

# EmuNoC: Hybrid Emulation for Fast and Flexible Network-on-Chip Prototyping on FPGAs

Yee Yang Tan\*, Felix Staudigl\*, Lukas Jünger\*, Anna Drewes†, Rainer Leupers\*, and Jan Moritz Joseph\*

\*Institute for Communication Technologies and Embedded Systems, RWTH Aachen University, Germany

{tan, staudigl, juenger, leupers, joseph}@ice.rwth-aachen.de

†Institute for Information and Communication Technologies, Otto-von-Guericke Universität Magdeburg, Germany

anna.drewes@ovgu.de

**Abstract**—Networks-on-Chips (NoCs) recently became widely used, from multi-core CPUs to edge-AI accelerators. Emulation on FPGAs promises to accelerate their RTL modeling compared to slow simulations. However, realistic test stimuli are challenging to generate in hardware for diverse applications. In other words, both a fast and flexible design framework is required. The most promising solution is hybrid emulation, in which parts of the design are simulated in software, and the other parts are emulated in hardware. This paper proposes a novel hybrid emulation framework called EmuNoC. We introduce a clock-synchronization method and software-only packet generation that improves the emulation speed by  $36.3\times$  to  $79.3\times$  over state-of-the-art frameworks while retaining the flexibility of a pure-software interface for stimuli simulation. We also increased the area efficiency to model up to an NoC with 169 routers on a single FPGA, while previous frameworks only achieved 64 routers.

**Index Terms**—Hybrid emulation, NoCs, FPGA

## I. INTRODUCTION

Today, massively-parallel multi-processors can be found in different forms in nearly any system. Their application ranges from conventional multi-core CPUs, e.g., in cloud servers, to massively scaled systems used as low-power neuromorphic edge AI accelerators. In any case, Networks-on-Chip (NoCs) have become the de facto communication infrastructure for their excellent scaling capability. Examples showing the nearly universal use of NoCs are data-flow or parallel processors, e.g., manycores [1], big data [2], server-scale AI [3], edge AI [4], data bases [5], in-memory computing [6], genome sequencing [7], medical applications [8].

Due to the wide use of NoCs, many design tools have been created, e.g., [9], [10]. They are used for design space exploration (DSE) with software simulations or hardware prototypes that evaluate key performance metrics (KPIs) and guide the architect. Traditionally, NoC simulators target multi-core CPUs. With the emergence of edge AI accelerators, there is a need for more versatile tools for changing applications or mappings, as contributed by this work.

A typical DSE flow comprises the following: Architectural simulators are fast and flexible, providing early critical insights. After implementing a register transfer level (RTL) model, these can be simulated. This is slow and practically only helpful in

Funded by the Federal Ministry of Education and Research (BMBF) and the Ministry of Culture and Science of the German State of North Rhine-Westphalia (MKW) under the Excellence Strategy of the Federal Government and the Länder (G:(DE-82)EXS-SF-Project No. StUpPD\_390-21)

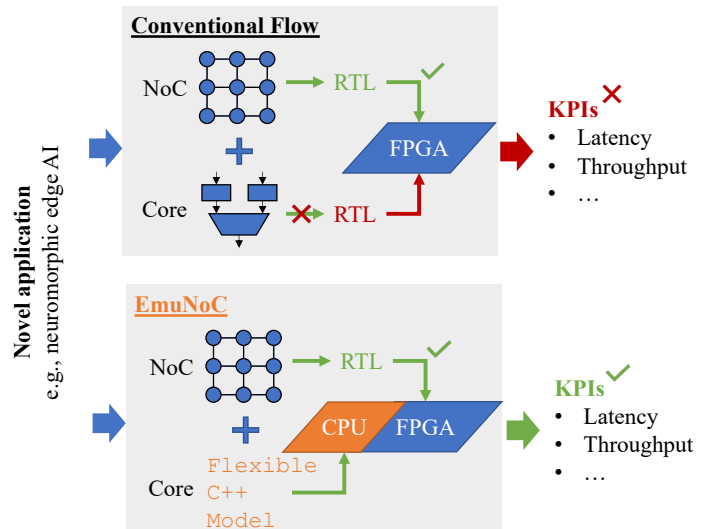


Fig. 1: EmuNoC: Toolflow

verifying or evaluating small parts of any design. Therefore, emulation on FPGAs is highly relevant as of massive speed improvements [11], [12] for any architecture optimization [13].

One key downside of FPGA emulation is its limited flexibility to adapt to novel use cases [12]. In other words, for every new NoC use case, the RTL design of all cores, etc., would be required. The availability of these models is unrealistic in an early design stage, hindering the deployment of NoCs for emerging use cases.

This paper provides an effective solution. We contribute an open-source framework<sup>1</sup> for NoC hybrid emulation, in which the traffic pattern can easily be switched by software models, but the NoC is emulated on the FPGA for high performance and accuracy (Fig. 1). We propose a novel clock-synchronization method and hardware-only packet generation that improves the emulation speed by more than one magnitude. Specifically, this paper yields the following novelties:

- We propose a hardware clock halting technique for faster hybrid emulation, which shows a speedup of  $79.3\times$  for synthetic traffic and  $36.3\times$  for Netrace [14] over the previous state-of-the-art emulation framework.

<sup>1</sup><https://github.com/ICE-RWTH/EmuNoC>

TABLE I: NoC-on-FPGA emulators.

Emulator	Model	TDM/DM	Transactor	Synthetic benchmark	Real benchmark	
					Netrace [14]	Edge-AI
Kouadri [16], [17]	Pure HW	-	-	✓	✗	✗
Wolkotte [18]	Hybrid	TDM	Bus	✓	✗	✗
Papamichael [19]	Hybrid	TDM	Bus	✓	✗	✗
HAsim [20]	Hybrid	TDM	Bus	✓	✗	✗
FNoC [21], [22]	Pure HW	TDM	-	✓	✗	✗
D'Hoore [23]	Pure HW	TDM	-	✓	✗	✗
Chu [12]	Pure HW	TDM	-	✗	✓	✗
AcENoCs [24]	Hybrid	DM	Bus	✓	✗	✗
Drewes et al. [11]	Hybrid	DM	Bus+Stream	✓	✓	✗
EmuNoC	Hybrid	DM	Stream	✓	✓	✓

- We contribute a novel concept of a single clock synchronous serializer used as a Network Interface (NI) for NoCs with virtual channels (VCs).
- We evaluate our system for different case studies (multi-core CPUs, edge AI accelerators) to show flexibility.

