

Torrent: A Distributed DMA for Efficient and Flexible Point-to-Multipoint Data Movement

Yunhao Deng*, Fanchen Kong*, Xiaoling Yi, Ryan Antonio, Marian Verhelst
MICAS-ESAT, KU Leuven

{yunhao.deng, fanchen.kong, xiaoling.yi, ryan.antonio, marian.verhelst}@esat.kuleuven.be

Abstract—The growing disparity between computational power and on-chip communication bandwidth is a critical bottleneck in modern Systems-on-Chip (SoCs), especially for data-parallel workloads like AI. Efficient point-to-multipoint (P2MP) data movement, such as multicast, is essential for high performance. However, native multicast support is lacking in standard interconnect protocols. Existing P2MP solutions, such as multicast-capable Network-on-Chip (NoC), impose additional overhead to the network hardware and require modifications to the interconnect protocol, compromising scalability and compatibility.

This paper introduces *Torrent*, a novel distributed DMA architecture that enables efficient P2MP data transfers without modifying NoC hardware and interconnect protocol. *Torrent* conducts P2MP data transfers by forming logical chains over the NoC, where the data traverses through targeted destinations resembling a linked list. This *Chainwrite* mechanism preserves the P2P nature of every data transfer while enabling flexible data transfers to an unlimited number of destinations. To optimize the performance and energy consumption of *Chainwrite*, two scheduling algorithms are developed to determine the optimal chain order based on NoC topology.

Our RTL and FPGA prototype evaluations using both synthetic and real workloads demonstrate significant advantages in performance, flexibility, and scalability over network-layer multicast. Compared to the unicast baseline, *Torrent* achieves up to a $7.88\times$ speedup. ASIC synthesis on 16nm technology confirms the architecture’s minimal footprint in area (1.2%) and power (2.3%). Thanks to the *Chainwrite*, *Torrent* delivers scalable P2MP data transfers with a small cycle overhead of 82CC and area overhead of $207\ \mu\text{m}^2$ per destination.

Index Terms—Direct Memory Access (DMA), Multicore SoC (MPSoC), AXI protocol, Network-on-chip (NoC)

I. INTRODUCTION

The widening imbalance between compute capabilities and on-chip communication requirement needs has emerged as a fundamental bottleneck in modern Systems-on-Chip (SoCs). While transistor density continues to increase following Moore’s Law, middle- and far-reach interconnect densities have remained nearly constant [1], leaving on-chip bandwidth as a scarce and critical resource. Although contemporary SoCs integrate increasingly powerful accelerators [2], [3], the performance of data-intensive workloads—such as Transformer-based large language models (LLMs)—is often constrained by memory bandwidth rather than processing elements [4]. Consequently, maximizing data reuse [5] has become a central strategy for

This project has been funded by the European Research Council (ERC) under grant agreement No. 101088865, the Flanders AI Research Program, and long term structural Methusalem funding by the Flemish Government.

*Both authors contributed equally to this research.

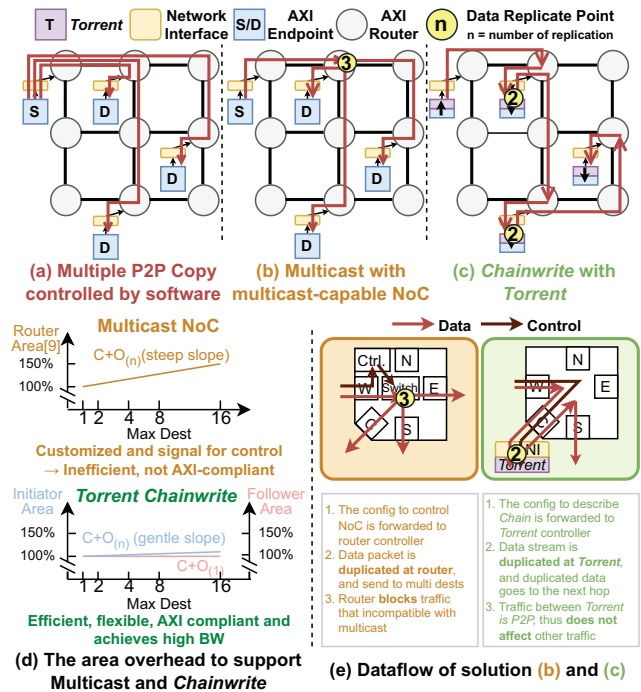


Fig. 1: The working principle, area overhead, and the dataflow of three P2MP data copy mechanisms with a 2D Mesh NoC

sustaining high accelerator performance. Data reuse can be achieved at both the accelerator level [6] and the SoC level. One example of SoC-level data reuse is tiled matrix multiplication, where one operand is tiled and the other operand needs to be distributed to multiple accelerators, creating a point-to-multipoint (P2MP) data movement pattern.

Software-based P2MP approaches [7] are conducted by issuing multiple independent requests from the source (Fig. 1(a)). However, these methods incur redundant memory accesses. To mitigate these inefficiencies, hardware support for multicast is being integrated into commercial [8] and academic [9] SoCs. By performing data replication within the interconnect (Fig. 1(b)), hardware multicast avoids redundant memory traffic and improves data reuse.

The Advanced eXtensible Interface (AXI) is a widely adopted on-chip communication protocol. Unfortunately, AXI lacks native support for P2MP communication [10]. One approach to deploying multicast on top of AXI is to modify the Network-on-Chip (NoC) routers, equipping them with multicast capabilities as illustrated in Fig. 1(e). This allows a single input stream

to be replicated to multiple output ports. While this approach is feasible, it requires protocol modifications and introduces extra control logic to prevent deadlocks [11]. More critically, as shown in Fig. 1(d), the complexity of NoC routers and the link width grow with the maximum number of multicast destinations [9]. Thus, the scalability of network-layer multicast is often limited.

