

DESIGN, AUTOMATION AND TEST IN EUROPE

THE EUROPEAN EVENT FOR ELECTRONIC SYSTEM DESIGN & TEST

31 MARCH – 2 APRIL 2025 LYON, FRANCE

CENTRE DE CONGRÈS DE LYON



Multi-Mode Borderguard Controllers for Efficient On-Chip Communication in Heterogeneous Digital/Analog Neural Processing Units

Hong Pang^{*}, Carmine Cappetta[†], Riccardo Massa[‡], Athanasios Vasilopoulos[§], Elena Ferro^{*}[§], Gamze Islamoglu^{*}, Angelo Garofalo^{*}¶, Francesco Conti[¶], Luca Benini^{*}¶, Irem Boybat[§], Thomas Boesch[◊]

*ETH Zurich, Switzerland [†]STMicroelectronics, Cornaredo, Italy [‡]STMicroelectronics, Agrate, Italy [§]IBM Research - Zurich, Switzerland [¶]University of Bologna, Italy [§]STMicroelectronics, Geneva, Switzerland



Innovate UK Schweizerische Eidgenossenschaft Confédération suisse Confederazione Svizzera Confederaziun svizra









ALMA MATER STUDIORUM Università di Bologna

A powerful solution for edge Al...

Performs MAC (dot product) operation with high energy efficiency (> 90% of overall operations in Al inference)

AINC (NVM-based analog inmemorycomputing)

Non-volatile memory-based AIMC is particularly attractive for:

- high on-chip weight storage capacity
- high power efficiency

We also need other digital units...



AINC (NVM-based analog inmemorycomputing) Performs MAC (dot product) operation with high energy efficiency (> 90% of overall operations in Al inference)

[2] S. Jain et al., "A Heterogeneous and Programmable Compute-In-Memory Accelerator Architecture for Analog-Al Using Dense 2-D Mesh," in IEEE Transactions on Very Large Scale Integration (VLSI) Systems, vol. 31, no. 1, pp. 114-127, Jan. 2023.

We also need other digital units...

CPU

General-Purpose Core

(Accuracy critical OP, activation function not supported by DPU)

Digital Processing Unit

(Convolutions, pooling, arithmetic, and activation functions) DPU

AIMC

(NVM-based analog inmemorycomputing) Performs MAC (dot product) operation with high energy efficiency (> 90% of overall operations in Al inference)

Local Storage Unit

SRAM

[2] S. Jain et al., "A Heterogeneous and Programmable Compute-In-Memory Accelerator Architecture for Analog-Al Using Dense 2-D Mesh," in IEEE Transactions on Very Large Scale Integration (VLSI) Systems, vol. 31, no. 1, pp. 114-127, Jan. 2023.

We also need other digital units...

CPU

General-Purpose Core

(Accuracy critical OP, activation function not supported by DPU)

Digital Processing Unit

(Convolutions, pooling, arithmetic, and activation functions)

Local Storage Unit

DPU

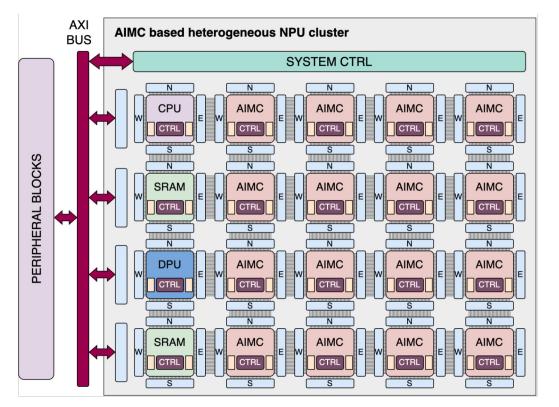
AIMC

(NVM-based analog inmemorycomputing)

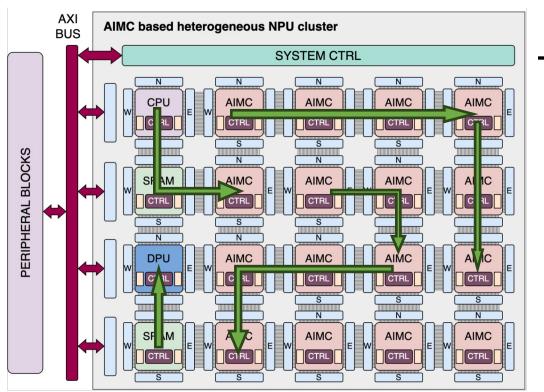
Performs MAC (dot product) operation with high energy efficiency (> 90% of overall operations in Al inference)