The paper is organized as follows. In Sec. II we will introduce the background for hybrid emulation of NoCs and discuss the related works. In Sec. III, we will explain our architecture. In Sec. IV we will analyze the system performance. Finally, the paper is concluded.

## II. BACKGROUND AND RELATED WORKS

When building efficient NoCs, they must be evaluated against benchmarks. For this, stimuli are to be generated. Synthetic traffic enables system validation, e.g., through fuzzy testing using random traffic, but it is not helpful for architectural exploration as it does not reflect real workloads. Full-system simulators (FSS), such as gem5, execute the whole system, including the operating system, cores, caches, and the NoC. The method offers the highest-precision benchmarks but often is unacceptably slow. Traces, recorded with an FSS and replayed later, provide a useful middle-ground.

For modeling the NoC, the accuracy of traces is often sufficient, as demonstrated by Netrace [14] for multi-core CPUs. Netrace provides trace files and a dependency-driven C-based player. The traces are generated by running PARSEC benchmarks on gem5 for a 64-core system. The dependency tracking between packets boosts the accuracy. For AI systems, trace generation is often more straightforward than for CPUs, because of the deterministic execution of DNNs (e.g., no situational caching). Most mapping strategies of DNNs enable mathematical trace modeling. One example is NewroMap [15], which showed that the feed-forward characteristic of DNNs can be exploited for mapping neurons to neuromorphic, memory-bound accelerators. The resulting traffic patterns in the NoC yield high locality and a low number of dependencies.

Cycle-accurate simulators, RTL simulation, or FPGA hybrid emulation offer different options to understand the performance of NoCs before their deployment to a system, hence enabling design space exploration.

Cycle-accurate simulators offer a decent compromise between accuracy and speed. Booksim [25], Noxim [10], and Ratatoskr [26] are the most popular. Booksim [25] uses a channel model that implements a two-phase evaluate-update protocol for routers. Noxim [10] is implemented using the cycle-accurate modeling of the SystemC library. It allows

varying NoC parameters with configuration files. Ratatoskr [26] also builds upon SystemC but uses TLM for system-level benchmarks. A 64-node mesh network yields a simulation frequency between 1000 Hz and 10 000 Hz (cf. Fig. 8).

For even higher speed, emulation on FPGAs became the industry standard. A whole NoC might be too large for a single FPGA, so Kouadri et al. [16], [17] emulate it through partitioning the NoC and mapping it onto multiple FPGAs. However, the measurement accuracy is limited by the off-chip data transmission architecture. When using only a single FPGA, one can directly map the NoC (if the FPGA is large enough) or use time-division multiplexing (TDM). TDM [18]–[20] implements a single router in the FPGA's programmable logic. The status of the routers are stored in the off-chip memory. This method can emulate a large-scale NoC with over 1000 nodes, but it requires huge off-chip memory and deteriorates performance. Chu et al. [21], [22] partition the NoC into multiple virtual clusters, each cluster containing multiple routers. Each cluster is emulated sequentially, increasing emulation speed. 3D NoCs can also be implemented using TDM [23].

TDM is rarely used in industry, since FPGAs can easily be clustered to provide sufficient resources to accommodate the whole NoC. In other words, directly-mapped (DM) approaches became viable. This method can achieve the highest performance because all routers run in parallel. Different framework designs will degrade the emulated NoC frequency. AcENoC [24] ran a  $5 \times 5$  NoC using a separate bus system to transmit the data to/from the NoC and achieved a maximum frequency of 23kHz. For larger  $8 \times 8$  NoCs, Drewes et al. [11] used an AXI4-Stream bus to collect the data from the NoC with the technique of an asynchronous serializer, and achieved 16kHz.

To further boost performance, Netrace's dependency tracking between packets was implemented in hardware by [12] achieving up to 12MHz emulation speed. However, this limits the system's flexibility, as the benchmark cannot be replaced easily.

A complete comparison of all NoC emulation systems is given in Tab. I. Only EmuNoC provides the flexibility for changing applications at high frequency.

## III. EMULATION SYSTEM ARCHITECTURE

The previous state-of-the-art emulation frameworks using directly-mapped NoCs are limited in performance from a) the bus system data transaction [11], [24], and b) the software clock halting technique for cycle-accurate emulation [11]. This paper presents an improved approach tackling both downsides to achieve an even faster emulation speed.

The architecture of EmuNoC is shown in Fig. 2. It consists of the software (virtual platform/green), executed in the CPU cores of our FPGA, the hardware (RTL design/blue), executed in the programmable logic, and a transactor (orange) that connects both of them. We propose a novel transactor to overcome the performance bottleneck with the AXI4-Stream data transaction and hardware-based clock halting technique. On the software side, a virtual platform generates and sends packets to the RTL model at a defined time quantum (*injection cycle*) through the transactor and places it in the virtual hardware buffer. A clock halter enables to stop the execution of the RTL model at any