In this paper, we propose *Torrent*¹, an alternative architecture enabling efficient P2MP data transfer in AXI-compatible SoCs. Instead of altering the NoC infrastructure, we embed P2MP capabilities directly within the Direct Memory Access (DMA) endpoints (Fig. 1(c)(e)). When a P2MP request is sent to one *Torrent* (called the initiator *Torrent*), it collaborates with all other *Torrents* attached to destination memories (called follower *Torrents*). Together, they form a data chain—an architecture resembling a doubly linked list that allows data to flow from the first node to the last node. This multicast mechanism is called *Chainwrite*². By shifting the data replication job from centralized NoC routers (at the network layer) to distributed *Torrents* (at the application layer), *Torrent* preserves full interoperability with AXI while allowing flexible P2MP capability. This approach supports a virtually unlimited maximal number of destinations ($N_{\text{dst,max}}$) without increasing router area or link width, while imposing ignorable hardware complexity (Fig. 1(d)).

Our main **contributions** are as follows:

- We propose *Torrent*, a distributed DMA architecture that supports efficient and scalable P2MP data transfers at the application layer, a mechanism called *Chainwrite*.
- *Chainwrite* supports an unbounded number of destinations without modifying the AXI standard and NoC routers by organizing endpoints into a *doubly linked list*.
- We develop scheduling algorithms that enable *Chainwrite* to be comparable with network-layer multicast in total number of hops for one P2MP data transfer task.
- We validate *Torrent* using both synthetic and real workloads and compare against SoTA P2P and multicast solutions. Synthetic results show *Torrent* achieves P2MP efficiency on par with multicast solutions. The real workloads extracted from the DeepSeek-V3 self-attention layers are executed on a *Torrent*-enhanced SoC in an FPGA, showing up to $7.88\times$ speedup over its unicast baseline. ASIC synthesis in 16nm and power analysis results show *Torrent* incurs only 1.2%/2.3% of the system area/power. The hardware overhead for adding one maximal destination with *Chainwrite* is minimized to $207\mu\text{m}^2$.

Compared to the SoTA DMAs and NoCs from industry and academia, as shown in Table I, we observe that *Torrent* is distinctive from other P2MP mechanisms that use the NoC network layer to replicate data [9] [12] [11]: by performing replication at the endpoints and using *Chainwrite*, *Torrent*

¹**Torrent Frontend** is open-sourced as one component in *SNAX*, while **Torrent Backend** is open-sourced separately.

²In this paper, to ensure a clear terminology, the conventional P2MP data transfer with the help of a NoC router is referred to as *Multicast* and our P2MP data transfer using *Torrent* is referred to as *Chainwrite*

	Arch.	Addr. Gen	AXI Comp.	P2MP Method	Area Scaling	Open Sourced
Xilinx DMA [13]	Mono. DMA	1D	Yes	SW	N/A	No
HyperDMA [14]	Dist. DMA	ND	No	SW	N/A	No
iDMA [15]	Mono. DMA	ND	Yes	SW	N/A	Yes
XDMA [16]	Dist. DMA	ND	Yes	SW	N/A	Yes
FlexNoC [12]	NoC	N/A	Yes	Multicast	N/A	No
ESP NoC [9]	NoC	N/A	No	Multicast	O(N)	Yes
Pulp XBar [11]	XBar	N/A	Yes	Multicast	$\sim O(1)$	Yes
<i>Torrent</i>	Dist. DMA	ND	Yes	<i>Chainwrite</i>	$\sim O(1)$	Yes

TABLE I: *Torrent* comparison with SoTA DMAs and NoCs

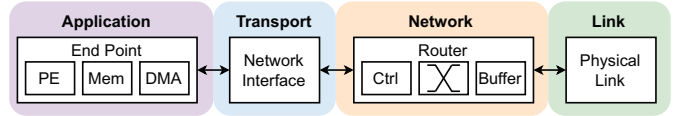


Fig. 2: Layered on-chip network architecture

achieves SoC-level data reuse without complicating the shared interconnect fabric. Furthermore, *Torrent* is highly flexible due to its distributed architecture, offering the ND affine access capability with low power and area overhead.

II. BACKGROUND

A. Network-on-Chip Architecture

Traditional hierarchical on-chip interconnects struggle to meet the performance requirements in complex multi-core SoCs, thus, NoCs with flattened structures (like mesh, torus, etc.) have become appealing substitutes [17].

NoCs can be segmented into four layers [18], referencing the OSI model (Fig. 2): (1) The application layer includes hardware components that act as the requesters and responders for I/O requests; (2) The transport layer converts transactions from the application into routable packets and delivers them to the network layer; (3) The network layer forwards the packets to the port designated in the routing table; (4) The link layer is the physical P2P link between two routers to transfer data.

B. Multicast on NoC

Similar to multicast in computer networks, the intuitive way to perform multicast on-chip is to rely on the NoC router. The multicast packet traverses the common four-stage [19] [20] router pipelines where special handling occurs at each phase: In the Route Computation (RC) phase, the head flit determines multiple output ports based on the pre-configured multicast destination set. During Virtual Channel Allocation (VA), the multicast packet requests available virtual channels for each identified output port simultaneously, which may stall if some VCs are unavailable. In Switch Allocation (SA) and Switch Traversal (ST) phases, interconnects are altered to connect multiple output ports to a single input port, duplicating the packet multiple times. The same mechanism repeats until one packet is delivered to each destination.

