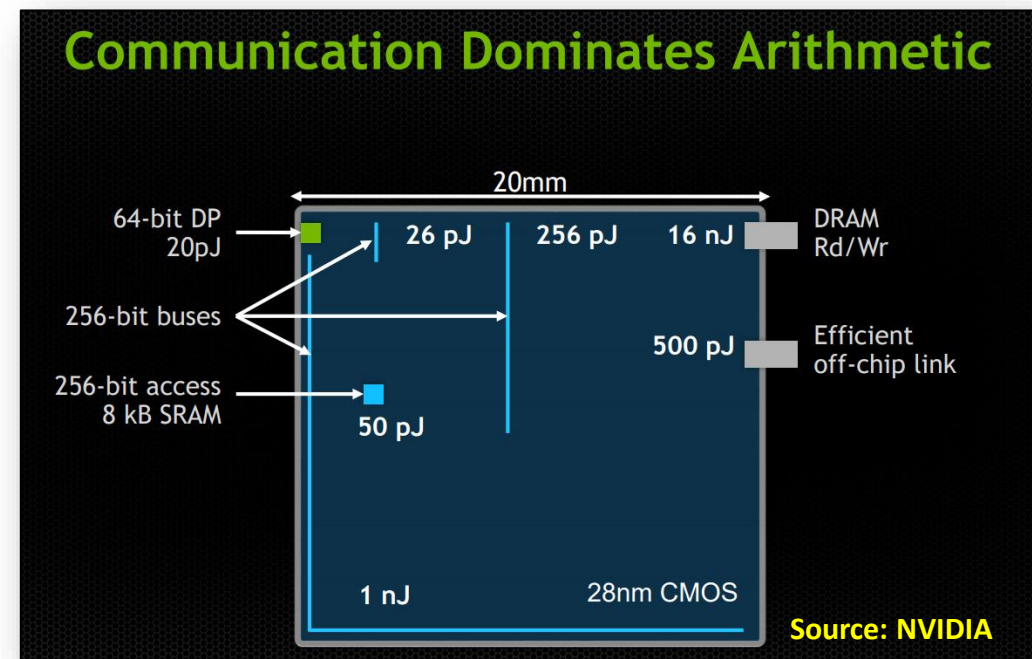
- **System performance is limited by system power.**
  - ✓ Reached tens of MW
    (Fugaku: 30MW, Frontier: 21MW for HPL)
  - ✓ Not easy to further increase
    (100MW is not real for SDGs & cost.)

- **With capped power budget, need to increase performance per power**

As of Nov, 2022

[https://top500.org]

# What Eats Power?

- **Data movement** rather than computing
  - ✓ We should remove unnecessary data movement, and make it shorter.

- **Unsuitable architecture**

  **with low efficiency and scalability**
  - ✓ von-Neumann architectures (CPU & GPU) cannot efficiently scale due to
    - ➤ *memory-bottlenecked structure;* such as register files and NoC w/ LLC for multiple cores
    - ➤ *Extra mechanisms* consuming power just to increase IPC, such as out-of-order and branch predictor.

- **Semiconductor scaling cannot save it.**
  - ✓ Power improvement per generation is limited while can still increase transistors per area for advanced tech nodes like 5, 3, 2, and 1.5nm …



**Communication Dominates Arithmetic**

20mm
64-bit DP 20pJ — 26 pJ — 256 pJ — 16 nJ — DRAM Rd/Wr
256-bit buses
256-bit access 8 kB SRAM — 50 pJ
500 pJ — Efficient off-chip link
1 nJ — 28nm CMOS
Source: NVIDIA



Inst.
x
x
x
von
Neumann + 
many-core +

Inefficient mechanisms
Reg. file (memory element)
Inst decoder
ALU
mem-update cycle    control cycle
Core $ | Core $ | Core $ | Core $

# What can Save Us?

- **Data-flow computing architecture**
  - ✓ Localized data-movement
  - ✓ No memory bottleneck;
    distributed and pipelined ALUs with
    regular/simplified memory access
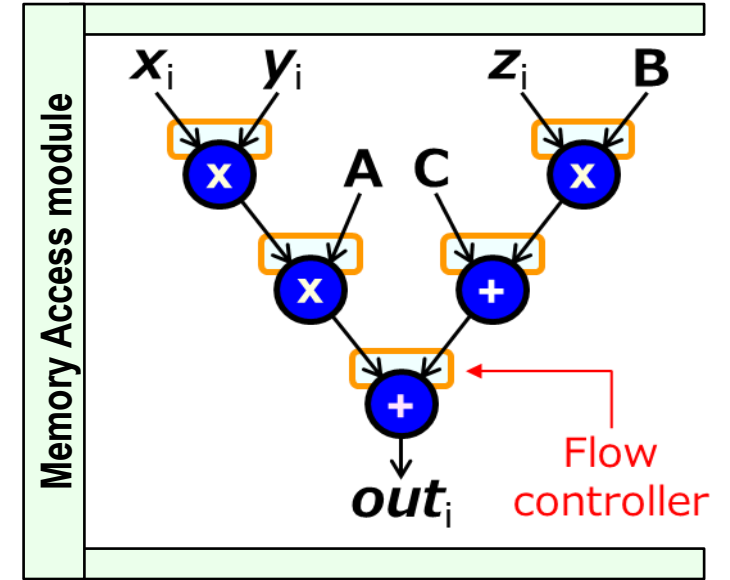  - ✓ No extra mechanisms for non-computing

- **Circuit reconfigurability is key.**
  - ✓ Giving programmability
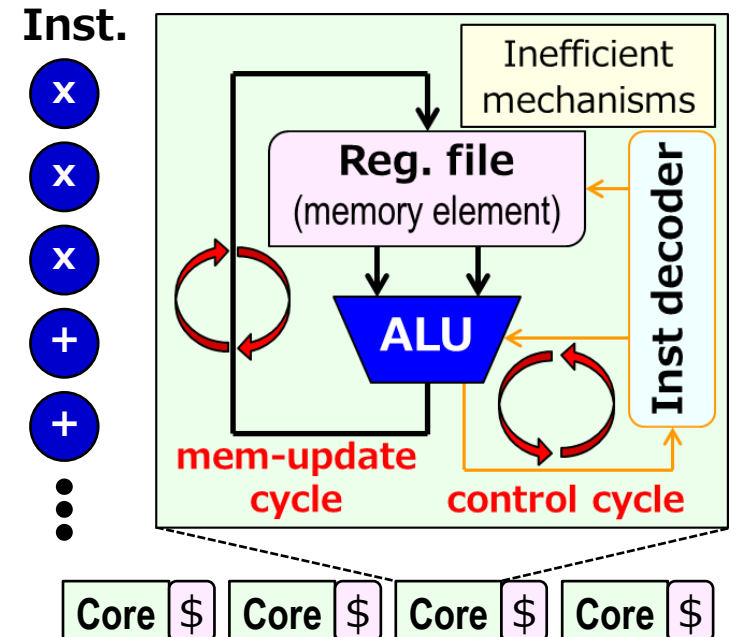  - ✓ Higher efficiency for target problems

- **What candidate technologies for reconfigurable data-flow?**
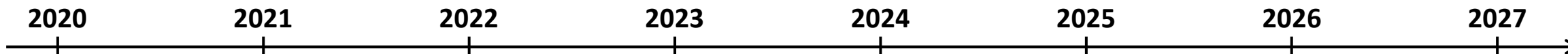  **FPGA and CGRA?**



**Data-flow computing**

**von Neumann many-core**

# Goal and Roadmap of Processor Research Team

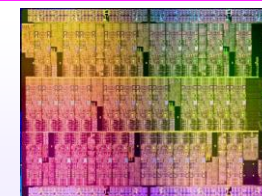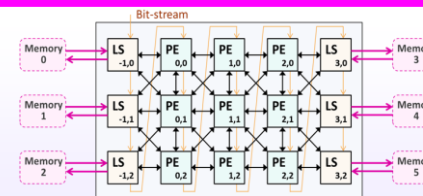## Goal: Establish HPC architectures suitable in Post-Moore Era

2020 — 2021 — 2022 — 2023 — 2024 — 2025 — 2026 — 2027

### 1. Advancement of Fugaku
- ✓ **Functional extension with FPGAs** (FPGA cluster, ESSPER)
- ✓ SoC, system software, applications

**ESSPER**
Elastic and Scalable System for High-Performance Re-configurable Computing
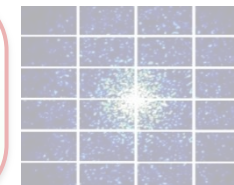
### 2. Exploration of New HPC Architectures
- ✓ Data-flow-based accelerators (**CGRA**)
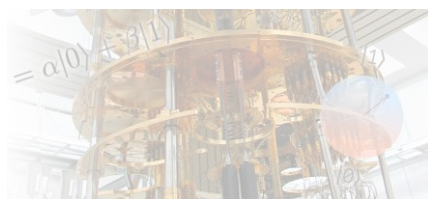- ✓ **Next-generation HPC systems**

### 3. Near-sensor / Near-storage Processing
- ✓ FPGA-based processing for **X-ray imaging detector** (Citius @Spring-8)
- ✓ Planning collaboration with next-generation radiation facility in Tohoku

### 4. Backend of Fault-Tolerant Quantum Computers
- ✓ Specialized hardware in digital circuits for **quantum error correction**
- ✓ FPGA cluster "ESSPER2" with Intel Agilex-M FPGA (7nm)

**We are hiring researchers! Contact me.**

# Prototype FPGA Cluster
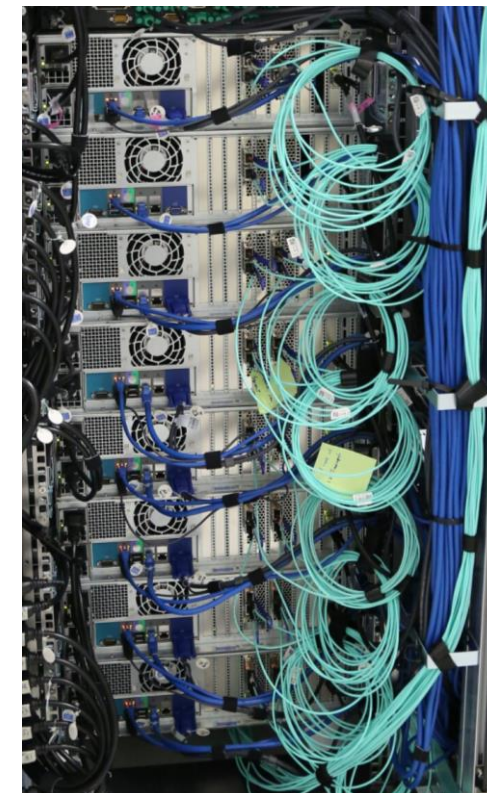# for Supercomputer Fugaku

Open-Access paper

# This Work

> **Goal : Design & demonstrate a proof-of-concept FPGA cluster for HPC research**

- **ESSPER** : Elastic and scalable FPGA-cluster system
  for high-performance reconfigurable computing

Open-Access paper

- **Contributions**
  - ✓ **Design concept** of FPGA cluster for HPC
  - ✓ **Classification** of FPGA cluster architectures
  - ✓ **Proposed system stack** with software-bridged APIs
  - ✓ **Implementation and evaluation** for FPGA-based extension
    of the world's top-class supercomputer, Fugaku

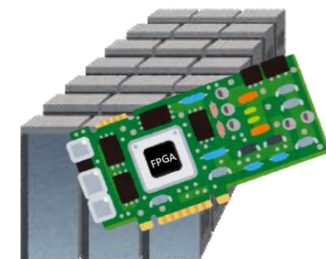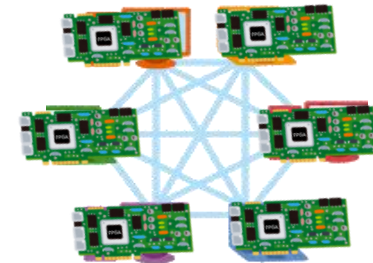# FPGAs have yet to be Mainstream in HPC.

System architecture not matured yet.