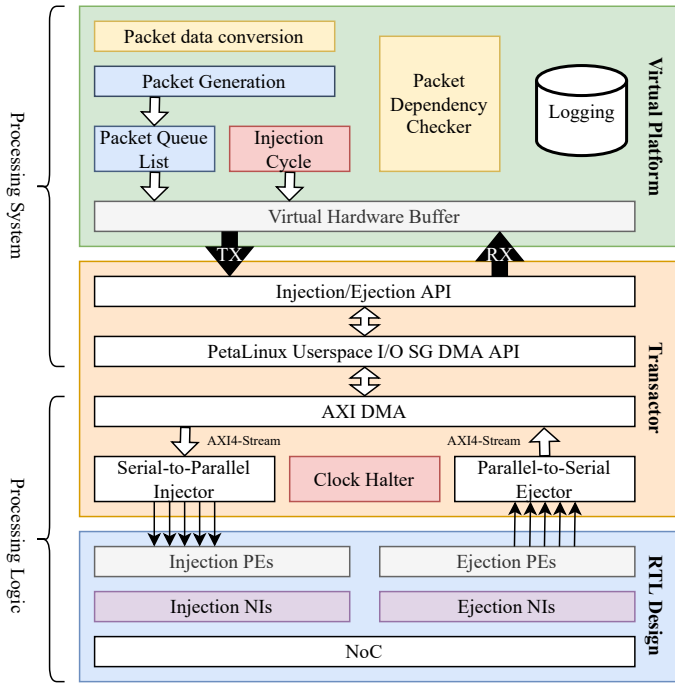


Fig. 2: EmuNoC system architecture.

given time. This unit is used in two scenarios: For packet injection, the transactor sets the time quantum step to the clock halter through the serial-to-parallel injector and enables RTL execution until the next given *injection cycle*. When the NoC RTL model is executed (*i.e.*, not halted), the serial-to-parallel injector injects packets into the source PEs. The parallel-to-serial ejector halts the RTL model via the clock halter when the packets arrive at a destination processing element (PE) for packet ejection. Then, it sends the halted time step (recorded by the clock halter) and the arrived packets via DMA to the software-side buffer through the transactor. The virtual platform checks, removes, and logs received packets.

#### A. Hardware Architecture

Fig. 3 shows the hardware architecture of our novel transactor with its adjacent units.

1) *Clock Halter*: The clock halter stops the execution of the emulated NoC at any time for synchronization with the virtual platform. It is a central logic connected to nearly all other components (see Fig. 3). Fig. 4 shows its block diagram. The *halting clock* is generated by a buffer driven by the *global clock* and enabled by a *ctrl* signal ( $=1$ ). An *injection cycle* is stored using the *write enable* signal. The *counter* counts cycles. The *ctrl* signal enables the *halting clock* when the value is smaller than the stored *injection cycle* value. The signal *ctrl* halts the buffer and the *counter* stops counting whenever the *halt* signal is 1. If the *counter* equals to the *injection cycle*, the *stop* signal is set and the *halting clock* is disabled.

2) *Serial-to-Parallel Injector*: The serial-to-parallel injector receives packets from the software side and injects them into the NoC. For communication with the software, the injector contains an AXI4-Stream slave port. When a transaction starts,

the first stream data determines the *injection cycle* until which the RTL design is executed. Once reached, the injector converts the packet into the header flits (the communication unit within the NoC; *conv* in Fig. 3) and sends them to the respective PE's FIFO using the packet's source address.

3) *Parallel-to-Serial Ejector*: When a destination PEs receives a complete packet (*i.e.*, all of its flits), the parallel-to-serial ejector instructs the clock halter to halt the RTL design through the *halt* signal ( $=1$ ). This module converts the header flit back (*iconv* in Fig. 3) to packet data and sends it via DMA (AXI4-Stream) to the software side. The sent data contain the clock cycle at which the packet was received. Fig. 5 shows the implementation logic of the single clock serializer. The blue part refers to the corresponding PE's FIFO. If a complete packet arrives at the destination, only the header flit is stored in the FIFO, the corresponding signal *read valid* becomes 1. The block *or reduce* (purple) tells the FSM (yellow) to initiate the stream transaction. The FSM updates the *round robin arbiter* (green) with the signal *ctrl* ( $=1$ ) to decide which FIFO to read. When data is read from the AXI4-stream port (*tready* $=1$ ), the corresponding FIFO is read through the multiplexer's output signal (red). If all header flits in the FIFOs are ejected (*halt* $=0$ ), the round-robin arbiter (selecting the ejection FIFO order) will be updated, and the RTL emulation will continue.

4) *Injection PE*: The injection PE subsequently injects flits of each packet into the NoC. For each PE, there is a network interface (NI) that handles the assignment of flits to VCs and sends them into the connected router (see below).

The injection PE contains a FIFO and an FSM, managing the AXI4-Stream transactions and packet injection via a NI. When the PE's FIFO holds the first flits of a packet (header flit), it starts the transaction by injecting it. The packet contains dummy payload flits as our implementation does not transmit "useful" payloads and handles the packet assignment at the destination via the software's virtual buffers.

5) *Injection NI*: The NI accepts a complete packet in one AXI4-Stream transaction. If the NoC uses multiple VCs, packets will be assigned via round-robin arbitration. The flits travel through the NoC until they are received in the target PE's NI. In the implementation at hand, we use the Ratatoskr router [13]; it can be exchanged by any other NoC that implements an AXI4-compatible interface.

6) *Ejection NI*: The ejection NI contains one FIFO per VC; the FIFO is long enough to store a whole packet, *i.e.*, all of its flits. When all packet flits are received, the NI starts the AXI4-Stream transaction and sends the packet to its PE. The number of flits is checked via a comparator in Fig. 5.

7) *Ejection PE*: This PE functions like the injection PE. It receives flits from its NI. Only the head flit is stored to be transmitted to the software side. This flit is put into a 1-flit FIFO connected to the serializer when a packet is completed.