III. TORRENT ARCHITECTURE

Torrent proposes a novel decentralized architecture to conduct P2MP data transfers through *Chainwrite*. Unlike traditional DMAs [13] [15] that perform internal loopback of two memory requests, *Torrent* adopts a distributed DMA

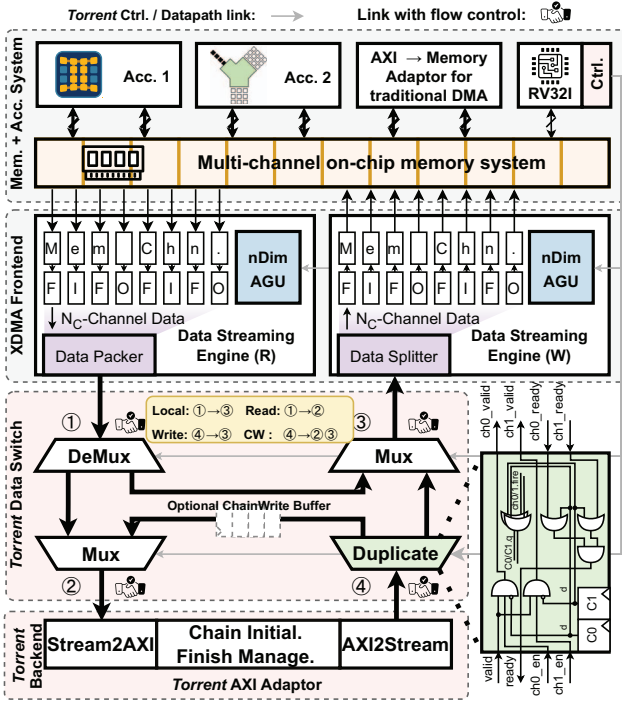


Fig. 3: The architecture of *Torrent*

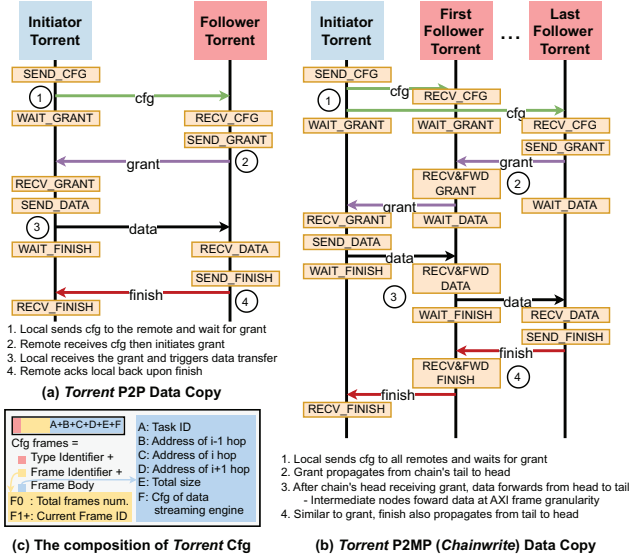


Fig. 4: *Torrent* Chainwrite four-phase orchestration

architecture: multiple *Torrents* attached to the source and destination regions are involved in one data transfer task.

The *Torrent*'s microarchitecture is shown in Fig. 3: The *Torrent* Frontend is built on the open-source XDMA framework [16] and its data streaming engine (DSE) [21], which can perform ND-affine memory accesses. This data then enters the *Torrent* data switch, which duplicates and forwards the data to different ports. Finally, the *Torrent* Backend encapsulates data into AXI requests.

A. *Torrent* Chainwrite Orchestration

Since *Torrent* adopts a **distributed orchestration strategy** to accomplish a task, we design a dedicated control flow to

avoid data corruption and deadlock. Specifically, we introduce a four-phase control flow, shown in Fig. 4 for *Torrent* P2P data copy and P2MP *Chainwrite*:

(1) **Configuration Dispatch**: When the initiator *Torrent* receives a P2P or P2MP task, it forwards the corresponding configuration settings (cfg, Fig. 4(c)) to the involved remote *Torrents*. For P2P data copy, only one remote *Torrent* is involved; therefore, one cfg is sent out; for *Chainwrite*, cfgs are forwarded to all participating *Torrents* in parallel, with each cfg describing the address of the previous node and the next node. Hence, there is no theoretical limit on the length of the chain. These cfgs specify a doubly linked list on the SoC through which data can flow from head to tail (for data) or tail to head (for control signals).

(2) **Grant Signaling Backward Propagation**: When the last *Torrent* receives the cfg packet, it generates the Grant signal and sends it backward. Every intermediate *Torrent* only forwards the **Grant** to the previous node when it is ready for this new *Chainwrite* task.

(3) **Data Transfer**: Once the initiator *Torrent* receives the Grant signal, it begins sending data into the chain. Each intermediate *Torrent* stores and forwards every received data frame to the next hop as soon as it receives it from the previous hop, such that the data finally traverses through all *Torrents*.

(4) **Finish Signaling Backward Propagation**: Similar to (2), a **Finish** signal is also generated by the last node and is propagated backward to the initiator *Torrent* to indicate that the *Chainwrite* on the chain has completed.

B. The Cross-Torrent Configuration

Torrent dispatches multi-field cfg packets (Fig. 4(c)) to orchestrate with other *Torrents*: **Type Identifier** defines whether this frame is a read/write request; **Frame Identifier** describes the total number of frames (for the first frame) and current frame ID (for the remaining frames) of a cfg packet. The cfg is split into multiple **Frame Bodies** to support interconnect with variable lengths. Each frame body contains six fields: fields A to D describe the data chain; field E is for the *Torrent* Backend to generate AXI requests with corresponding size; and field F describes the access pattern for the DSE.