Still have **system-level challenges** for FPGA-based HPC

**Productive customizability** for computing HW

**Performance scalability** with multiple FPGAs

**Interoperability** with existing HPC systems

# Challenges and Approaches for FPGA-based HPC

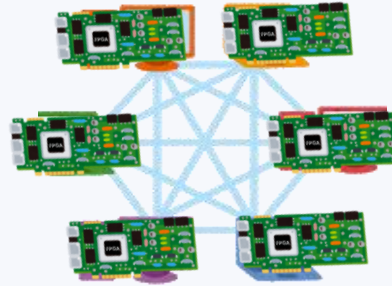## Productive customizability for computing HW

✓ Able to implement various hardware (algorithms) on FPGA



➤ No OpenCL (not limit computing models)
➤ FPGA Shell & HLS/HDL programming, where any hardware can be easily implemented

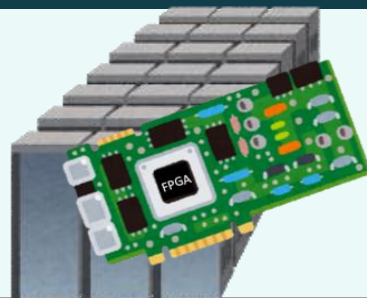## Performance scalability with multiple FPGAs

✓ Inter-FPGA communication available
✓ Allow users to easily try multi-FPGA applications



➤ FPGA Shell supporting high-bandwidth and low-latency network dedicated to FPGAs
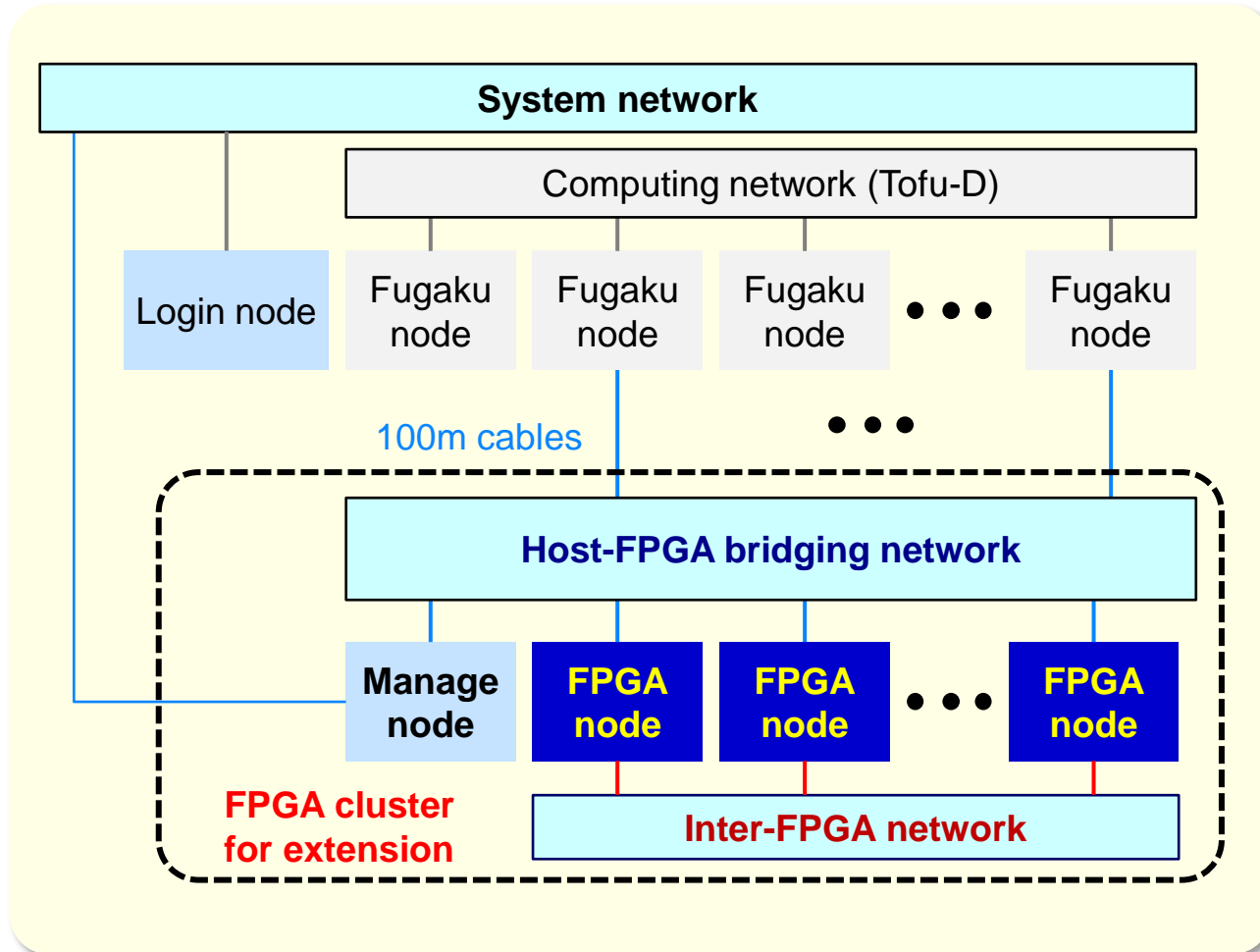
## Interoperability with existing HPC systems

✓ Able to easily extend existing systems with FPGAs
✓ Can we extend Supercomputer Fugaku?



➤ Software-bridged APIs to access FPGAs remotely through host-FPGA bridging network

# Architecture of ESSPER



- ✓ **Productive customizability**
  - ➤ No OpenCL (not limit computing models)
  - ➤ FPGA Shell & HLS/HDL programming, where any hardware can be easily implemented
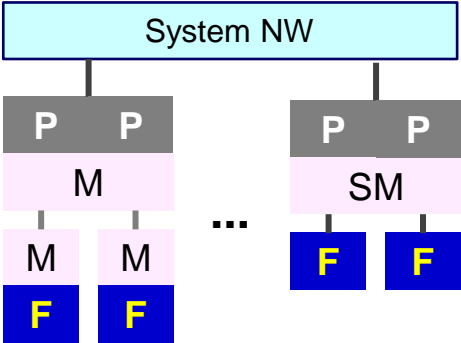
- ✓ **Performance scalability**
  - ➤ FPGA Shell supporting high-bandwidth and low-latency network dedicated to FPGAs
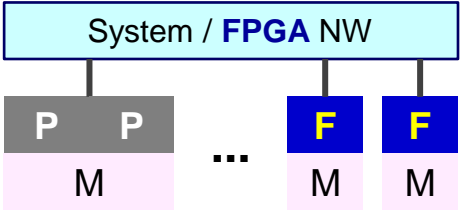
- ✓ **Interoperability**
  - ➤ Software-bridged APIs to access FPGAs remotely through host-FPGA bridging network
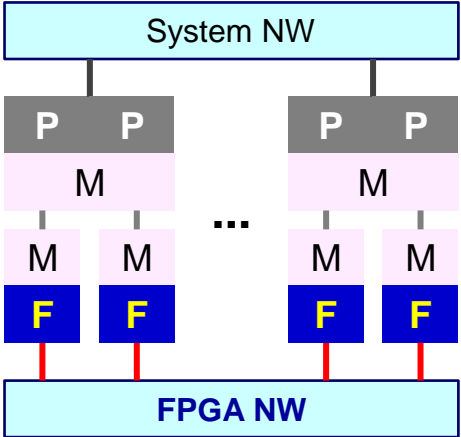
# Architecture Classification



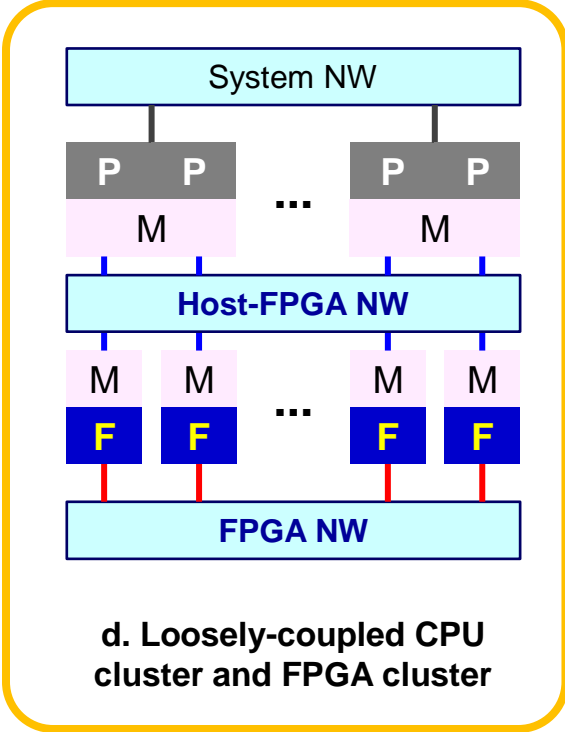(S)M — (Shared) memory  P — CPU  F — FPGA  NW — Network

a. Cluster of CPUs with FPGAs (distributed or shared memory)

b. Cluster of CPUs and FPGAs

c. Clusters of CPUs with inter-connected FPGAs

d. Loosely-coupled CPU cluster and FPGA cluster

**Our architecture**

# Related Work : FPGAs Clusters in HPC/DC

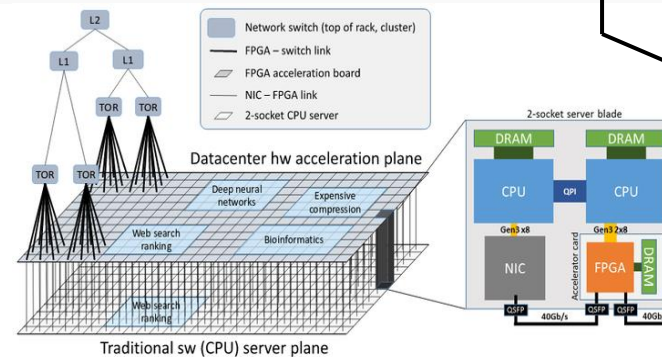| FPGA NW Type | Direct network | Indirect network | Indirect circuit-switching nw |
|---|---|---|---|
| **Characteristics** | p2p-connection without switches, typical: torus | connection with switches, typical: Ethernet | connection with optical switch (MEMS) |
| **Switching** | circuit or packet (w/ router) | packet | circuit or packet (w/ router) |
| **Pros** | **low latency** | **flexibility, small diameter** | **low latency, flexibility** |
| **Cons** | inflexibility, large diameter | higher latency, complex | expensive, signal attenuation |

**Representative systems**



**Cygnus @ U of Tsukuba**

**Archi-C**

**Novo-G# @ U of Florida**

**Archi-C**

**Catapult @ Microsoft**

**Archi-B**

**Noctua @ Paderborn U**

**Archi-C**

# System Design

# Hardware Organization of ESSPER



**FPGA Shell**

**Computing nodes of Supercomputer Fugaku**

10GBASE-SR

Frontend Server (fsarchitgate)

CAD Server #1

CAD Server #2

File Server

ARM Server #0 (armserv0)

ARM Server #1 (armserv1)

ARM Server #2 (armserv2)

10GBASE-T Ethernet Switch

LCD, Keyboard switch

Infiniband 100G Switch

1000BASE-T Switch — for IPMI

FPGA Server #1 (fpgaserv1) — FPGA #1, FPGA #2

FPGA Server #2 (fpgaserv2) — FPGA #3, FPGA #4

FPGA Server #3 (fpgaserv3) — FPGA #5, FPGA #6

FPGA Server #8 (fpgaserv8) — FPGA #15, FPGA #16

3m 100G IB

10GBASE-T x 10

QSFP28 cables

100GbE Switch or Ring network

Management

**FPGA Shell** — EMIF 1, EMIF 2, EMIF 3, EMIF 4, CCIP, Avalon-MM, Read DMA, Write DMA, HLS Computing Core, Avalon-ST, 4x4 Crossbar, 512-bit, usr_clk 250+ MHz, FC #1, FC #2, CycleCounter, DC FIFO, Stream Computing Core, 256-bit, WidthConv, 377 MHz, SLIII #1, SLIII #2, AFU, FIM