[2] S. Jain et al., "A Heterogeneous and Programmable Compute-In-Memory Accelerator Architecture for Analog-AI Using Dense 2-D Mesh," in IEEE Transactions on Very Large Scale Integration (VLSI) Systems, vol. 31. no. 1. pp. 114-127. Jan. 2023. Hong Pang / ETH Zurich

SRAM



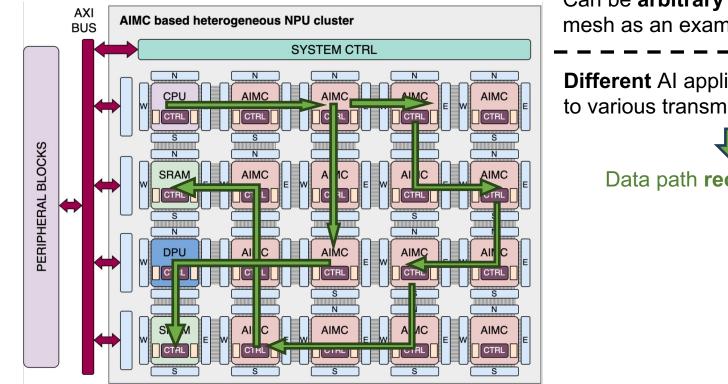
Can be **arbitrary** topology, we take mesh as an example



Can be **arbitrary** topology, we take mesh as an example

Different AI application mapping leads to various transmission patterns

Application 1 Data Flow

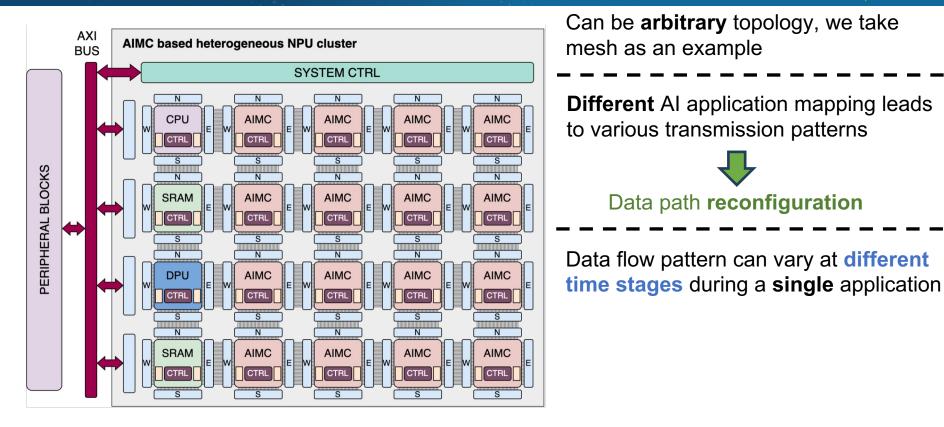


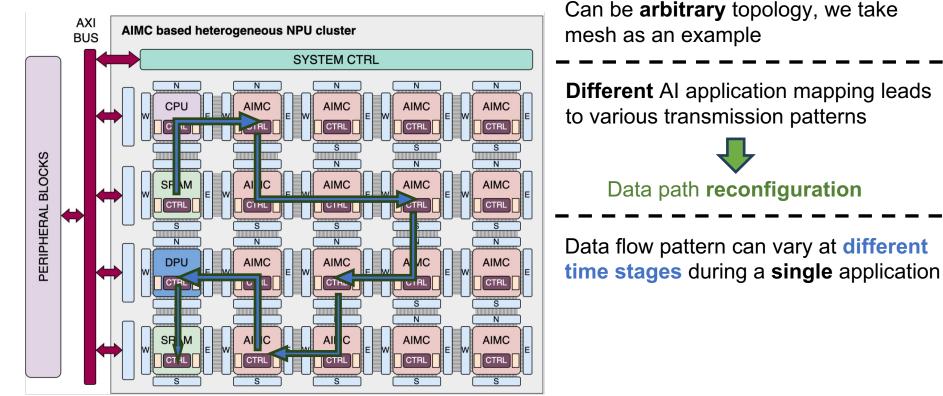
Can be **arbitrary** topology, we take mesh as an example