C. The *Torrent* Datapath

The ***Torrent* Data Switch**, depicted in Fig. 3, forwards and/or duplicates the data for different *Torrent* working modes: Local loopback, Read, Write, and *Chainwrite*.

When the source and destination addresses are in the same memory, *Torrent* works in local loopback mode. In this case, *Torrent* is regarded as a dedicated data reshuffling accelerator, and the data flows from ① to ③. When the source and destination addresses are in different memories, then the *Torrent* attached to the source memory is in read mode (① to ②) and the one attached to the destination memory is in write mode (④ to ③). In *Chainwrite* mode, the data switch duplicates the incoming data from ④ into two copies; one copy is sent to the next hop (②) and another copy is sent to the local DSE (③). This behavior is reflected in the **RECV&FWD DATA** State

of Fig. 4(b). The stream duplicator in the data switch enables the on-the-fly data duplication with no temporary storage of data, maximizing the energy efficiency of *Chainwrite*.

The **Backend** of *Torrent* serves as the bridge between the *Torrent* Frontend and the AXI interconnect, establishing lightweight “virtual tunnels” across *Torrents* on top of the AXI.

D. Torrent Chainwrite Sequence Algorithm

Algorithm 1 Chain Write Greedy Optimization Algorithm

Require: Destinations list D , NoC dimensions noc_x, noc_y
Ensure: Ordered traverse destination list $order$

```

1:  $remaining_{dest} \leftarrow D$   $\triangleright$  Init remaining destinations
2:  $start \leftarrow \min(remaining_{dest})$   $\triangleright$  Start from dest closest to C0
3:  $order \leftarrow [start], remaining_{dest} \leftarrow remaining_{dest} \setminus \{start\}$ 
4:  $used_{path} \leftarrow XYpath((0,0), start)$   $\triangleright$  Initial path from C0
5: while  $remaining_{dest} \neq \emptyset$  do
6:    $best \leftarrow null, best_{hops} \leftarrow noc_x + noc_y$ 
7:   for each  $cand$  in  $rem\_dest$  do
8:      $path \leftarrow XYpath(order[-1], cand)$ 
9:     if no_overlap( $used_{path}, path$ ) &  $|path| < best_{hops}$  then
10:       $best \leftarrow cand, best_{hops} \leftarrow |path|, best_{path} \leftarrow path$ 
11:    end if
12:  end for
13:  if  $best = null$  then  $\triangleright$  Fallback to shortest path
14:     $best \leftarrow \operatorname{argmin}_{c \in remaining_{dest}} |path(order[-1], c)|$ 
15:     $best_{path} \leftarrow XYpath(order[-1], best)$ 
16:  end if
17:   $order.append(best)$ 
18:   $used_{path} \leftarrow used_{path} \cup best_{path}$ 
19:   $remaining_{dest} \leftarrow remaining_{dest} \setminus \{best\}$ 
20: end while
21: return  $order$ 

```

Chainwrite, in contrast to network-layer multicast, exposes destination traversal order explicitly. Our experiments demonstrate its significance for high performance. Thus, we design two complementary strategies to optimize this sequence:

(1) **Greedy heuristic (Alg. 1)** iteratively selects the next destination such that the routing path does not overlap with previously used links, while also minimizing path length. This approach balances efficiency with computational cost, making it well-suited for just-in-time optimization.

(2) The scheduling problem can be modeled as an open-path variant of the **Traveling Salesman Problem (TSP)**. Using Google OR-Tools [22], we construct distance matrices based on XY-routing and solve for the globally optimal path. This approach guarantees a global optimum with higher computational overhead, making it a candidate for ahead-of-time optimization.

IV. EVALUATION

We evaluate *Torrent* in a multi-accelerator SoC. We first compare the P2MP efficiency of *Torrent* against the SoTA P2P DMA and the network-layer multicast solution (§IV-B). We then demonstrate the effectiveness of the software *Chainwrite* sequence optimization to boost the P2MP performance (§IV-C). Then, the configuration overhead is quantified in §IV-D, showing linear overhead scaling with the number of destinations. To validate the adaptability to real workloads, we prototype

a multicore SoC on the AMD Versal™ VPK180 FPGA and evaluate its performance in the self-attention layers of DeepSeek-V3 (§IV-E). Finally, we synthesize³ a SoC with *Torrent* for area and power breakdowns (§IV-F).

A. System Evaluation Setup

In our system under evaluation, each compute cluster encompasses a 1MB, 32-bank, 64-bit-per-bank memory, two RV32I cores [23], a GeMM accelerator with 1024 8-bit MACs, and a *Torrent*. The GeMM accelerator is optimized for LLM workloads and has two operating modes: (1) multiply two matrices of 16×8 and 8×8 (for the prefill stage); (2) multiply a vector of 1×64 with a matrix of 64×16 (for the decode stage). With this, we set up a 20-cluster SoC derived from Occamy [24], interconnected by FlooNoC [25] with a 4×5 2D mesh topology. The NoC uses XY-routing and the bandwidth is 64 bytes per cycle.

Our primary comparison baseline is ESP [26], a SoC platform with an in-house NoC that supports network-layer multicast [9]. For fair comparison, the ESP SoC is configured with the same 4×5 2D mesh NoC, XY-routing, and 64 bytes/CC bandwidth. We additionally select iDMA [15] for the software-based P2MP data copy condition.