**Service servers**
- CAD servers
- Storage server
- ARM servers

**CPU - FPGA network**
- 100G Infiniband
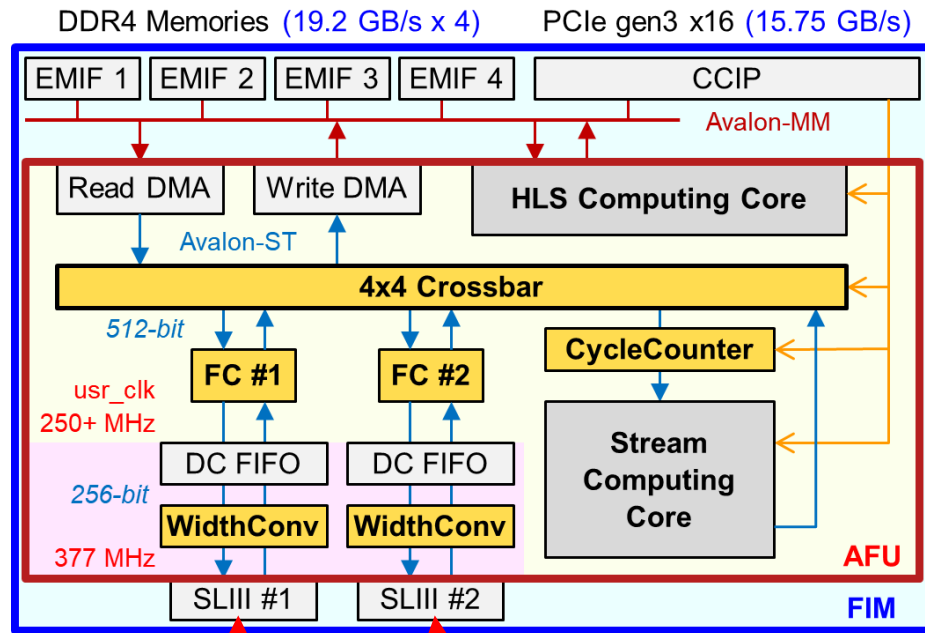- Software-bridged driver (R-OPAE)

**FPGA cluster**
- x86 host servers
- FPGA boards
- Inter-FPGA network

**FPGA Shell (SoC)**
- AFU Shell design
- User HW modules can be embedded for custom computing.
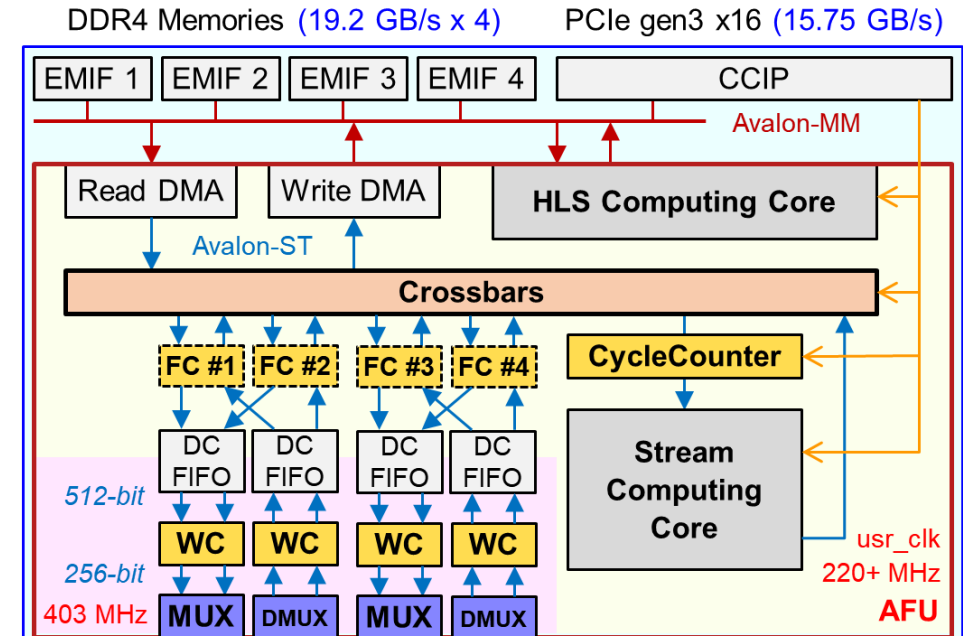
# FPGA Shells for Direct and Indirect Networks

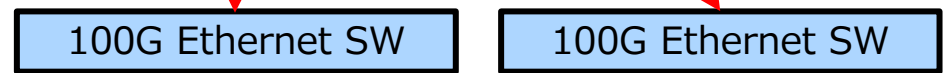## Direct connection network (DCN)



## Indirect network (VCSN)



Tomohiro Ueno, Atsushi Koshiba, Kentaro Sano, "Virtual Circuit-Switching Network with Flexible Topology for High-Performance FPGA Cluster," Procs. of ASAP, pp.41-48, 2021.
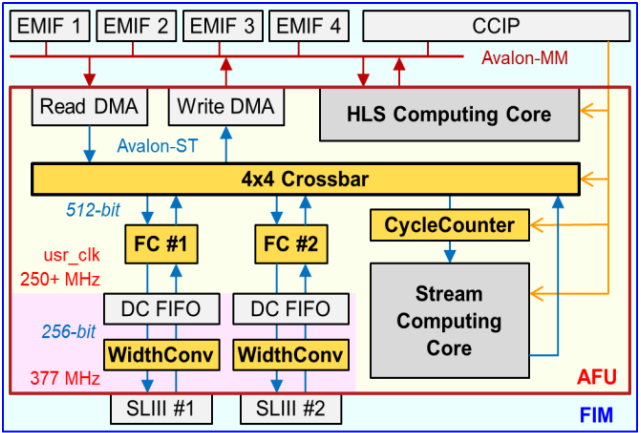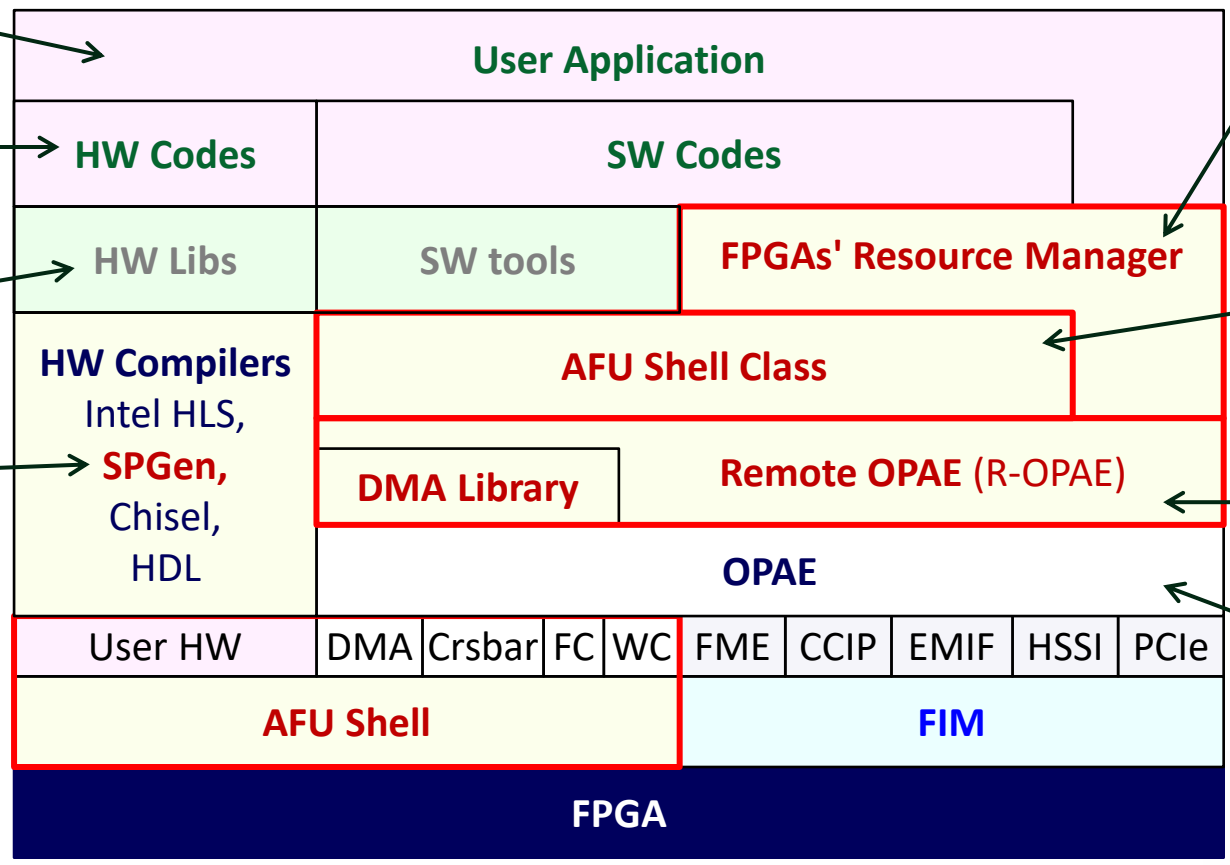
WRC2024: FPGA or CGRA

# System Stack of ESSPER

**High-level application**
- Including both SW & HW

**Low-level application**
- Separately-written HW & SW

**Hardware, software libs/tools**
- Pre-implemented functions

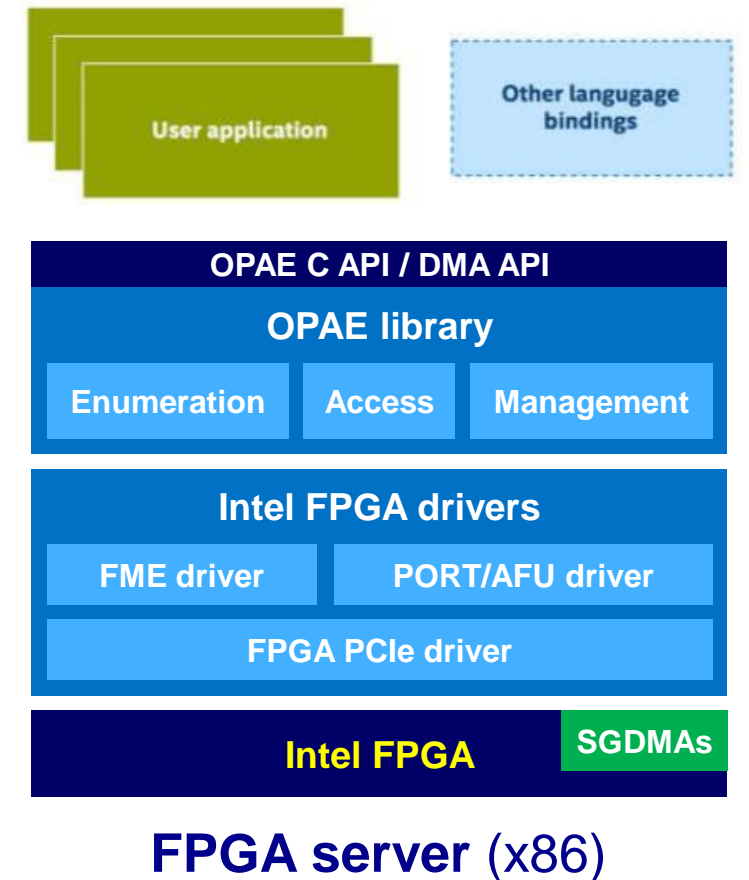**Programming tools**
- SPGen : DSL for stream/ systolic computing

**Resource manager**
- Search and allocate resources of multiple FPGAs
- FPGA network management / control

**AFU Shell class**
- Object of AFU shell
- Abstraction of HW

**Remote OPAE**
- Software bridge using Infiniband Verbs

**OPAE**
- Low-level driver



| User Application | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **HW Codes** | **SW Codes** | | | | | | | |
| HW Libs | SW tools | **FPGAs' Resource Manager** | | | | | | |
| **HW Compilers** Intel HLS, **SPGen,** Chisel, HDL | **AFU Shell Class** | | | | | | | |
| | **DMA Library** | **Remote OPAE** (R-OPAE) | | | | | | |
| | OPAE | | | | | | | |
| User HW | DMA | Crsbar | FC | WC | FME | CCIP | EMIF | HSSI | PCIe |
| **AFU Shell** | | | | | **FIM** | | | |
| **FPGA** | | | | | | | | |