#### B. Software Architecture

Drewes et al. [11] used the simple DMA mode to transfer data to the NoC. EmuNoC uses Scatter Gather (SG) DMA mode to improve performance. On the software side, the SG

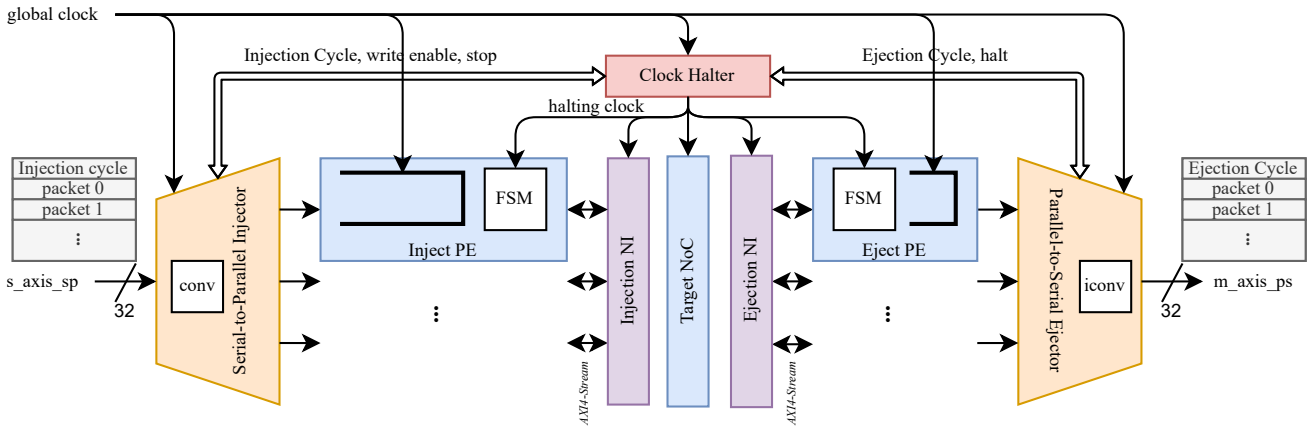


Fig. 3: EmuNoC hardware architecture.

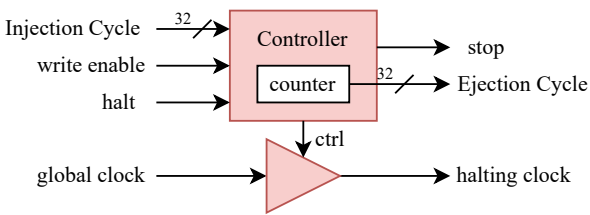


Fig. 4: Block diagram of the clock halter.

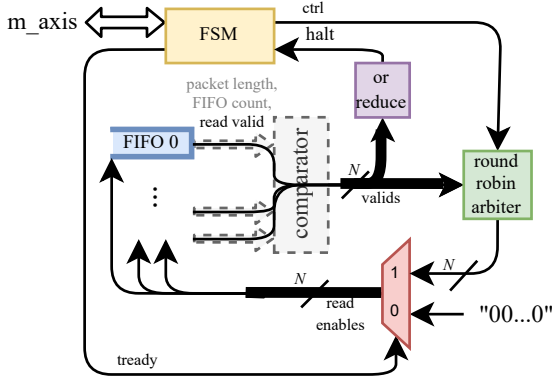


Fig. 5: Serializer as used in the ejector; gray-shaded part for multi-VC NoCs.

DMA API for PetaLinux Userspace I/O [27] is implemented (Fig. 2). We also used compiler optimizations, *e.g.*, `-Ofast`.

Fig. 6 shows the software design of the framework. It is the detailed view of the virtual platform in Fig. 2. The variables *icyc* is the *injection cycle*, *src* the packet’s source address, *dst* the destination address, and *len* the packet length/flit count.

In general, our software is executed in the following six steps (highlighted in Fig. 6):

- (1) Packet data is generated using the flexible software interface (example see below).
- (2) The program searches for the earliest available packets and puts them into the queue.

- (3) The *virtual hardware buffer* sends the packet data to the Injection PE’s FIFO. It also keeps a copy to handle the assignment of the received flits to the correct packet, as explained above. The *injection cycle* and the packets are sent to the NoC and transmitted there to the target PE.
- (4) When the packets have arrived at their destination, the DMA stored them in the main memory. Then, the program compares the received packet, which is matched with its counterpart in the *virtual hardware buffer* to enable dependency tracking.
- (5) After injecting the previous time quantum in (3), the program needs to determine the next time quantum to inject the packets. Then, the program checks whether the next time quantum has exceeded the user-defined maximum cycle to run. If the program reaches the maximum cycle or no more packet to inject, it goes to (6); else to (2).
- (6) If the virtual buffer is empty, the emulation will stop.

We will demonstrate the flexibility of this software architecture with three different traffic scenarios (see below). Any other use case can be implemented easily by modifying the software code. A simple example code exemplifies this (Listing 1). Line 1 refers to the yellow blocks in Fig. 6, different modes are decided, and the metadata is generated. The rest of the program refers to the six steps, in which a packet list is iterated and injected to the NoC.

Listing 1: Example packet generation.

```

1 metadata = initialization(mode);
2 pkt_cyc_list = generate_packets(metadata);
3 do
4 {
5     if (cyc < max_cyc)
6         put_packet_to_queue(cyc, pkt_cyc_list, queue_lists);
7     hw_list = copy_to_hw_buffers_and_create_hw_list(hw_buffers, queue_lists);
8
9     inject(cyc, hw_list);
10    eject(hw_buffers);
11
12    cyc = calculate_next_injection_cycle(pkt_cyc_list);
13 } while (cyc < max_cyc);
14 check_lost_packets(hw_buffers);

```

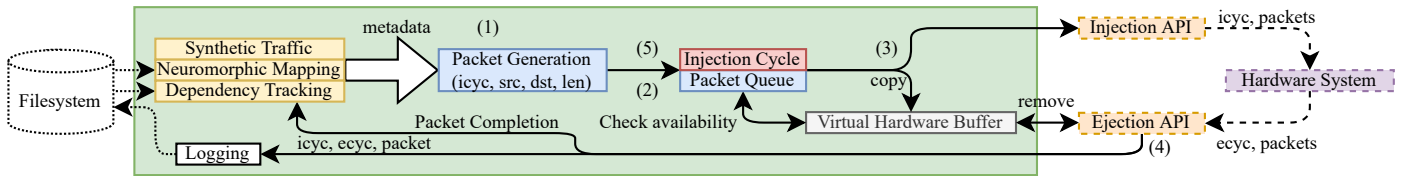


Fig. 6: EmuNoC software architecture.