B. P2MP Copy Efficiency

We evaluate the P2MP efficiency η_{P2MP} under three approaches: (1) repeated P2P copies using iDMA, (2) Multicast in the NoC on ESP, and (3) *Chainwrite* with *Torrent*. Data sizes range from 1–128 KB and the number of destinations (N_{dst}) ranges from 2–16, yielding 192 test points.

The latencies are retrieved from hardware counters for all conditions. For iDMA, cycles equal the sum of all P2P transfers; *Torrent* measures the cycles from task dispatch to the DSE until the initiator *Torrent* receives the finish signal. Since the ESP platform only supports multicasting to accelerators, we implement dummy accelerators and integrate a hardware counter in the DMA for the latency.

Based on the measured latency⁴, we define the P2MP efficiency η_{P2MP} as follows:

$$\eta_{P2MP} = \frac{lat_{P2P,theo}}{lat_{measured}} = \frac{N_{dst} \cdot \frac{Size_{data}}{BW_{P2P,ideal}}}{lat_{measured}} \quad (1)$$

where $lat_{P2P,theo} = N_{dst} \cdot \frac{Size_{data}}{BW_{P2P,ideal}}$ is the theoretical P2P copying latency for N_{dst} destinations, and $BW_{P2P,ideal} = 64B/CC$ (system AXI bandwidth). By definition, iDMA will have $\eta_{P2MP} \leq 1$ since there is no data duplication and all data needs to be fetched from SRAM. Multicast or *Chainwrite* can reach $\eta_{P2MP} > 1$ and the ideal condition is $\eta_{P2MP} = N_{dst}$, which means all destinations receive the data at the $BW_{P2P,ideal}$ with no delay.

³Silicon synthesis targets TSMC 16FFC technology (600 MHz/0.8 V) using Synopsys Design Compiler®. Power analysis is carried out on the synthesized netlist with gate-level switching activity via Synopsys PrimeTime®.

⁴The iDMA and *Torrent* results are retrieved from RTL simulation. We conduct RTL simulation on ESP for sizes of 1-8KB and extrapolate the results for large tasks due to long simulation time.

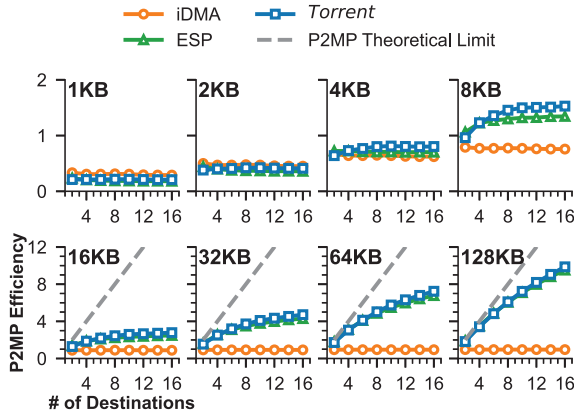


Fig. 5: The η_{P2MP} comparison for iDMA (Unicast), ESP (Multicast), and *Torrent* (*Chainwrite*)

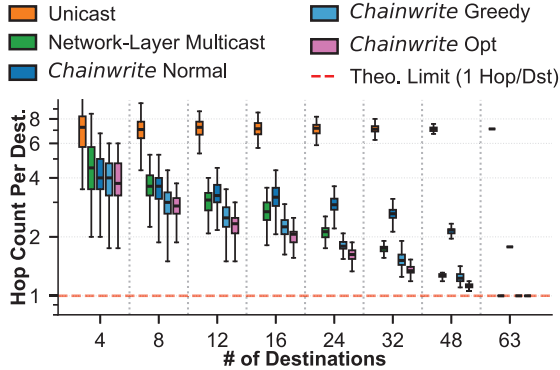


Fig. 6: The comparison of average hops per destination

Fig. 5 shows that η_{P2MP} for *Torrent* and the ESP platform approach the ideal P2MP limit with increasing data sizes and N_{dst} . At small sizes (1KB–4KB), the control overhead dominates but is amortized with larger transfers. ESP outperforms *Torrent* for few-destination scenarios (2–4) due to lower link setup overhead, but its configuration complexity grows faster with N_{dst} compared to *Torrent*. iDMA approaches $\eta_{P2MP} = 1$ when transferring larger data (>8KB), indicating near-ideal P2P link utilization. This evaluation shows *Torrent* can achieve better performance than multicast in the NoC.

C. Efficiency Impact of *Chainwrite* Sequence

Chainwrite offers greater universality and flexibility than multicast, since it does not require modifications to AXI and can write data to different addresses with varying patterns. However, it is intuitively less efficient compared to multicast, in which data can be transferred along multiple paths in parallel instead of hop-by-hop. We compare *Chainwrite* against multicast by quantifying the average number of hops per destination, an implementation-agnostic metric that reflects energy consumption and latency, across all experiments. We set up an 8×8 -cluster mesh NoC and scan N_{dst} from 4 to 63 (8 groups). Every group selects destinations randomly and repeats this 128 times (1024 test points in total). For example, one possible destination set for the $N_{dst} = 4$ test initiated by C_0 is $\{C_3, C_7, C_{21}, C_{63}\}$.

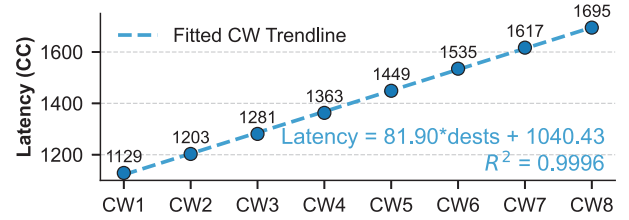


Fig. 7: The configuration overhead to *Chainwrite* 64KB data to 1-8 destinations

For unicast, the Manhattan distances between the source and all destinations are used to obtain the average number of hops. For multicast, one packet is routed following the standard XY-routing, and is divided when routes to different destinations do not overlap. For *Chainwrite*, we evaluate the two proposed scheduling algorithms introduced in §III-D against a naive ordering following cluster IDs. Average hop counts are defined as the number of edges the data traverses divided by N_{dst} .