# Remote-OPAE (for remote FPGA Access)
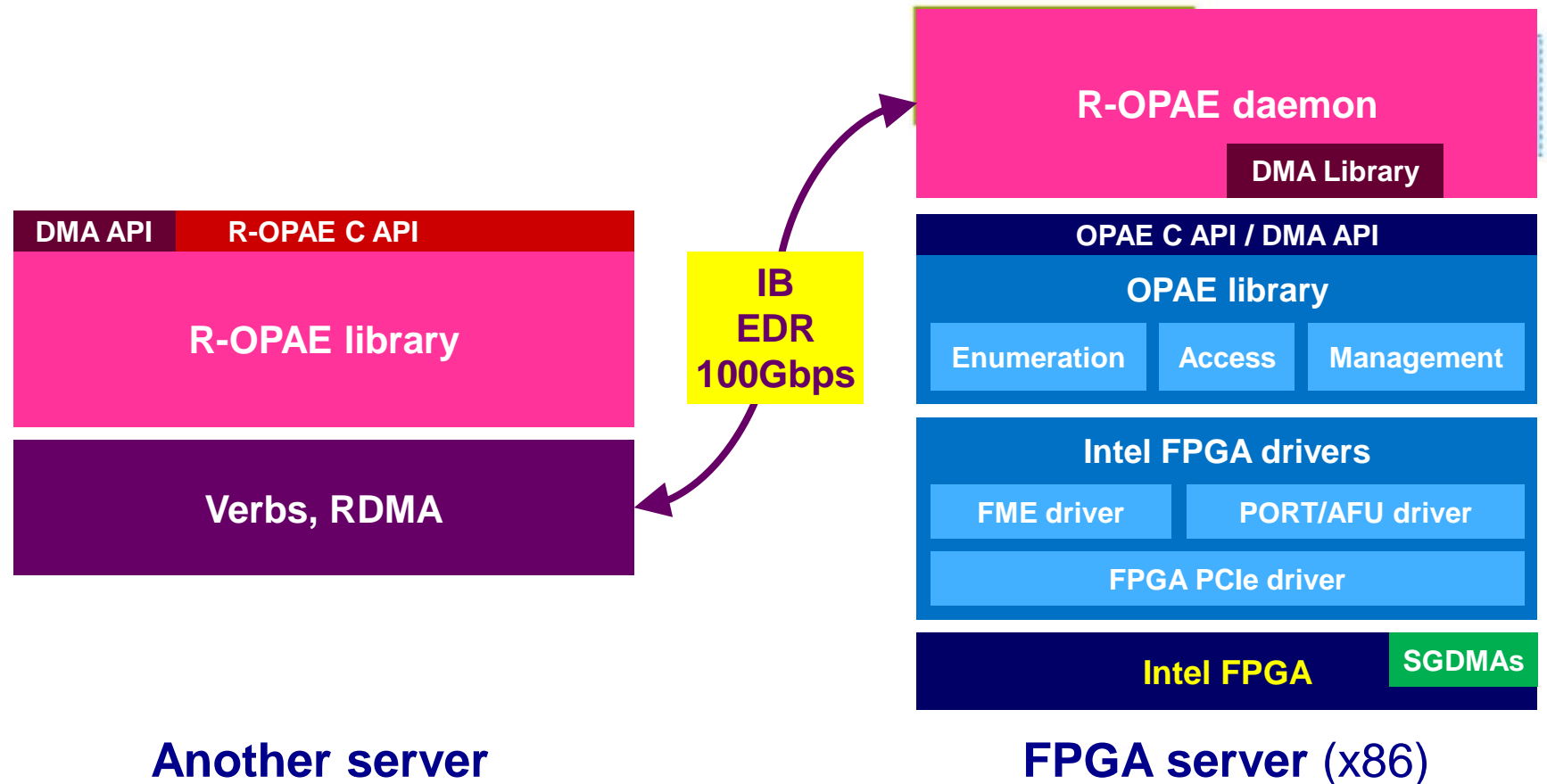
## Software bridge for FPGAs over Infiniband

✓ **OPAE**: Open Programmable Acceleration Engine
(PCIe FPGA driver)



**FPGA server** (x86)

# Remote-OPAE (for remote FPGA Access)

## Software bridge for FPGAs over Infiniband

- ✓ **OPAE**: Open Programmable Acceleration Engine (PCIe FPGA driver)

- ✓ 99% of OPAE APIs are supported.

- ✓ We can use any FPGAs in a system via IB as if they were locally installed.

| DMA API | R-OPAE C API |
| --- | --- |

**R-OPAE library**

**Verbs, RDMA**

**Another server**

**IB EDR 100Gbps**

**R-OPAE daemon**

**DMA Library**

**OPAE C API / DMA API**

**OPAE library**

| Enumeration | Access | Management |
| --- | --- | --- |

**Intel FPGA drivers**

| FME driver | PORT/AFU driver |
| --- | --- |

**FPGA PCIe driver**

**Intel FPGA** — **SGDMAs**

**FPGA server** (x86)
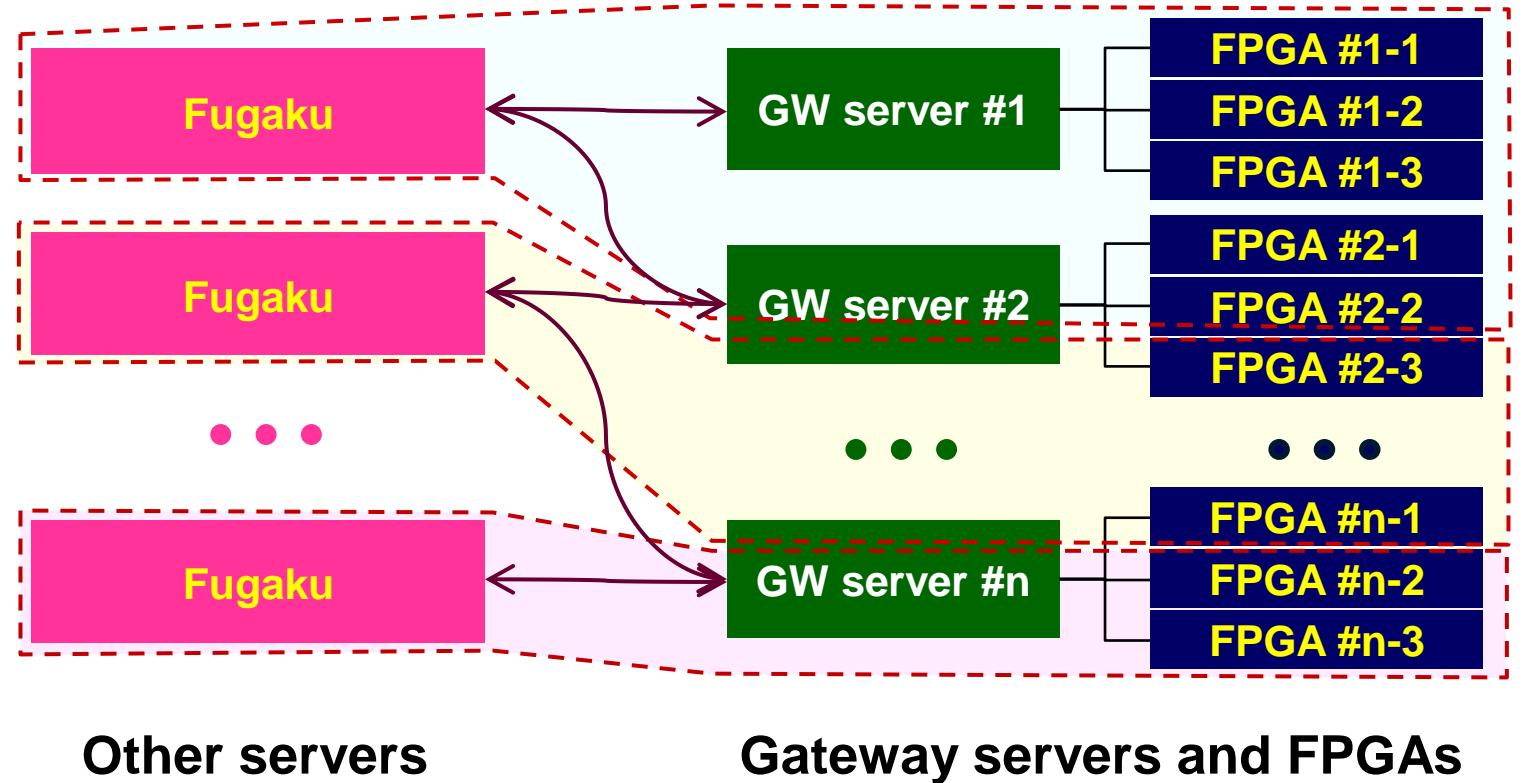
# R-OPAE as Software-based Resource Disaggregation

**Transparent access to remote FPGAs**

**Flexible utilization:**

✓ Can use any available FPGA resources

**Inter-operability and extensibility:**

✓ Vendor/ISA-independent

✓ Operable with various architectures such as Fugaku (ARM)
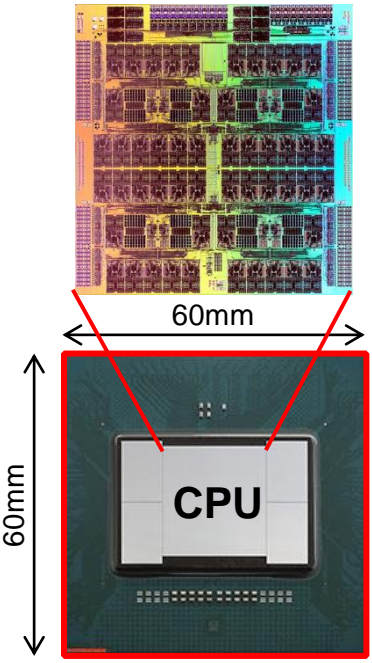


**Other servers**                    **Gateway servers and FPGAs**

ESSPER
Elastic and Scalable System for High-Performance Re-configurable Computing
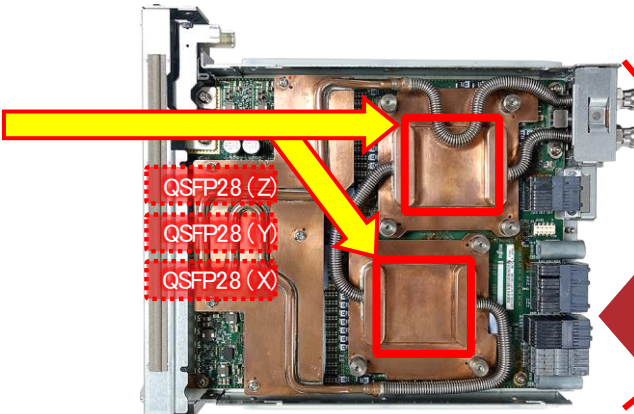
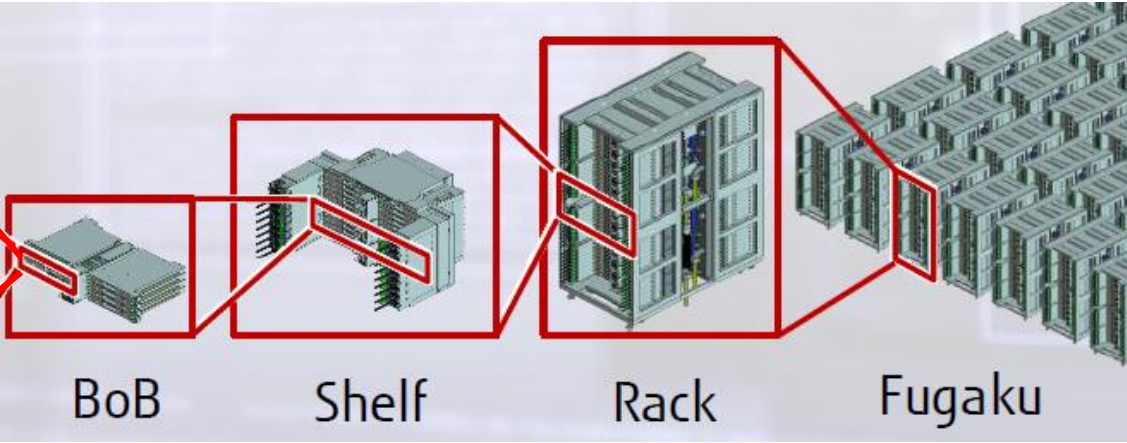# Proof-of-concept and evaluation