TABLE II: Hardware resources of different NoC configurations of each emulator. (FB: flit buffer)

Emulator	NoC	VC	FB	LUT	LUTRAM	FF	BRAM	BUFG	Slices
EmuNoC	5×5	2	8	40750	6150	45944	39.5	3	-
EmuNoC	8×8	2	3	100148	14934	105960	98	3	-
EmuNoC	13×13	2	4	267230	39174	274644	255.5	3	-
EmuNoC	13×13	-	-	34261	1110	38666	86.5	3	-
AcENoC [24]	5×5	2	8	52520	840	19569	-	-	-
Drewes [11]	8×8	2	3	97515	-	78361	36.5	-	-
Chu [12]	8×8	-	-	133327	-	111576	830	-	43903

#### IV. RESULTS

##### A. Hardware Costs

Table II shows the used FPGA resources of EmuNoC, AcENoC [24], Drewes et al. [11] and Chu [12]. EmuNoC's resources were obtained from Vivado 2018.2 for a Zynq UltraScale+ MPSoC ZCU102. The *global clock* (Fig. 3) is set to 80MHz and the FIFOs in the NoC and the transactor use the Xilinx's FIFO IP [28] to support larger setups.

As we can see, the FPGA resources of EmuNoC increase approximately linearly with the router count. EmuNoC used more memory (LUTRAM, BRAM) than AcENoC [24] and Drewes et al. [11], which is required for our better emulation performance. Still, we can host up to 169 routers, more than triple the router count than in the previous DM framework (even with larger single routers). We achieved this by enabling them to use larger standard FPGAs; [11] relied on a custom clock halter that only was possible in their FPGA. Chu [12] consumes more resources than EmuNoC (8×8) because they implement the dependency tracking in hardware.

##### B. Emulation Performance

EmuNoC is validated with uniform random traffic (uniform random source-destination pairs and injection times). This traffic allows to fuzzy-test the NoC as random traffic is sent though the network.

To compare EmuNoC's performance with the state-of-the-art emulation, the NoC is set to the same configuration as AcENoC [24] (5×5 mesh with 2 VCs and 8-flit buffer) and Drewes et al. [11] (8×8 mesh with 2 VCs and 3-flit buffer). The largest NoC that can be emulated on our FPGA is 13×13 mesh NoC with 2 VCs and 4-flit buffer. The emulation performance of these configurations are shown in Fig. 7. We observe a performance degradation with NoC size and injection rate.

Table III shows the emulation frequency at 5% flit injection rate for synthetic traffic. EmuNoC achieves 2221 kHz for 5×5 mesh, 1319 kHz for 8×8 mesh. EmuNoC yields a 96.6× speedup over AcENoC [24], and a 79.3× speedup over Drewes

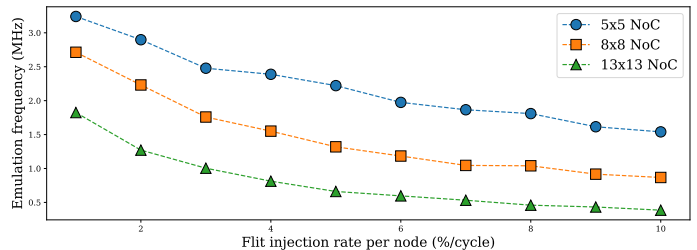


Fig. 7: Emulation performance with uniform random traffic.

et al. [11]. We have achieved faster emulation speed than any other flexible, directly-mapped framework.

##### C. Performance Scaling

As stated, the emulation performance drops with NoC size and traffic load. We compare our system against simulators to analyze this scaling effect as they show the same behavior from the similar root of traffic injection.

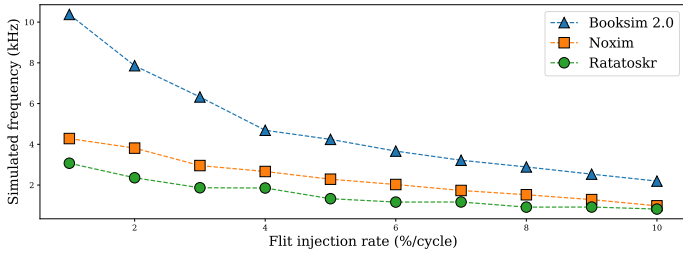
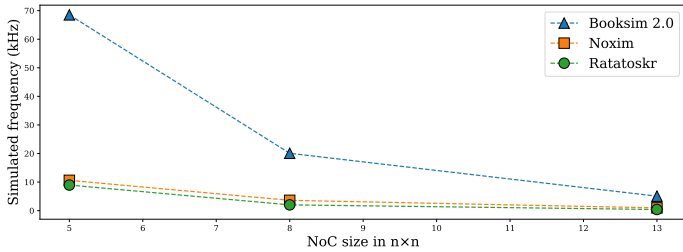
We simulate a 13×13 mesh NoC (2 VCs and 4-flit buffer) with the simulators Booksim 2.0 [25], Noxim [10], and Ratatoskr [26] using dimension-ordered routing, 5-flit packets. We measure the median of 10 simulations on a WSL Ubuntu 20.04.2 LTS using one Intel i7-5700HQ core at 2.7 GHz and show the results in Fig. 8. By increasing the flit injection rate, we can see the simulation performance decreases (Fig. 8(a)). Fig. 8(b) shows three different NoC sizes at fixed 5% injection rate. A larger injection rate or NoC size degrades the simulation performance because the simulation needs to generate traffic and simulate the routers sequentially.

Analyzing the performance drop of these simulators and EmuNoC, we found out that Booksim 2.0 [25] has the highest performance loss from 1 to 10% injection rate for 13×13 NoC (78.9%; EmuNoC is 78.8%; Noxim [10] is 77.1%; Ratatoskr [26] is 73.3%). Emulation behaves similarly to the simulations because of the software-side traffic generation.

If the NoC size is increases from 5×5 to 13×13 (fix 5% flit injection rate per cycle), Ratatoskr [26] has lost the most performance (95.4%) compared to Booksim 2.0 [25] (92.6%), Noxim [10] (90.8%) and EmuNoC (70.2%). Here we observe an advantage of emulation, as EmuNoC yields the lowest performance drop.