As shown in Figure 6, as N_{dst} increases, unicast converges to the average Manhattan distance of the network while both multicast and *Chainwrite* converge to smaller values since packets exploit shared edges. Simple *Chainwrite* suffers from redundant paths and performs worse than multicast. However, *Chainwrite* with the greedy heuristic optimizer is comparable to multicast, while the TSP-based scheduling successfully helps *Chainwrite* surpass multicast at scale. For $N_{dst} = 63$, both *Chainwrite* and multicast with optimizers converge to the theoretical limit of 1 hop per destination.

D. *Chainwrite* Configuration Overhead

As discussed in §III-A, *Chainwrite* introduces extra overhead due to the four-phase handshaking through the chain. To quantify the relationship between the overhead and the number of destinations, we measure latency for a 64KB data copy using *Chainwrite* with 1 to 8 destinations.

Fig. 7 shows that *Chainwrite* setup incurs an overhead of 82 cycles per additional destination, following a linear trend. Although the latency value depends on the underlying NoC implementation, the scaling behavior remains consistent.

E. DeepSeek-V3 Self-attention Layers Data Movement Evaluation on FPGA

We implement the SoC with 3×3 clusters on an AMD Versal™ VPK180 FPGA. Among them, C_0 is a full cluster with the GeMM accelerator, *Torrent*, and XDMA; and the others are GeMM-less clusters due to FPGA resource constraints. Fig. 8 details the annotated floorplan, frequency, and resource utilization of the FPGA.

We extract three workloads from the self-attention layers of Deepseek-V3 [27] and the required data layouts from the architecture of the GeMM accelerator as shown in Table II: (1) When calculating $Q \cdot K^T$ at the prefill stage, the K matrix (MNM16N8 layout) is the output of the previous matrix multiplication. Since the Q matrix is large and will be tiled to multiple accelerators, the K matrix needs to be copied to multiple accelerators; (2) The output of $Q \cdot K^T$

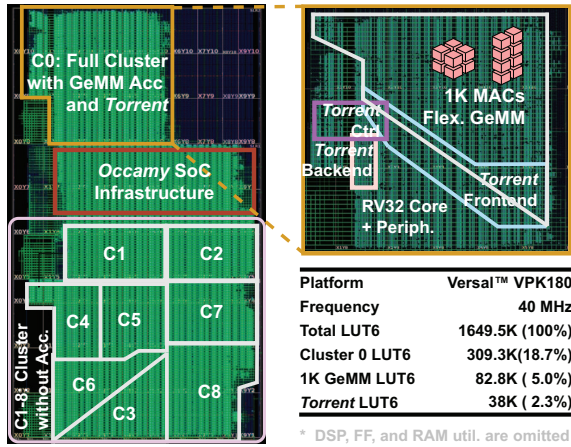


Fig. 8: The FPGA implementation result of a 9-cluster SoC

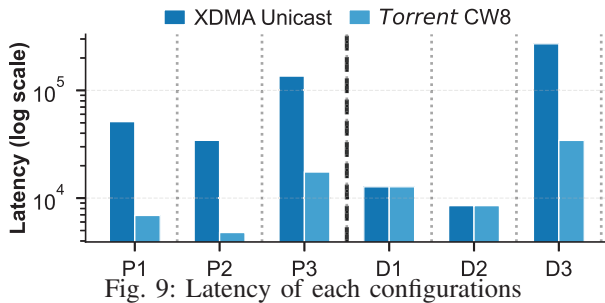


Fig. 9: Latency of each configurations

Experiment Setup	Shape	I/O Layout	Multicast
P1:QKT_Single_Head	2048×192	MNM16N8/ MNM8N8	TRUE
P2:SV_Single_Head	2048×128	MNM16N8/ MNM8N8	TRUE
P3:KV_Matrix_MLA_Recovery	2048×512	MNM16N8/ MNM16N8	TRUE
D1:QKT_Single_Head	4096×192	MNM16N8/ MNM64N16	FALSE
D2:SV_Single_Head	4096×128	MNM16N8/ MNM64N16	FALSE
D3:KV_Matrix_MLA_Recovery	4096×512	MNM16N8/ MNM16N8	TRUE

TABLE II: Configuration Details

Fig. 10: DeepSeekV3 Evaluation: Performance results (top) and configuration details (bottom)

is then multiplied with the V matrix. Thus, we need to transform its data layout and copy it to multiple destinations, similar to condition (1); (3) The recovery of the KV matrix requires copying the KV -cache to all accelerators, but no layout transformation is required. The same workloads are also evaluated at the decode stage.

Fig. 9 shows that *Torrent*'s *Chainwrite* P2MP data copying achieves up to $7.88\times$ speedup over XDMA [16] thanks to its highly-efficient ND affine access and *Chainwrite* capability. Overall, *Torrent* provides order-of-magnitude speedups in data copying operations of the self-attention layers, especially when P2MP and layout transformations are required.