# Supercomputer Fugaku



**Modern Supercomputers are based on Many-core CPUs (& GPUs).**

60mm

60mm

**CPU**

48+ cores / 1 node
2.7+ TF

QSFP28 (Z)
QSFP28 (Y)
QSFP28 (X)

**CPU-Memory Unit (CMU)**
2 nodes
5.4+ TF

| BoB | Shelf | Rack | Fugaku |
|---|---|---|---|
| 16 nodes 43+ TF | 48 nodes 129+ TF | 384 nodes 1+ PF | 158,976 nodes 537 PF @ FP64 (414 racks) |

Photos & figs by Fujitsu

# Elastic and Scalable System for High-PErformance Reconfigurable Computing

**Experimental prototype for research on functional extension with FPGAs**



ESSPER

Supercomputer Fugaku

Connected w/ 100m cables

# Elastic and Scalable System for High-PErformance Reconfigurable Computing
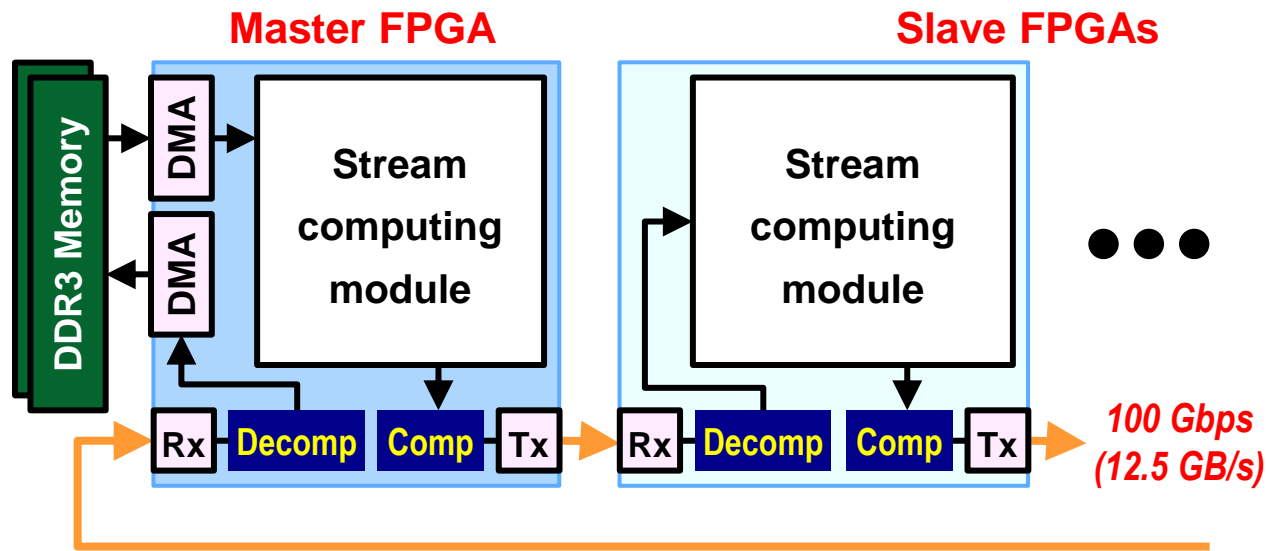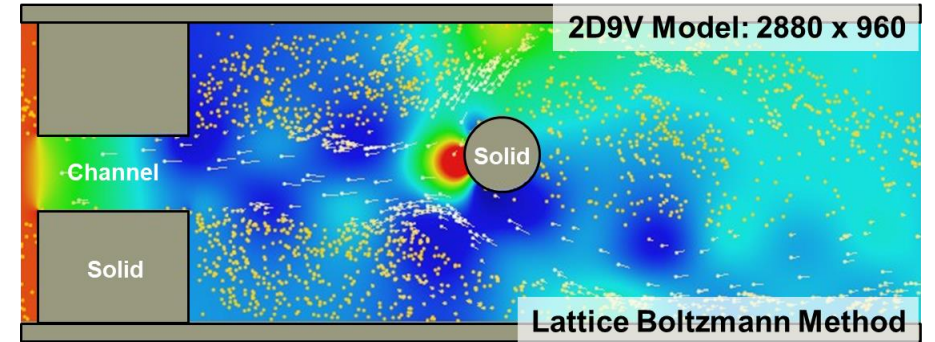
Open-Access paper

**ESSPER**

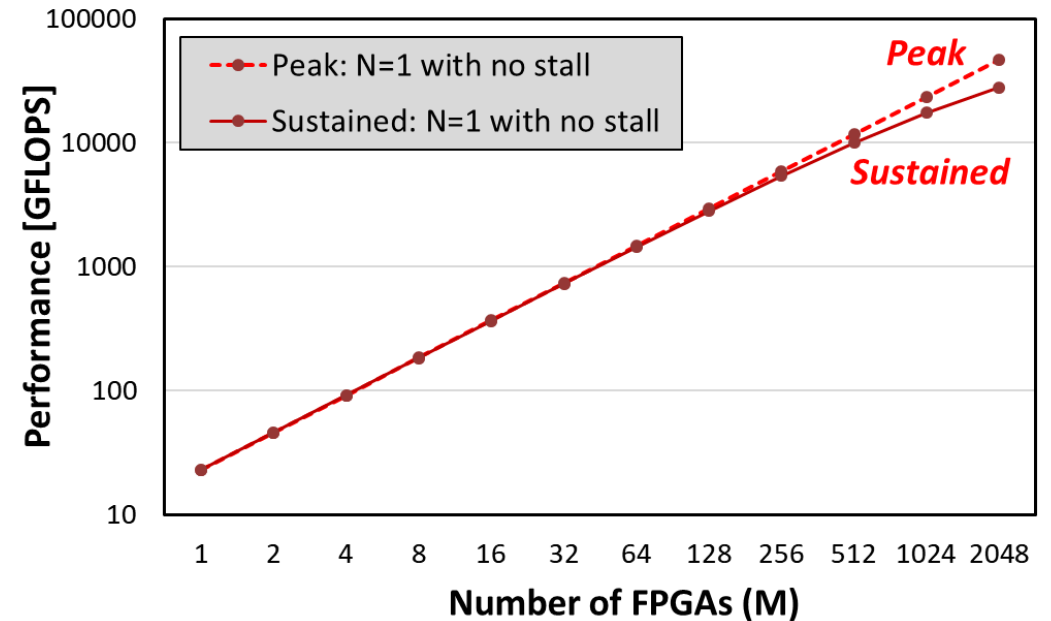Elastic and Scalable System for High-Performance Re-configurable Computing

# Applications, Joint Research Projects

# Ringed FPGAs for Deeper Pipelining

- **Deeply-pipelined FPGAs with 1D ring**
  - ✓ Linear array of Stratix10 FPGAs
  - ✓ Pipelining works well for almost linear speedup if data stream is sufficiently large.
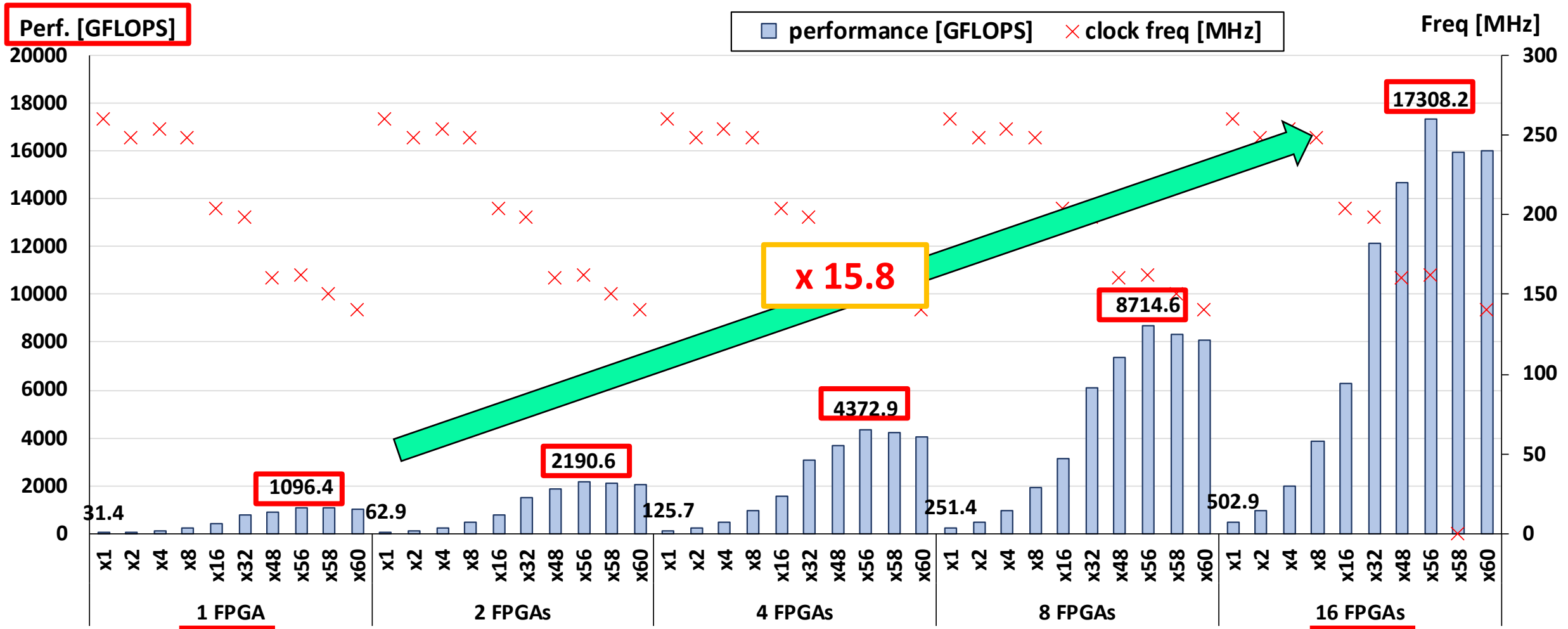


2D9V Model: 2880 x 960

Channel

Solid

Solid

**Lattice Boltzmann Method**



**Master FPGA**

**Slave FPGAs**

DDR3 Memory — DMA — Stream computing module — Rx — Decomp — Comp — Tx

Rx — Decomp — Comp — Tx

*100 Gbps (12.5 GB/s)*

**Block diagram of FPGAs in a ring**



Peak: N=1 with no stall
Sustained: N=1 with no stall

*Peak*

*Sustained*

Performance [GFLOPS]

Number of FPGAs (M)

**Performance model for Arria10 FPGAs**

# Performance of 2D LBM with 100Gbps Ring NW

Computational performance (FLOPS) when processing about 2GB data

WRC2024: FPGA or CGRA
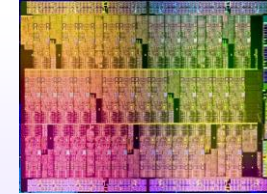Jan 17, 2024
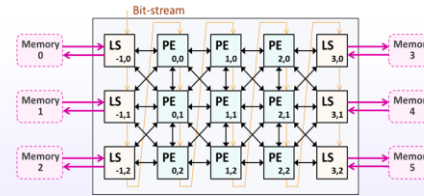
# Lessons Learned with ESSPER

Open-Access paper

- **FPGA-based reconfigurable computing works.**

- **Productivity is not high, especially for multiple FPGAs.**
  - ✓ Even HLS requires know-how on optimizing computation and memory access.
  - ✓ Lack of debugging tool, and simulation environment.