##### D. Case Study I: Multi-core Processors

Netrace [14] provides 64-core processor's traces with dependency tracking. It contains five phases from the PARSEC benchmarks: the startup, warmup, region of interest (ROI), result output, and post benchmark. Drewes et al. [11] has run


 (a) Simulation performance for different flit injection rates in a  $13 \times 13$  NoC.


(b) Simulation performance for different size NoC sizes at 0.5% flit injection rate.

Fig. 8: Simulation performance with NoCs.

TABLE III: Performance comparison of frameworks.

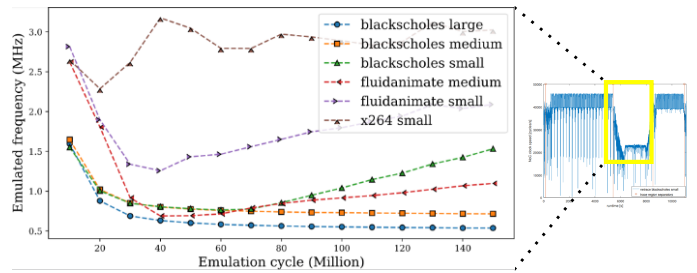
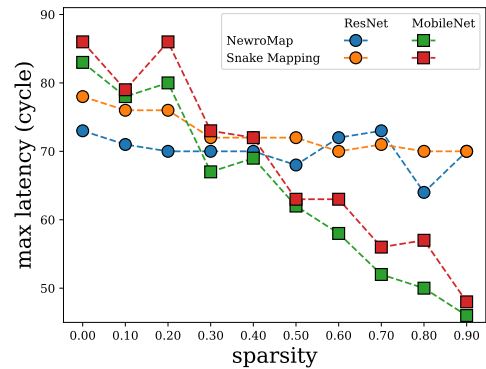
Emulator	NoC	Traffic	Frequency [kHz]	Speedup
AcENoC [24]	$5 \times 5$	Synthetic	23	96.6
Drewes et al. [11]	$8 \times 8$	Synthetic	16.639	79.3
Drewes et al. [11]	$8 \times 8$	Netrace	39.243	36.3
Chu [12]	$8 \times 8$	Netrace	12979	0.11
EmuNoC	$5 \times 5$	Synthetic	2221.464	-
EmuNoC	$8 \times 8$	Synthetic	1319.333	-
EmuNoC	$8 \times 8$	Netrace	1426.404	-
EmuNoC	$13 \times 13$	Synthetic	661.291	-

the whole benchmark (Fig. 9 right-side yellow box). Their results show a performance drop in the ROI (yellow highlighted) because this region contains the highest traffic workload. As the ROI is hence the exciting part, only it is investigated in our experiment. The result is shown in Fig. 9, left-hand side. Similar to [11], we first observe a performance drop then followed by a performance recovery. In average, we achieve 1426 kHz (see Table III), which has achieved  $36.3 \times$  speedup compared to Drewes et al. [11]. Our framework is slower than Chu [12] by  $0.11 \times$ . The reason is that Chu [12] implemented Netrace-specific dependency-tracking hardware. While this method has a high performance, it is not flexible for application-centered engineers to adopt different use cases easily. Therefore, our software-based dependency tracking offers a compelling trade-off between performance and flexibility for many practitioners.

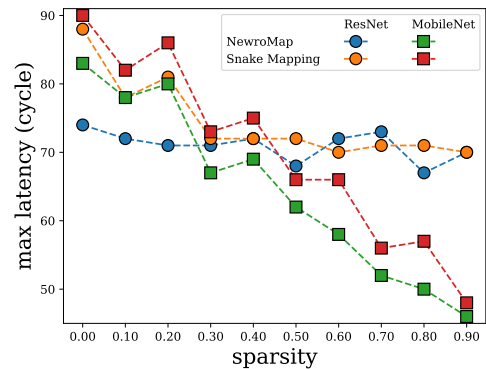
### E. Case Study II: Neuromorphic Edge-AI Accelerator

Studying edge AI-accelerator architectures [29]–[32] shows that all of them demand a scalable NoC. This yields a workload mapping problem, *e.g.*, solved by NewroMap [15] for CNNs, in which the network activations are transmitted via the NoC.

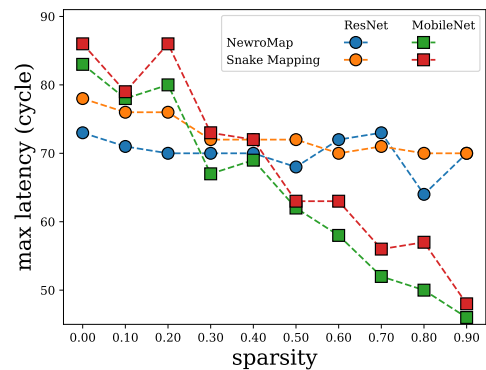
One well-known property of neural networks is their high sparsity, *i.e.*, 0-values, which need not be sent via the NoC [33]. Hence, the effective injection rate in the NoC for each


 Fig. 9: Emulation performance for PARSEC [14] ROI for different traces. ( $8 \times 8$  NoC with 2 VCs, 3-flit buffer). The right-hand-side image is cited from [11].


(a) 1 VC, 2-flit buffer



(b) 2 VCs, 1-flit buffer



(c) 2 VCs, 2-flit buffer

Fig. 10: Maximum latency measurement for different CNN mappings [15] for sparsity rates.

mapping is scaled by a) a sparsity factor and b) the target framerate in our exemplary case of video applications. This gives the formula for the injection per PE rate as

$$irate = \frac{map_{neurons} * (1 - sparsity) * framerate}{frequency_{NoC}},$$

where  $map_{neurons}$  is the number of neurons mapped to a core. We use the  $framerate$  and  $frequency_{NoC}$  from one commercially available accelerator called NeuronFlow [4] as 30 FPS and 1 GHz.

