

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

1 Motivation

3 Architecture

- NVM-based AIMC heterogeneous NPU cluster is a promising solution for edge AI applications:
- AIMC: Performs MVMs directly inside memory, where data resides;
- Other digital units: Accuracy critical operations, operations not supported by AIMC



- We need an efficient infrastructure for **inter-tile communication**:
- High generality: Can be applied to arbitrary topologies and a wide range of AI applications
- High routing flexibility: Be free to switch data paths at different time stages during a single application
- High energy efficiency & low area overhead

We propose **Borderguard Controller (BG-CTRL):** A **multi-mode**, **distributed**, **light-weight** and **flexible routing controller** for inter-

• Time-scheduling mode

- Instruction integrates operation, data source/destination, transmission timestamp
- Supports (nested) loops to ease stall_of programming and trim code size
- Data-driven mode
- Highest energy efficiency



- Configuration register-based, fixed data path during runtime
- Data is transmitted as long as source-sink handshaking passes
- Time sliced data-driven mode
- Mixed advantages of routing flexibility and energy efficiency
- Extended data-driven mode to support multiple sets of data path configuration, each time-multiplexed in a defined period

tile data transmission in heterogeneous NPU cluster.

2 Integration and Evaluation



- Context switching enables application switching without system reset
- Double-bank code memory (BGMEM) hides program reloading latency

 TABLE I INSTRUCTION SET AND FORMAT FOR TIME-SCHEDULING MODE.

 Field 2
 Field 1
 Field 0
 Functionality

 INST[19:16]
 INST[15:12]
 INST[11:0]
 TS
 TS
 TS
 TS
 Configures the upper 20 bits of the active timestamp (0 by de

	[]			
SET_TS	TS	TS	TS	Configures the upper 20 bits of the active timestamp (0 by default).
SET_OTS	NA	NA	Offset TS	Configures the implicit active offset timestamp value (1 by default).
INC_TS	NA	NA	TS	Increments the upper 20 bits of the active timestamp by 1.
FWIM	DIR	NA	TS	Forwards data from the specified source (DIR).
\mathbf{FW}	DIR	NA	Offset TS	DIR: OFIFO ID (for data popping) or side (for data forwarding/pushing).
POPUSHIM	RP	RP	TS	DOBC: Pops data from an OFIFO; DIBC: Pushes data into an IFIFO.
POPUSH	RP	RP	Offset TS	
REPEATIM	NR	RP	TS	Forms a loop with NR instructions, repeating for RP iterations.
REPEAT	NR	RP	Offset TS	RP = 0 indicates infinite loop.
REPEATL	NR[9:6]	RP [9:6]	$\{NR[5:0], RP[5:0]\}$	REPEATL active timestamp is set using SET_OTS.
WAITIM	NA	NA	TS	Keeps BG-CTRL idle, with data-selecting signal unchanged.
WAIT	NA	NA	Offset TS	
RESTART	RP	RP	TS	Reruns the entire program for RP iterations. ($RP = 0$ indicates infinite iterations.)
DONE	NA	NA	TS	Signals the end of the program.
NA: Unused field TS: Active timestamp DIR mapping configuration: 0: West, 1: North				guration: 0: West, 1: North, 2: East, 3: South, (N+4): OFIFO #N

4 Conclusion

Operation

INST[23:20]

We present **BG-CTRL**, a **multi-mode**, **compact**, **flexible**, and **energy-efficient** (0.24 pJ/B/hop) **routing controller** for heterogeneous NPU cluster:

- Without data-packeting or conversion circuits, it shows low
- Integration: We integrate BG-CTRLs into node wrapper and build a 3x3 mesh structure for evaluation.
- Implementation: We synthesize our design using STM FD-SOI 28nm, worst-case (SS/0.9V/-40°C)
- BG-CTRL can run up to 640 MHz under path-through fashion (limited by long data transmission path)
- Aggregate throughput of 984 Gb/s
- Only 204 kGE area overhead (5% of node area), can be lower!
- Compared to the state-of-the-art, our design is 2.5x smaller than work [1], and 2.5x more energy efficient than work [2].







area overhead of 204 kGE (can be min. 145 kGE)

- With programmability, it supports arbitrary deterministic routing algorithms, and adapts to various AI applications
- With **distributed** property, it can work on **arbitrary topologies** and supports **data multicasting.**

Reference

[1] T. Fischer, M. Rogenmoser, M. Cavalcante, F. K. G⁻⁻urkaynak and L. Benini, "FlooNoC: A Multi-Tb/s Wide NoC for Heterogeneous AXI4 Traffic," in IEEE Design & Test, vol. 40, no. 6, pp. 7-17, Dec. 2023

[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.