- **Can obtain scalability, but**

  **absolute performance in FP is lower than competitors (GPUs)**
  - ✓ FPGA-bases system development takes time while GPUs are being further advanced.
  - ✓ For fixed domain of computing *(such as HPC in FP and AI workloads)*, FPGAs are redundant with more area, more power, and lower frequency with low memory bw.

- **Concept should be Okay for reconfigurable data-flow computing, but implementation approach could be improved :** CGRA instead of FPGA?

**2. Exploration of New HPC Architectures**

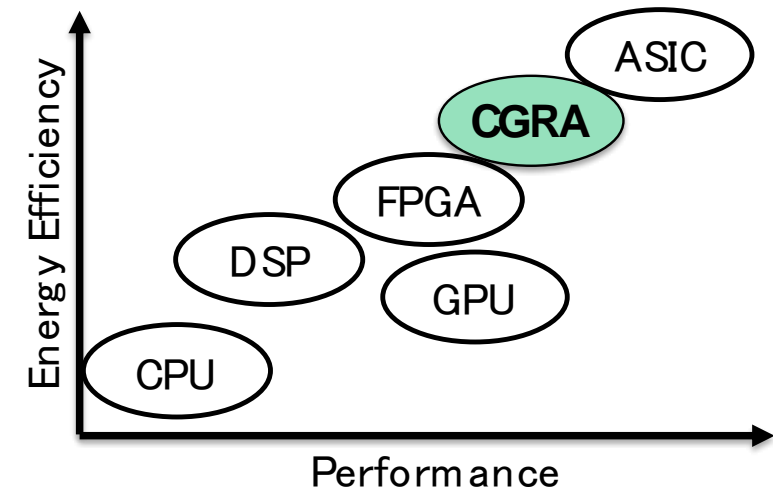- ✓ Data-flow-based accelerators (**CGRA**)
- ✓ **Next-generation HPC systems**
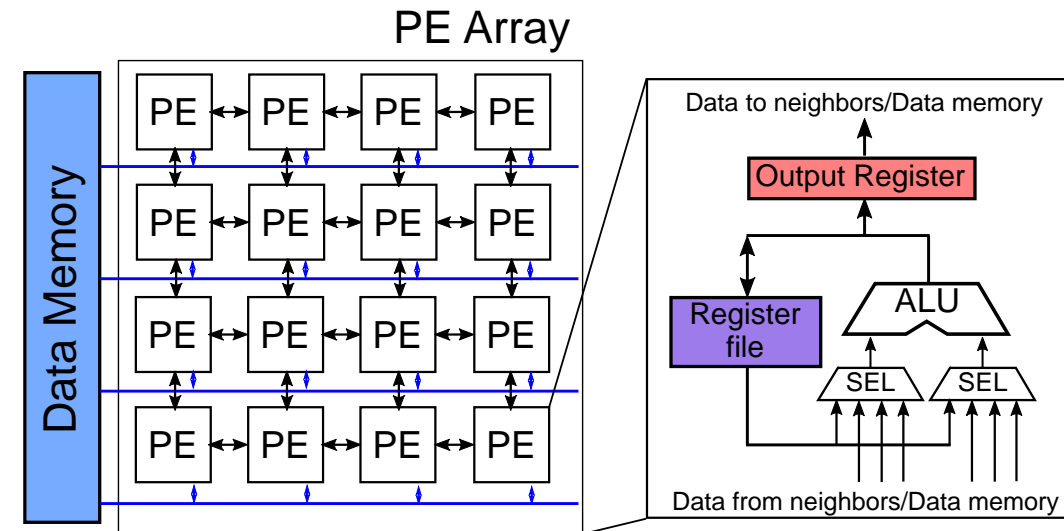
# Exploration of New HPC Architectures

Data-flow-based accelerators (CGRA)

# Coarse-Grained Reconfigurable Array (CGRA)

- **Architecture for reconf. data-flow computing**
  - ✓ Composed of an array of processing elements (PEs), where we can map DFGs for computing
  - ✓ Provide a word-wise reconfigurability (e.g., 32-bit)
  - ✓ Higher energy efficiency than FPGAs (of bit-level)
  - ✓ Performance close to ASIC-based accelerators

- **Application area of CGRAs**
  - ✓ Traditionally, targeted for lower-power embedded apps, e.g., image processing
  - ✓ Recently, expected for hi-performance AI

- **Questions**
  - ✓ CGRAs also promising for HPC?
  - ✓ What architecture/design decision required HPC?



**Comparison with other architectures [1]**

PE Array
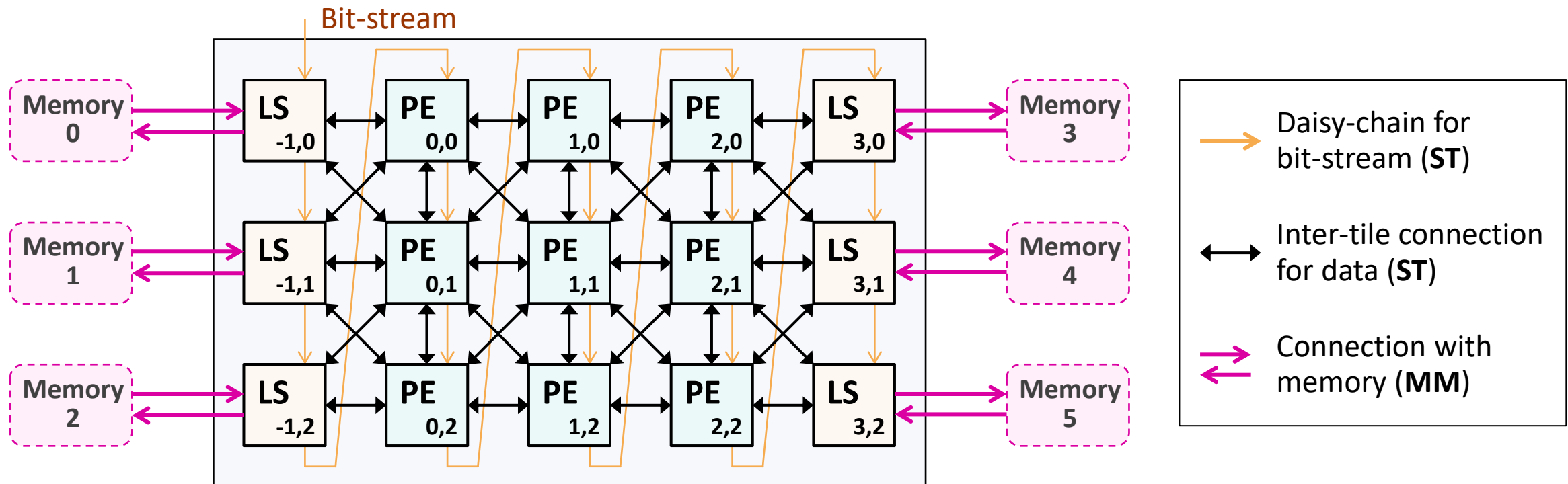


**General structure of the CGRAs [2]**

[1] Liu, Leibo, et al. "A survey of coarse-grained reconfigurable architecture and design: Taxonomy, challenges, and applications." *ACM Computing Surveys (CSUR)* 52.6 (2019): 1-39.
[2] Takuya Kojima, et al., "Exploration Framework for Synthesizable CGRAs Targeting HPC: Initial Design and Evaluation," Procs. CGRA4HPC, May 30-June 3, 2022.

# RIKEN CGRA Architecture (baseline)

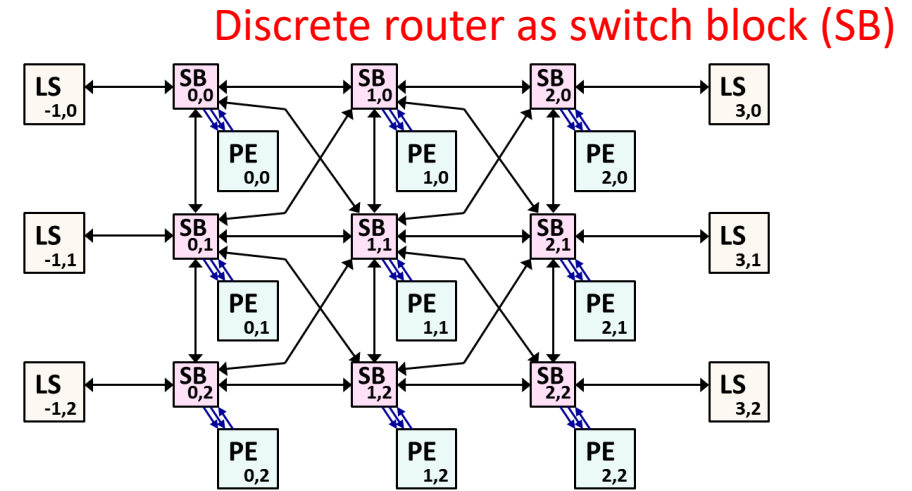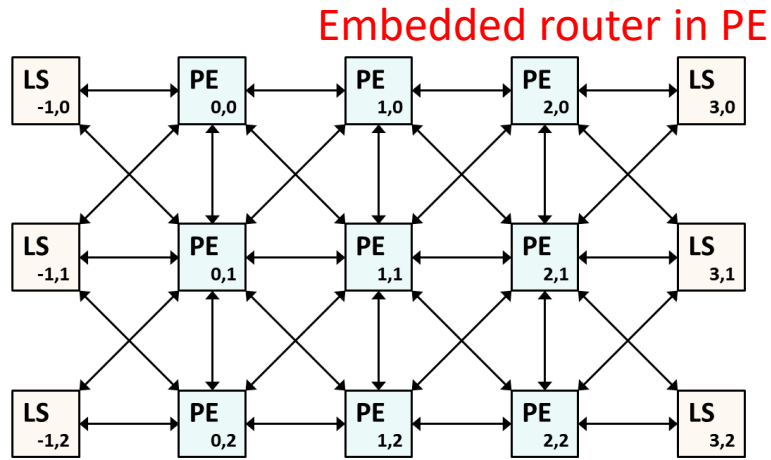- **HPC-oriented CGRA with the following design philosophy**
  - ✓ Modular design for design space exploration with various architecture configuration and sizes
  - ✓ Isolation between computation in a PE array and memory access with load-and-store (LS) tiles
  - ✓ Capability of floating-point operations for HPC apps

# Past work:
# CGRA Designs with Embedded routers (ER) or Discrete routers (DR)

# Design Decision on Intra-CGRA Interconnects

Embedded router in PE

Discrete router as switch block (SB)