We analyzed different mappings and NoC architectures to demonstrate the flexibility of EmuNoC. We plot the maximum packet latency in Fig. 10. We used very lightweight NoCs with/without VCs and multi-flit buffers. As edge AI workloads tend to have particular traffic patterns with high locality, these architectures are promising to investigate (even though they would be not useful for conventional multi-core CPUs).

We observe that the maximum packet latency decreases with higher sparsity for all architectures. This effect is as expected as less traffic yields less congestion. We further observe that the optimized mappings proposed by NewroMap can improve latency vs. snake mapping, which was previously only demonstrated in simulations in [15].

Another interesting finding is that a NoC without VCs and 2-flit buffers has a lower maximum latency than the architectures with 2 VCs and 1-flit buffers (see Fig. 10(a) vs. Fig. 10(b)). Both architectures yield approximately the same area costs. At first glance, this is surprising since VCs promise better peak performance. However, the high locality of edge-AI traffic patterns effectively removes the need for VCs in many routers. Also, adding additional buffers (Fig. 10(c)) does not improve performance. These findings advocate for very lightweight NoCs in edge-AI systems. However, the authors would like to mention that a VC-less router will not be the best architecture to choose in all cases. While it has the lowest area costs and best performance, it effectively prohibits multi-thread processing of CNN layers on single cores. Hence, the higher NoC costs might be worth it from a system design perspective.

## V. CONCLUSION

This paper proposed a fast and flexible FPGA-based NoC hybrid emulation called EmuNoC. These features are achieved by a novel transactor architecture and a programmable software interface. The transactor contains a novel clock-synchronization method and hardware-only packet generation unit. The programmable software interface enables mapping different system benchmarks to the NoC, which is highly relevant as NoCs are widely used today. EmuNoC has achieved  $36\times$  to  $96\times$  speedup compared to the comparable previous DM method. We also increased the area efficiency and were able to emulate a NoC with 169 routers on a single FPGA with a state-of-the-art size at time of writing this paper. We used the emulator in two case studies to demonstrate its practical use for architects.

## REFERENCES

- [1] A. Bakhoda, J. Kim, and T. M. Aamodt, "Throughput-effective on-chip networks for manycore accelerators," in *2010 43rd Annual IEEE/ACM International Symposium on Microarchitecture*. IEEE, 2010, pp. 421–432.
- [2] C. Blochwitz, J. M. Joseph, R. Backasch, T. Pionteck, S. Werner, D. Heinrich, and S. Groppe, "An optimized radix-tree for hardware-accelerated dictionary generation for semantic web databases," in *2015 International Conference on ReConfigurable Computing and FPGAs (ReConFig)*. IEEE, 2015, pp. 1–7.
- [3] M. Davies, N. Srinivasa, T.-H. Lin, G. China, Y. Cao, S. H. Choday, G. Dimou, P. Joshi, N. Imam, S. Jain *et al.*, "Loihi: A neuromorphic manycore processor with on-chip learning," *Ieee Micro*, vol. 38, no. 1, pp. 82–99, 2018.
- [4] O. Moreira, A. Yousefzadeh, F. Chersi, A. Kapoor, R.-J. Zwartenkot, P. Qiao, G. Cinserin, M. Khoei, M. Lindwer, and J. Tapson, "Neuronflow: A hybrid neuromorphic – dataflow processor architecture for ai workloads," in *2020 2nd IEEE International Conference on Artificial Intelligence Circuits and Systems (AICAS)*, 2020, pp. 01–05.
- [5] C. Blochwitz, J. Wolff, J. M. Joseph, S. Werner, D. Heinrich, S. Groppe, and T. Pionteck, "Hardware-Accelerated radix-tree based string sorting for big data applications," in *International Conference on Architecture of Computing Systems*. Springer, 2017, pp. 47–58.
- [6] R. Guirado, A. Rahimi, G. Karunaratne, E. Alarcón, A. Sebastian, and S. Abadal, "Wireless On-Chip Communications for Scalable In-memory Hyperdimensional Computing," *arXiv preprint arXiv:2205.10889*, 2022.
- [7] S. Sarkar, G. R. Kulkarni, P. P. Pande, and A. Kalyanaraman, "Network-on-chip hardware accelerators for biological sequence alignment," *IEEE Transactions on Computers*, vol. 59, no. 1, pp. 29–41, 2009.
- [8] D. Passaretti, J. M. Joseph, and T. Pionteck, "Survey on FPGAs in medical radiology applications: Challenges, architectures and programming models," in *2019 International Conference on Field-Programmable Technology (ICFPT)*. IEEE, 2019, pp. 279–282.
- [9] J. M. Joseph, L. Bamberg, I. Hajjar, B. R. Perjikolaie, A. García-Ortiz, and T. Pionteck, "Ratatoskr: An open-source framework for in-depth power, performance, and area analysis and optimization in 3d nocs," *ACM Trans. Model. Comput. Simul.*, vol. 32, no. 1, sep 2021. [Online]. Available: <https://doi.org/10.1145/3472754>
- [10] V. Catania, A. Mineo, S. Monteleone, M. Palesi, and D. Patti, "Noxim: An open, extensible and cycle-accurate network on chip simulator," in *2015 IEEE 26th international conference on application-specific systems, architectures and processors (ASAP)*. IEEE, 2015, pp. 162–163.
- [11] T. Drewes, J. M. Joseph, and T. Pionteck, "An FPGA-based prototyping framework for Networks-on-Chip," in *2017 International Conference on ReConfigurable Computing and FPGAs (ReConFig)*. IEEE, 2017, pp. 1–7.
- [12] T. V. Chu, K. Kise, and K. Tanaka, "Dependency-Driven Trace-Based Network-on-Chip Emulation on FPGAs," in *Proceedings of the 2020 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, ser. FPGA '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 211–221. [Online]. Available: <https://doi.org/10.1145/3373087.3375309>
- [13] J. M. Joseph, L. Bamberg, D. Ermel, B. R. Perjikolaie, A. Drewes, A. García-Ortiz, and T. Pionteck, "Nocs in heterogeneous 3d socs: Co-design of routing strategies and microarchitectures," *IEEE Access*, vol. 7, pp. 135 145–135 163, 2019.
- [14] J. Hestness, B. Grot, and S. W. Keckler, "Netrace: dependency-driven trace-based network-on-chip simulation," in *NoCArc '10*, 2010.
- [15] J. M. Joseph, M. S. Baloglu, Y. Pan, R. Leupers, and L. Bamberg, "NEWROMAP: Mapping CNNs to NoC-Interconnected Self-Contained Data-Flow Accelerators for Edge-AI," in *Proceedings of the 15th IEEE/ACM International Symposium on Networks-on-Chip*, ser. NOCS '21. New York, NY, USA: Association for Computing Machinery, 2021, p. 15–20. [Online]. Available: <https://doi.org/10.1145/3479876.3481591>
- [16] A.-M. Kouadri-Mostéfaoui, B. Senouci, and F. Pétrot, "Scalable multi-fpga platform for networks-on-chip emulation," in *2007 IEEE International Conf. on Application-specific Systems, Architectures and Processors (ASAP)*. IEEE, 2007, pp. 54–60.
- [17] A. Kouadri-Mostéfaoui, F. Rousseau, and F. Pétrot, "Large scale on-chip networks: An accurate multi-fpga emulation platform," in *11th EUROMICRO Conference on Digital System Design Architectures Methods and Tools (DSD'08)*. IEEE Computer Society, 2008, pp. 3–9.
- [18] P. T. Wolkotte, P. K. Holzspies, and G. J. Smit, "Fast, accurate and detailed noc simulations," in *First International Symposium on Networks-on-Chip (NOCS'07)*. IEEE, 2007, pp. 323–332.
- [19] M. K. Papamichael, "Fast scalable fpga-based network-on-chip simulation models," in *Ninth ACM/IEEE International Conference on Formal Methods and Models for Codesign (MEMPCODE2011)*. IEEE, 2011, pp. 77–82.