F. ASIC Synthesis and Power Analysis

Finally, we synthesize the SoC in a 16 nm silicon technology, including 4 clusters with a 256KB scratchpad per cluster,

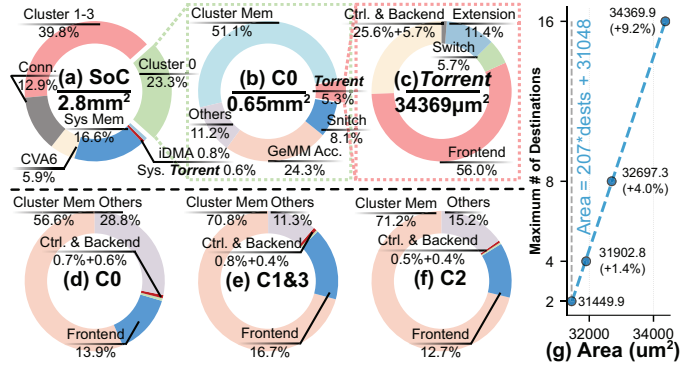


Fig. 11: The area breakdown of (a) the 4-cluster SoC, (b) the accelerator cluster, and (c) the *Torrent*; the power breakdown of (d) the initiator cluster, (e) the follower cluster at the middle of the chain, and (f) the follower cluster at the tail of the chain; (g) the relationship between the area and maximum number of destinations for the initiator *Torrent*

and 512KB global memory. The first cluster is the fully-featured cluster with a GeMM accelerator, while the accelerator is removed from the remaining clusters. Five *Torrents* are instantiated (one attached to the global SRAM and one for each cluster) in the SoC. We conduct post-synthesis simulations with a 64KB, 3-destination *Chainwrite* from cluster 0. We also synthesize *Torrent* with different $N_{dst,max}$.

1) *Area Analysis*: Fig. 11(a) shows the SoC consumes $2.8mm^2$ of silicon area. Among them, the CVA6 [28] core consumes 5.9%, cluster 0 consumes 23.3%, and the global SRAM consumes 16.6% of the area. At the cluster scope (Fig. 11(b)), *Torrent* consumes 5.3% of the area, only equivalent to $\sim 1/5$ of the size of the GeMM accelerator. *Torrent* has a comparable size to iDMA while providing efficient N-dimensional data copying and *Chainwrite*. At the SoC scope, the *Torrent* attached to global memory barely occupies 0.6% of the area. Thanks to the *Chainwrite* architecture, the total *Torrent* area scales with a nearly constant hardware overhead of 0.65% additional area per destination. (Fig. 11(g))

2) *Power Analysis*: The power of the initiator cluster is 175.7 mW (Fig. 11(d)). Among the follower *Torrents*, the ones in the middle consume more power (Fig. 11(e)) because they need to forward the data to the next hop. Our SoC reaches 4.68 pJ/B/hop energy efficiency.

V. CONCLUSION

We presented *Torrent*, a distributed DMA architecture that enables flexible application-layer multicast for large-scale SoCs. The FPGA implementation of *Torrent*-enhanced SoC demonstrates up to $7.88\times$ speedup over unicast baseline. Compared to network-layer multicast, *Torrent* offers full AXI compatibility and greater flexibility. Software scheduling further improves efficiency and ensures full NoC bandwidth utilization. Hardware results confirm minimal overhead of 1.2% in SoC area and 2.3% in system power. *Torrent* delivers scalable P2MP data transfers with a small cycle overhead of 82CC and area overhead of $207 \mu m^2$ per destination.