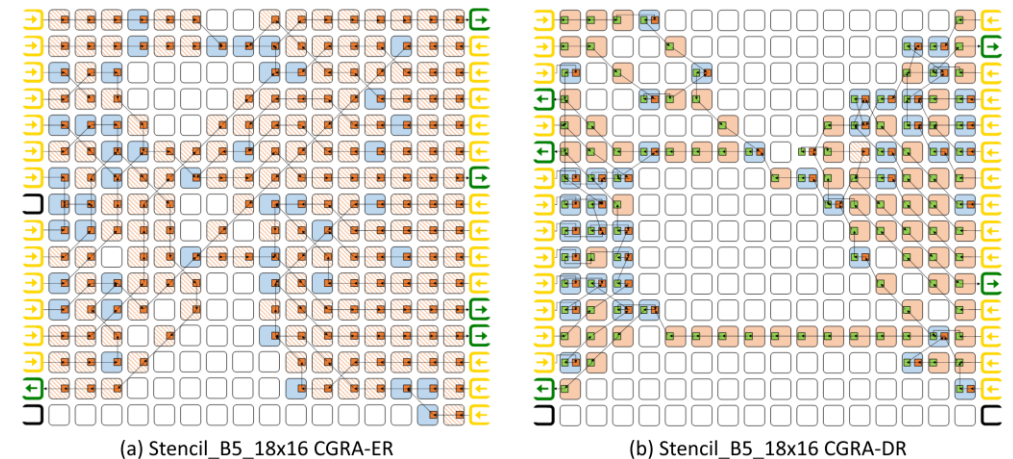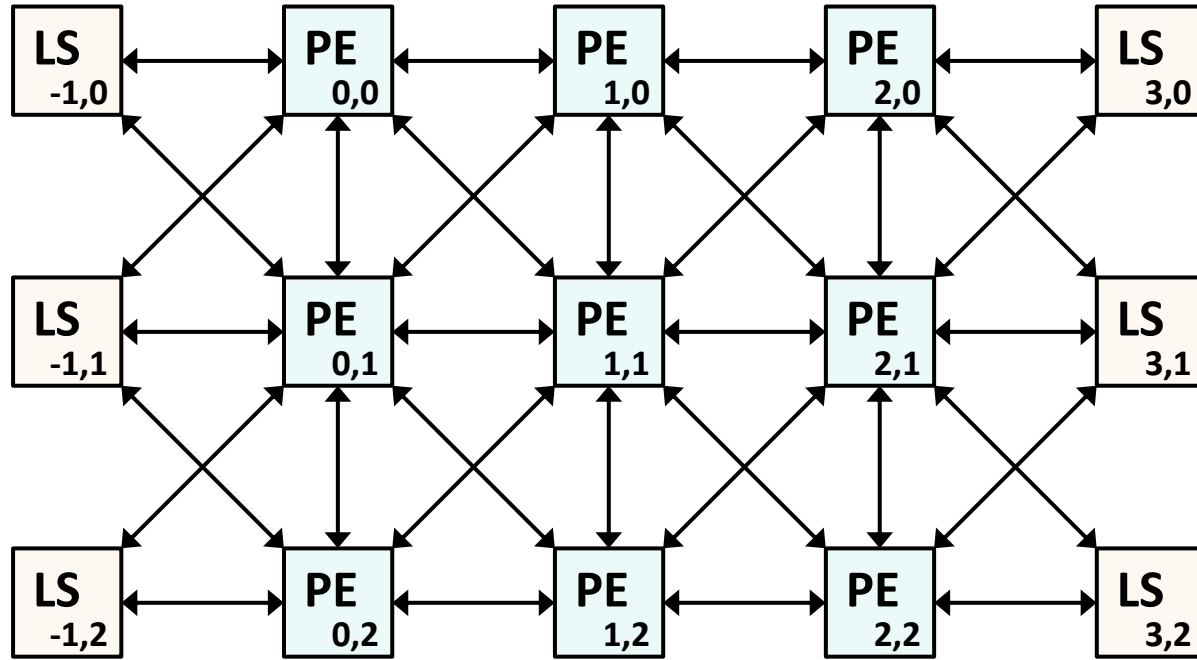| CGRA with Embedded Router (CGRA-ER) | CGRA with Discrete Router (CGRA-DR) |
|---|---|
| Routers in each PE mediate communication between PEs | Discrete switch blocks for communication between PEs |
| Simpler in design with smaller area and higher frequency | More complex design with large area and lower frequency |
| May lead to wastage of computing resources: When mapping complex graphs, some PEs must be configured to bypass data without computation | More efficient: allow more PEs to be used for computing: By relieving PEs from compulsory data bypass with better routability with switch blocks |
| Ex) ADRES, CGRA-ME, HiPreP, and MorphoSys | Ex) HyCube, RAW (& the Tilera Processor), and Plasticine |

## Which one is suitable for HPC considering computing efficiency and HW area?

# Objective and Contributions

- **Objective**    **Evaluate the positive and negative aspects of the different routing architectures: CGRA-ER and CGRA-DR for HPC apps**

- **Contributions**

  ✓ Parameterized implementation of CGRA-ER and CGRA-DR using our baseline architecture

  ✓ Evaluation and verification of benchmarks (DFGs) by RTL simulation
    with CGRA-evaluation framework (our previous work)

  ✓ Comparison between CGRA-ER and CGRA-DR

  - ➤ Difficulty in place-and-route of DFGs
  - ➤ PE utilization for computation
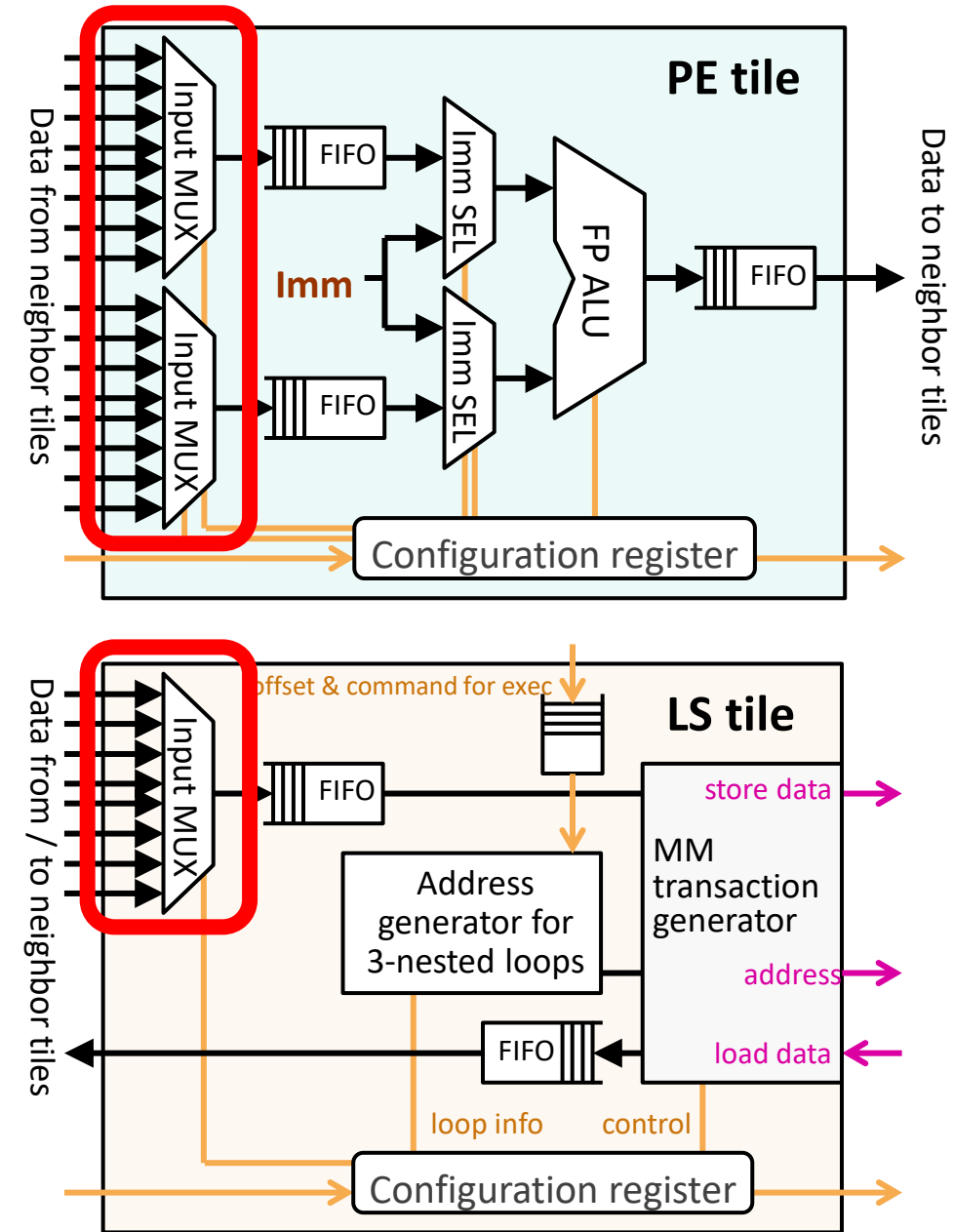  - ➤ Hardware resource consumption

(a) Stencil_B5_18x16 CGRA-ER        (b) Stencil_B5_18x16 CGRA-DR

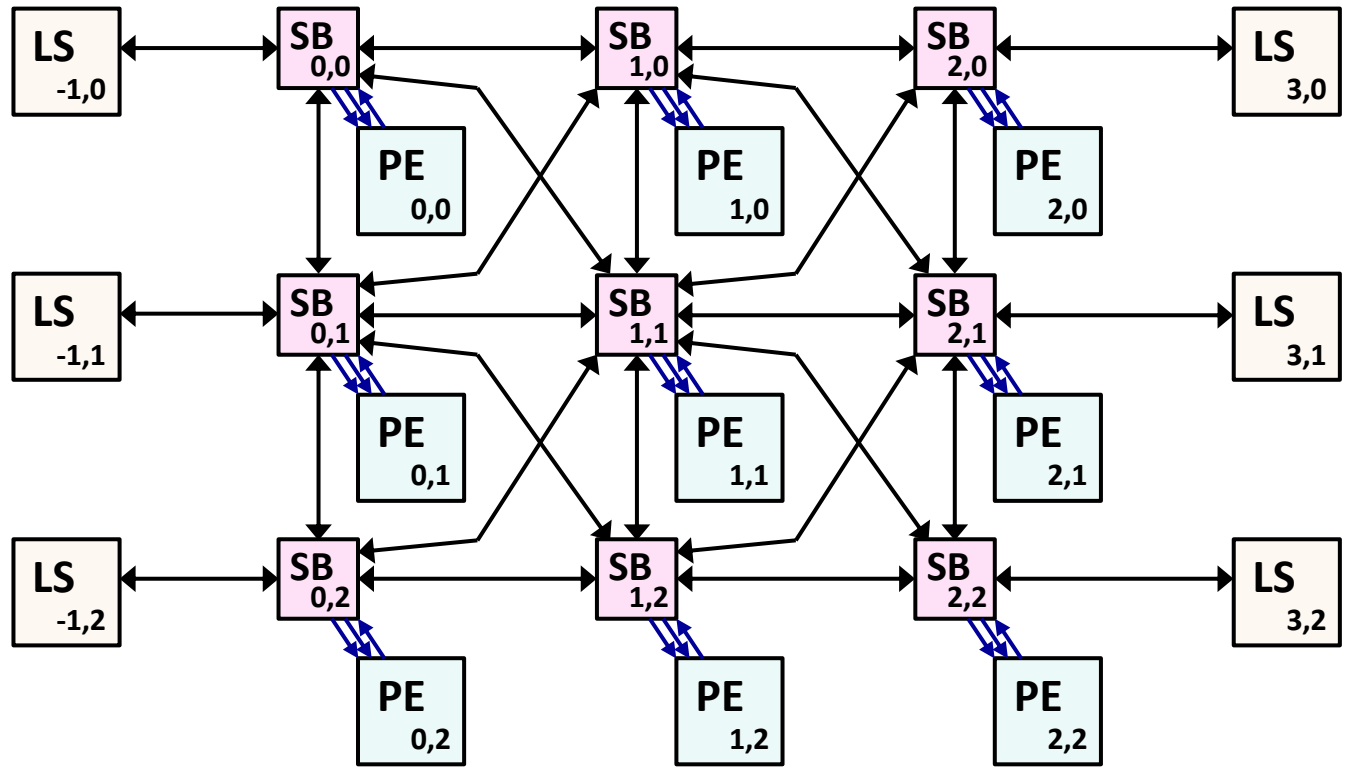# CGRA with Embedded Router (CGRA-ER)