- [20] M. Pellauer, M. Adler, M. Kinsky, A. Parashar, and J. Emer, "Hasim: Fpga-based high-detail multicore simulation using time-division multiplexing," in *2011 IEEE 17th International Symposium on High Performance Computer Architecture*. IEEE, 2011, pp. 406–417.
- [21] T. V. Chu, S. Sato, and K. Kise, "Fast and cycle-accurate emulation of large-scale networks-on-chip using a single fpga," *ACM Transactions on Reconfigurable Technology and Systems (TRETS)*, vol. 10, no. 4, pp. 1–27, 2017.
- [22] T. Van Chu, S. Sato, and K. Kise, "Ultra-fast noc emulation on a single fpga," in *2015 25th International Conference on Field Programmable Logic and Applications (FPL)*. IEEE, 2015, pp. 1–8.
- [23] J. D'Hoore, P. Bahrebar, and D. Stroobandt, "3D NoC emulation model on a single FPGA," in *Proceedings of the Workshop on System-Level Interconnect: Problems and Pathfinding Workshop*, 2020, pp. 1–8.
- [24] S. Lotlikar, V. Pai, and P. V. Gratz, "Acenocs: A configurable hw/sw platform for fpga accelerated noc emulation," in *2011 24th International Conference on VLSI Design*, 2011, pp. 147–152.
- [25] N. Jiang, D. U. Becker, G. Michelogiannakis, J. Balfour, B. Towles, D. E. Shaw, J. Kim, and W. J. Dally, "A detailed and flexible cycle-accurate network-on-chip simulator," in *2013 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, 2013, pp. 86–96.
- [26] J. M. Joseph, S. Wrieden, C. Blochwitz, A. García-Ortiz, and T. Pionteck, "A simulation environment for design space exploration for asymmetric 3d-network-on-chip," in *2016 11th International Symposium on Reconfigurable Communication-centric Systems-on-Chip (ReCoSoC)*, 2016, pp. 1–8.
- [27] L. Hans-Jürgen Koch Linux developer, "The Userspace I/O HOWTO," 2006. [Online]. Available: <https://www.kernel.org/doc/html/v4.12/driver-api/uio-howto.html>
- [28] Xilinx, "FIFO Generator v13.1 LogiCORE IP Product Guide," April 5, 2017.
- [29] A. Shafiee, A. Nag, N. Muralimanohar, R. Balasubramonian, J. P. Strachan, M. Hu, R. S. Williams, and V. Srikumar, "ISAAC: A Convolutional Neural Network Accelerator with in-Situ Analog Arithmetic in Crossbars," in *Proceedings of the 43rd International Symposium on Computer Architecture*, ser. ISCA '16. IEEE Press, 2016, p. 14–26. [Online]. Available: <https://doi.org/10.1109/ISCA.2016.12>
- [30] A. Ankit, I. E. Hajj, S. R. Chalamalasetti, G. Ndu, M. Foltin, R. S. Williams, P. Faraboschi, W.-m. W. Hwu, J. P. Strachan, K. Roy, and D. S. Milojevic, "PUMA: A Programmable Ultra-Efficient Memristor-Based Accelerator for Machine Learning Inference," in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 715–731. [Online]. Available: <https://doi.org/10.1145/3297858.3304049>
- [31] P. Chi, S. Li, C. Xu, T. Zhang, J. Zhao, Y. Liu, Y. Wang, and Y. Xie, "Prime: A novel processing-in-memory architecture for neural network computation in reram-based main memory," *ACM SIGARCH Computer Architecture News*, vol. 44, no. 3, pp. 27–39, 2016.
- [32] X. Liu, M. Mao, B. Liu, H. Li, Y. Chen, B. Li, Y. Wang, H. Jiang, M. Barnell, Q. Wu *et al.*, "RENO: A high-efficient reconfigurable neuromorphic computing accelerator design," in *Proceedings of the 52nd Annual Design Automation Conference*, 2015, pp. 1–6.
- [33] X. Liu, W. Wen, X. Qian, H. Li, and Y. Chen, "Neu-noc: A high-efficient interconnection network for accelerated neuromorphic systems," in *2018 23rd Asia and South Pacific Design Automation Conference (ASP-DAC)*, 2018, pp. 141–146.