Different AI application mapping leads to various transmission patterns



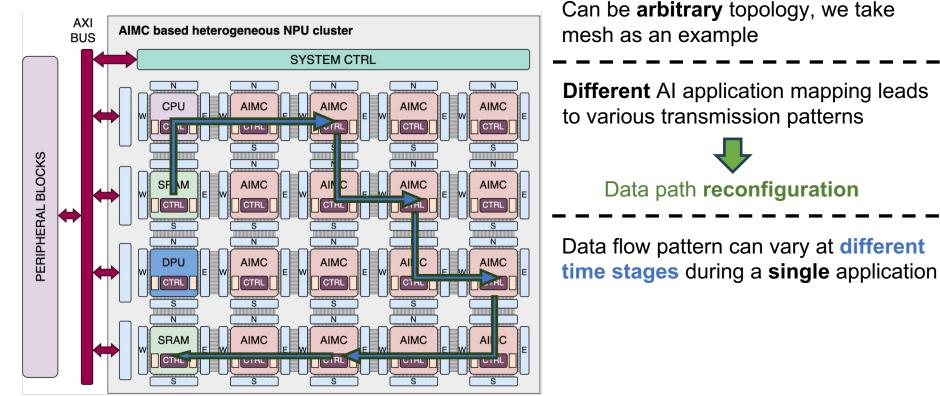
Application 2 Data Flow





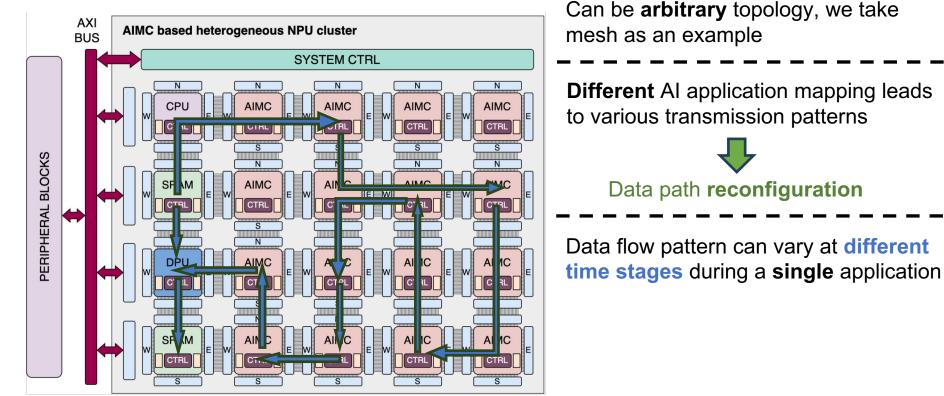
A single complex application Time stage 1

Hong Pang / ETH Zurich



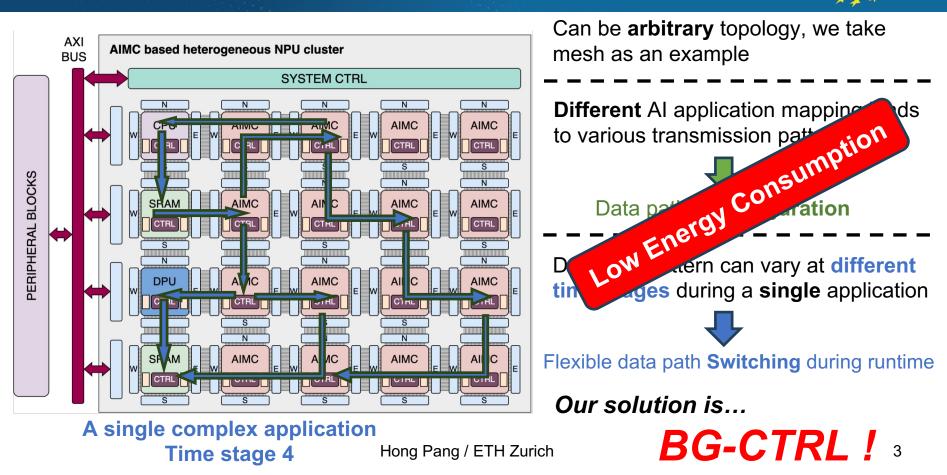
A single complex application Time stage 2