REFERENCES

- [1] M. Bohr, "A 30 year retrospective on dennard's mosfet scaling paper," *IEEE Solid-State Circuits Society Newsletter*, vol. 12, no. 1, pp. 11–13, 2009.
- [2] P. A. Hager, B. Moons, S. Cosemans, I. A. Papistas, B. Rooseleer, J. Van Loon, R. Uytterhoeven, F. Zaruba, S. Koumoussi, M. Stanisavljevic *et al.*, "11.3 metis aipu: A 12nm 15tops/w 209.6 tops soc for cost-and energy-efficient inference at the edge," in *2024 IEEE International Solid-State Circuits Conference (ISSCC)*, vol. 67. IEEE, 2024, pp. 212–214.
- [3] V. Schmulbach, J. Kim, E. Gao, N. Jha, E. Wu, O. Yu, B. Oliveau, X. Kong, B. Roberts, C. McMahon *et al.*, "Nectar and rasoc: Tale of two class socs for language model inference and robotics in intel 16," in *2024 IEEE Hot Chips 36 Symposium (HCS)*. IEEE Computer Society, 2024, pp. 1–1.
- [4] A. Gholami, Z. Yao, S. Kim, C. Hooper, M. W. Mahoney, and K. Keutzer, "Ai and memory wall," *IEEE Micro*, vol. 44, no. 3, pp. 33–39, 2024.
- [5] V. Sze, Y.-H. Chen, T.-J. Yang, and J. S. Emer, *Efficient processing of deep neural networks*. Springer, 2020.
- [6] H. Genc, S. Kim, A. Amid, A. Haj-Ali, V. Iyer, P. Prakash, J. Zhao, D. Grubb, H. Liew, H. Mao, A. Ou, C. Schmidt, S. Steffl, J. Wright, I. Stoica, J. Ragan-Kelley, K. Asanovic, B. Nikolic, and Y. S. Shao, "Gemmini: Enabling systematic deep-learning architecture evaluation via full-stack integration," in *Proceedings of the 58th Annual Design Automation Conference (DAC)*, 2021.
- [7] H. Xu, Y. Gui, and L. M. Ni, "Optimal software multicast in wormhole-routed multistage networks," *IEEE Transactions on Parallel and Distributed Systems*, vol. 8, no. 6, pp. 597–607, 1997.
- [8] J.-W. Jang, S. Lee, D. Kim, H. Park, A. S. Ardestani, Y. Choi, C. Kim, Y. Kim, H. Yu, H. Abdel-Aziz *et al.*, "Sparsity-aware and re-configurable npu architecture for samsung flagship mobile soc," in *2021 ACM/IEEE 48th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 2021, pp. 15–28.
- [9] J. Zuckerman, J.-D. Wellman, A. Vanamali, M. Shankar, G. Tombesi, K. Swaminathan, K. Lee, M. Kapur, R. Philhower, P. Bose *et al.*, "Towards generalized on-chip communication for programmable accelerators in heterogeneous architectures," *arXiv preprint arXiv:2407.04182*, 2024.
- [10] ARM, "Amba axi and ace protocol specification," <https://developer.arm.com/documentation/ih10022/hc/?lang=en>, 2021.
- [11] L. Colagrande and L. Benini, "A multicast-capable axi crossbar for many-core machine learning accelerators," *arXiv preprint arXiv:2502.19215*, 2025.
- [12] "FlexNoC Interconnect IP - Arteris — arteris.com," <https://www.arteris.com/products/non-coherent-interconnect-ip/flexnoc/>, [Accessed 29-08-2025].
- [13] A. Xilinx, "Axi dma logicore ip product guide," 2022.
- [14] M. Peng, H. Chen, Y. Zhang, and S. Liu, "Hyperdma: Enhancing high-performance computing and ai workflows with advanced data transfer capabilities," in *2024 9th International Conference on Integrated Circuits and Microsystems (ICICM)*. IEEE, 2024, pp. 636–644.
- [15] T. Benz, M. Rogenmoser, P. Scheffler, S. Riedel, A. Ottaviano, A. Kurth, T. Hoefler, and L. Benini, "A high-performance, energy-efficient modular dma engine architecture," *IEEE Transactions on Computers*, vol. 73, no. 1, pp. 263–277, 2023.
- [16] F. Kong, Y. Deng, X. Yi, R. Antonio, and M. Verhelst, "Xdma: A distributed, extensible dma architecture for layout-flexible data movements in heterogeneous multi-accelerator socs," 2025. [Online]. Available: <https://arxiv.org/abs/2508.08396>
- [17] L. Benini and G. De Micheli, "Networks on chips: A new soc paradigm," *computer*, vol. 35, no. 1, pp. 70–78, 2002.
- [18] T. Bjerregaard and S. Mahadevan, "A survey of research and practices of network-on-chip," *ACM Computing Surveys (CSUR)*, vol. 38, no. 1, pp. 1–es, 2006.
- [19] F. A. Samman, T. Hollstein, and M. Glesner, "Multicast parallel pipeline router architecture for network-on-chip," in *Proceedings of the conference on Design, automation and test in Europe*, 2008, pp. 1396–1401.
- [20] J. Zhao, A. Agrawal, B. Nikolic, and K. Asanović, "Constellation: An open-source soc-capable noc generator," in *2022 15th IEEE/ACM International Workshop on Network on Chip Architectures (NoCArc)*. IEEE, 2022, pp. 1–7.
- [21] X. Yi, Y. Deng, R. Antonio, F. Kong, G. Paim, and M. Verhelst, "Datamaestro: A versatile and efficient data streaming engine bringing decoupled memory access to dataflow accelerators," *arXiv preprint arXiv:2504.14091*, 2025.
- [22] L. Perron and F. Didier, "Cp-sat," Google. [Online]. Available: https://developers.google.com/optimization/cp/cp_solver/
- [23] F. Zaruba, F. Schuiki, T. Hoefler, and L. Benini, "Snitch: A tiny pseudo dual-issue processor for area and energy efficient execution of floating-point intensive workloads," *IEEE Transactions on Computers*, vol. 70, no. 11, pp. 1845–1860, 2020.
- [24] G. Paulin, P. Scheffler, T. Benz, M. Cavalcante, T. Fischer, M. Eggimann, Y. Zhang, N. Wistoff, L. Bertaccini, L. Colagrande *et al.*, "Occamy: A 432-core 28.1 dp-gflop/s/w 83% fpu utilization dual-chiplet, dual-hbm2e risc-v-based accelerator for stencil and sparse linear algebra computations with 8-to-64-bit floating-point support in 12nm finfet," in *2024 IEEE Symposium on VLSI Technology and Circuits (VLSI Technology and Circuits)*. IEEE, 2024, pp. 1–2.
- [25] T. Fischer, M. Rogenmoser, T. Benz, F. K. Gürkaynak, and L. Benini, "Floonoc: A 645-gb/s/link 0.15-pj/b/hop open-source noc with wide physical links and end-to-end axi4 parallel multistream support," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 2025.
- [26] M. C. Dos Santos, T. Jia, J. Zuckerman, M. Cochet, D. Giri, E. J. Loscalzo, K. Swaminathan, T. Tambe, J. J. Zhang, A. Buyuktosunoglu *et al.*, "14.5 a 12nm linux-smp-capable risc-v soc with 14 accelerator types, distributed hardware power management and flexible noc-based data orchestration," in *2024 IEEE International Solid-State Circuits Conference (ISSCC)*, vol. 67. IEEE, 2024, pp. 262–264.
- [27] A. Liu, B. Feng, B. Xue, B. Wang, B. Wu, C. Lu, C. Zhao, C. Deng, C. Zhang, C. Ruan *et al.*, "Deepseek-v3 technical report," *arXiv preprint arXiv:2412.19437*, 2024.
- [28] F. Zaruba and L. Benini, "The cost of application-class processing: Energy and performance analysis of a linux-ready 1.7-ghz 64-bit risc-v core in 22-nm fdsoi technology," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 27, no. 11, pp. 2629–2640, Nov 2019.