**Embedded multiplexors**



- ✓ PEs and LSs are directly connected to each other.
- ✓ Limited routing capability causes inefficient PE utilization.
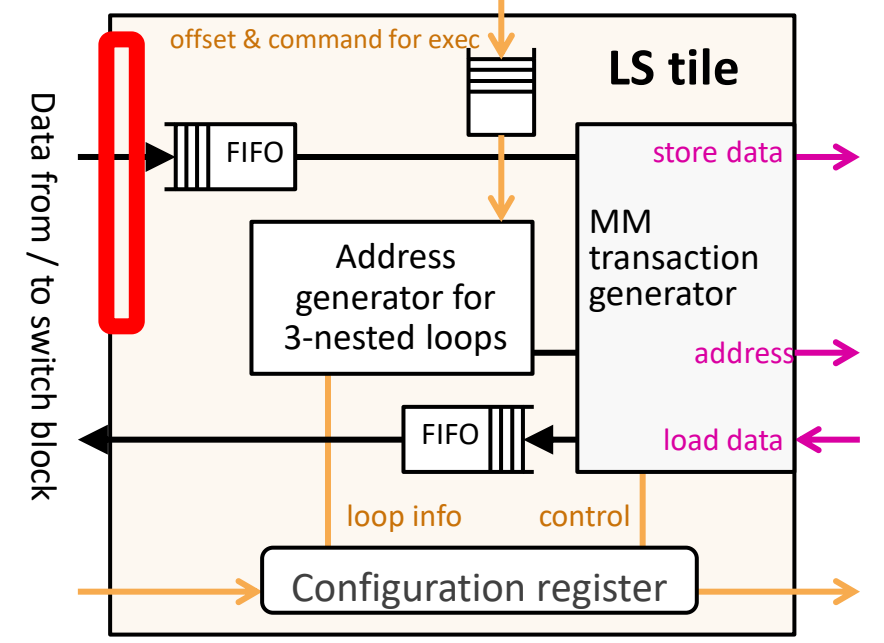- ✓ Wasted PEs for routing with NOP to bypass data through (No computing)

# CGRA with Discrete Router (CGRA-DR)
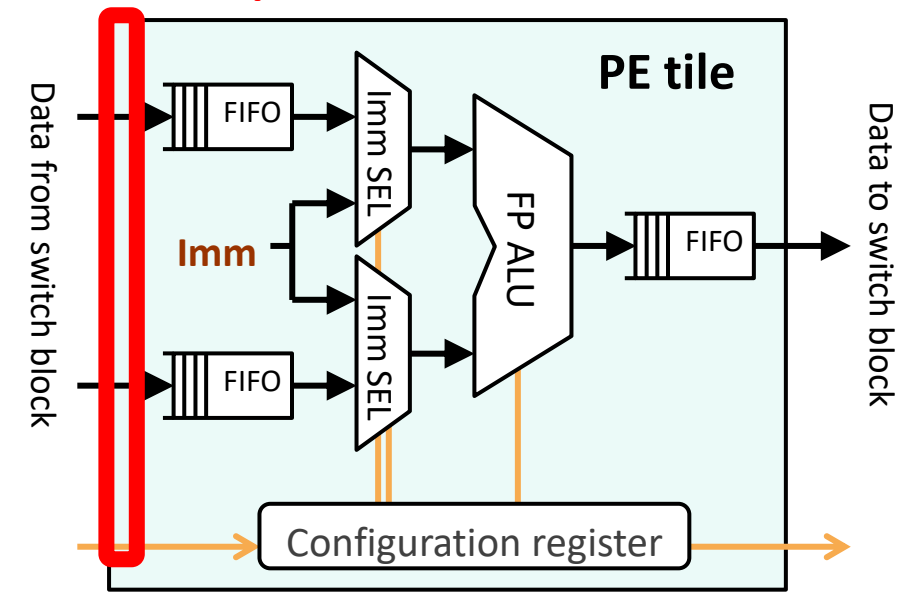
**No embedded multiplexor**
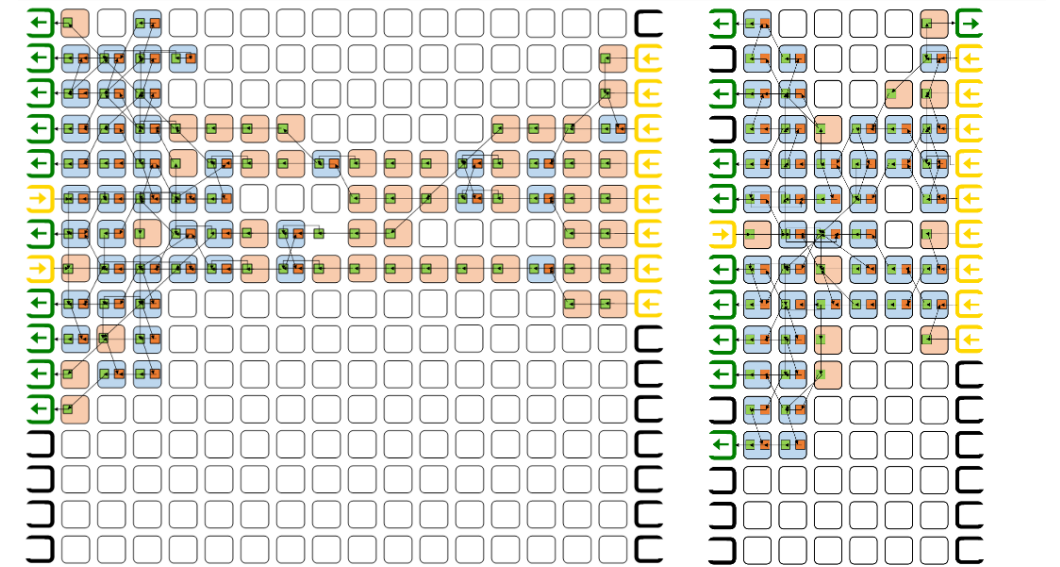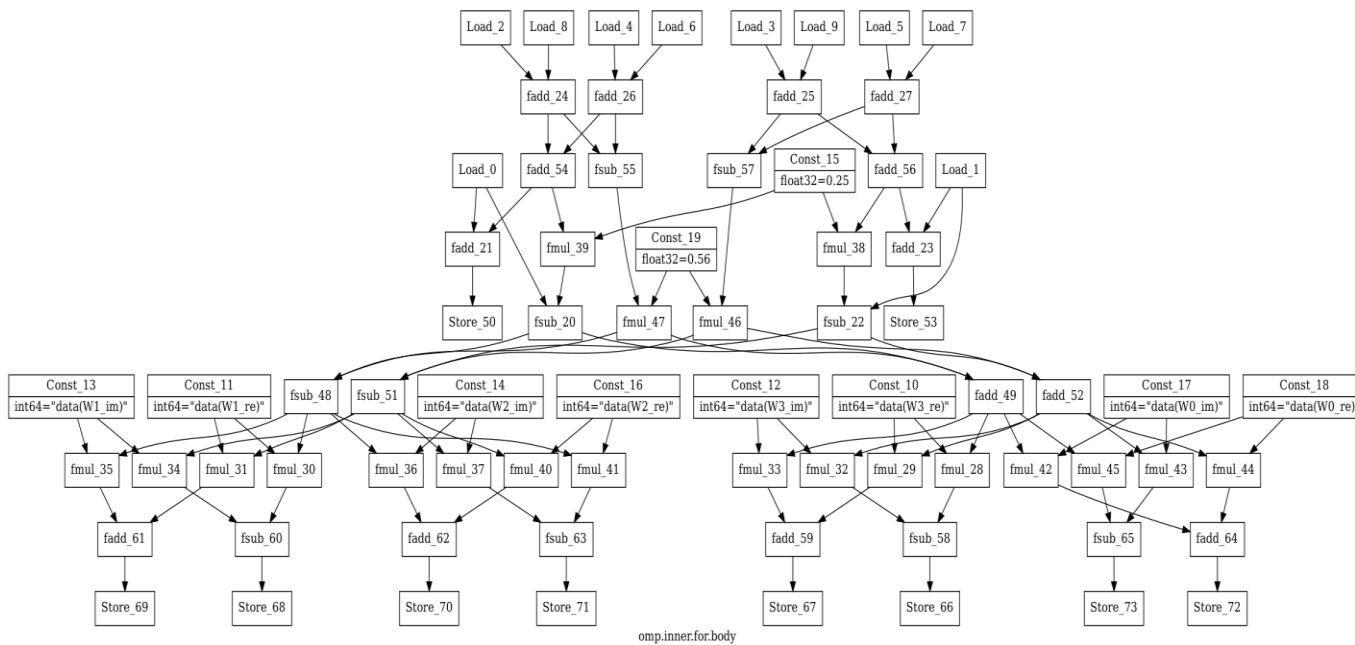


- ✓ Switch blocks as discrete router
- ✓ Simpler PE and LS tile with no multiplexor
- ✓ Higher hardware resources required for SBs

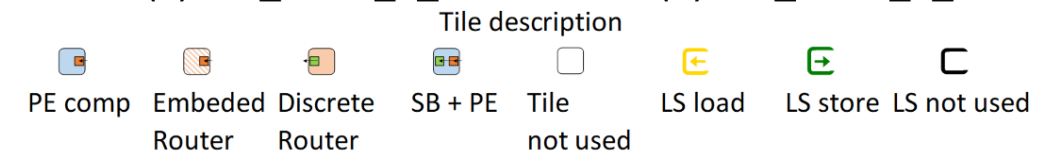# Routing Flexibility vs. Complex Kernel

- **Stockham Radix-5 FFT Kernel DFG has <span style="color:red">more edges connecting among nodes.</span>**
  - ✓ **<span style="color:red">Could not map onto CGRA-ER</span>** ; routing with PEs is not enough.
  - ✓ CGRA-DR allows it to be mapped onto 18×16 CGRA-DR and even smaller one (8×16).



**DFG of the Innermost Loop of Stockham Radix 5 FFT kernel**

(a) FFT_Radix_5_18x16     (b) FFT_Radix_5_8x16

Tile description

PE comp | Embeded Router | Discrete Router | SB + PE | Tile not used | LS load | LS store | LS not used

**Stockham Radix 5 FFT kernel mapped on CGRA-DR**

# Future Work on CGRA for HPC
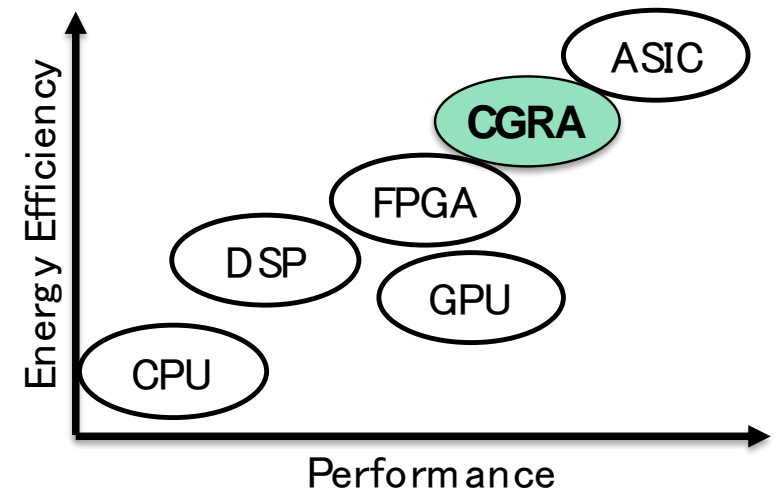


**Comparison with other architectures [1]**

- **Have not obtained conclusions yet.**

- **Extension of the baseline architecture for practical kernels**

  ✓ Predication for conditional execution

  ✓ Heterogeneous architecture with FP div, sqrt, log, and transcendental functions to cover wider range of applications

  ✓ Programmable buffer for data reuse (such like line/stencil buffer)

- **Evaluation of operation frequency and hardware resource consumption with FPGA-based and/or ASIC implementation**

  ✓ Initial rough evaluation results for ASIC

  ✓ FPGA implementation is also on-going.

    ➢ Supercomputer Fugaku – CGRA emulation on ESSPER.

# Summary

**Reconfigurable data-flow computing** should be promising for power-efficient HPC.

- ✓ ***FPGA-based HPC testbed; ESSPER*** *(prototype FPGA cluster)*
  - ➢ *Stratix 10 FPGAs*
  - ➢ *FPGA Shell with inter-FPGA network*
- ✓ ***RIKEN CGRA for HPC***
  - ➢ *Baseline architecture*
  - ➢ *Design space exploration for inter-tile connection*

## Future work

- ✓ ESSPER2 with Intel Agilex-M FPGA
- ✓ CGRA for HPC and AI, (design for ASIC and compiler)
- ✓ Feasibility study for next-gen supercomputers (conducted)