Hong Pang / ETH Zurich



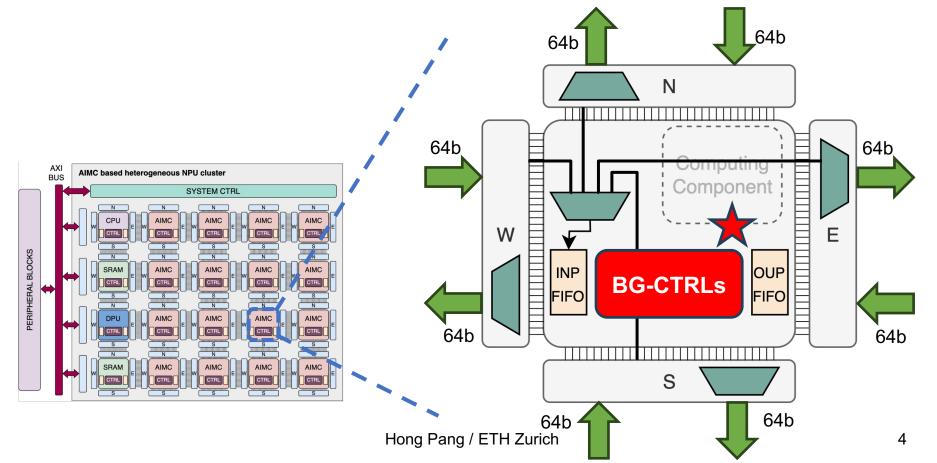
A single complex application Time stage 3

Hong Pang / ETH Zurich



NPU Node Architecture







Generality

The interconnect can work on various layer mapping and data flow scenarios.

High routing flexibility

Each node can be used for multiple data paths during a single application.

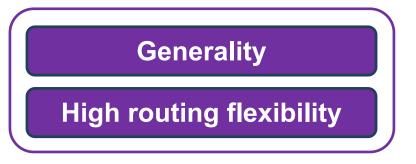
Synchronization

All nodes within the cluster should transmit data following dependency.

High Energy Efficiency

Challenges we have...



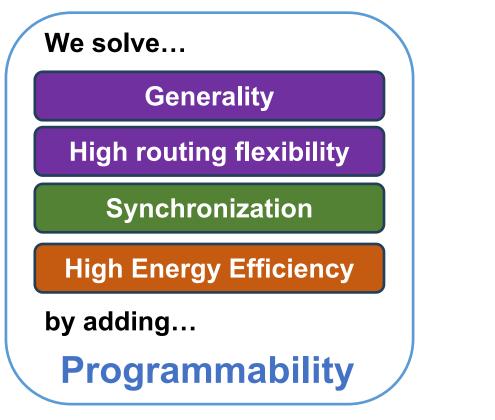


Fine-grained data path reconfiguration

Synchronization

A global reference timestamp telling when to transmit data

High Energy Efficiency



We solve...



Synchronization

High Energy Efficiency

by adding...

Programmability

We design a specialized instruction set...

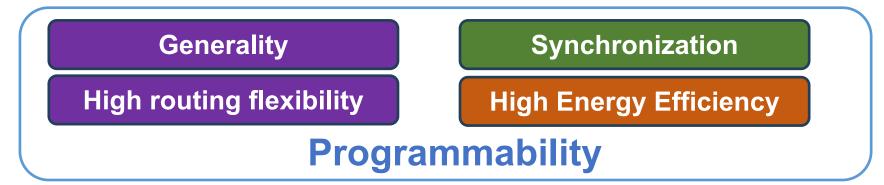
Instruction tells BG-CTRL...

- When to transmit data from where to where
- Loop repeated data transmission pattern

This enables...

- Cycle-wise data path formation/switching
- Precise timestamp of data transmission
- Less energy consumption on instr. fetching

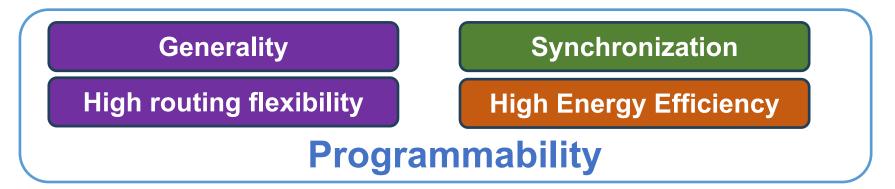




High Energy Efficiency

Multiple Operating Modes





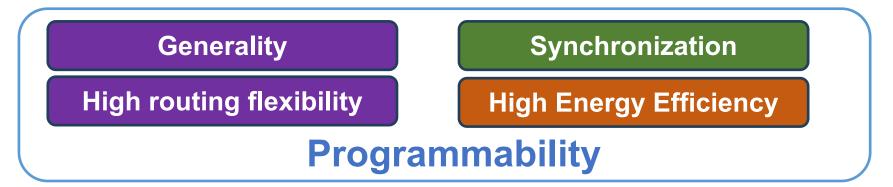
High Energy Efficiency

Multiple Operating Modes
Fixed Configuration

Form data path via configuration registers

- **Cannot** reconfigure data path
- Data is transmitted whenever available
- Higher energy efficiency than program
- Recommended for simple data flow





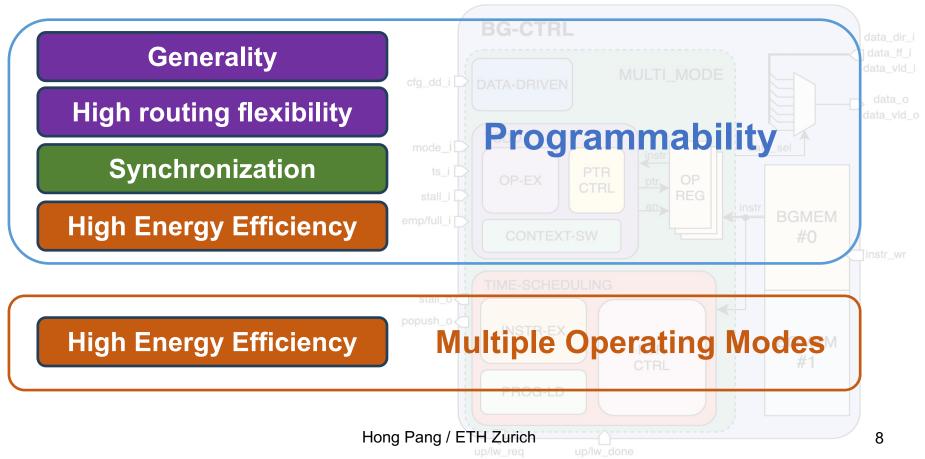
High Energy Efficiency

Multiple Operating Modes
Time-sliced Configuration

Extend fixed configuration with several **data path sets**, each activated in certain period

- More flexible than fixed configuration
- Higher energy efficiency than program
- Recommended for data flow with regular path switching





Hong Pang / ETH Zurich



STM FD-SOI 28nm for Evaluation

We integrate it into node and build a 3x3 mesh...

BG-CTRL data dir i data ff i data_vld_i MULTI MODE 640 MHz cfg_dd_i DATA-DRIVEN data o data_vld_o **@SS Corner** TSDD mode i data sel PTR CTRL ts i OP REG OP-EX ptr. stall_i instr BGMEM emp/full_i CONTEXT-SW #0 linstr wr TIME-SCHEDULING stall o **Aggregate throughput** opush_o **INSTR-EX** BGMEM 984 Gb/s PTR #1 CTRL PROG-LD up/lw rea up/lw_done

200 kGE ~5% of the node area

0.24-0.41 pJ/B/hop @typical corner



A multi-mode, flexible, compact and energy-efficient (0.24 pJ/B/hop) routing controller for heterogeneous NPU cluster, which is planned to be taped-out in the near future.

Welcome to the poster for more details!

- Support for arbitrary deterministic routing algorithms;
- Capability on data multicasting & path through;
- Capability on context switching;
- Adaptability to various topologies;
- Implementation & Evaluation in STM FD-SOI 28nm technology;

Hong Pang / ETH Zurich

• BG-CTRL integration into node.