# Hybrid Systolic Array Accelerator with Optimized Dataflow for Edge Large Language Model Inference

Chun-Ting Chen, HanGyeol Mun, Jian Meng, Mohamed S. Abdelfattah, and Jae-sun Seo

*School of Electrical and Computer Engineering, Cornell University*

{cc2793, hm632, jm2787, mohamed, js3528}@cornell.edu

*Abstract*—Edge inference for large language models (LLM) offers secure, low-latency, and cost-effective inference solutions. We emphasize that an edge accelerator should achieve high area efficiency and minimize external memory access (EMA) during the memory-bound decode stage, while maintaining high energy efficiency during the compute-intensive prefill stage. This paper proposes an edge LLM inference accelerator featuring a hybrid systolic array (HSA) architecture that optimizes inference efficiency in both stages. To further reduce EMA, we adopt MXINT4 weight quantization and propose an optimized dataflow tailored for HSA, ensuring negligible dequantization overhead and achieving 100% hardware utilization with minimal accuracy loss under edge DRAM bandwidth constraints. For non-linear operations, we incorporate optimized root mean square normalization (RMSNorm) and rotary position embedding (RoPE) units, reducing their latency, area, and memory access overhead while enabling end-to-end inference on our accelerator. Our solution achieves 247/117 (token/s/mm$^2$) while running a 1.3B LLM on long-input/long-output scenarios, providing >2.45×/13.5× improvement over existing approaches, while maintaining superior energy efficiency in token generation.

## I. INTRODUCTION

The demand for LLM inference continues to grow as a wide range of new applications are being developed. Compared to data centers, edge AI has emerged as a popular alternative, offering secure, low-power, and cost-effective inference capabilities. However, the high memory and compute demands of LLMs pose significant challenges for edge deployment.

First, edge devices have limited memory capacity, which restricts the size of deployable LLMs. Second, power consumption is a critical reliability concern, unlike in data centers where it is primarily an optimization problem. Furthermore, data centers typically benefit from batching multiple input queries to address memory bottlenecks during the decoding stage, which improves inference efficiency. In contrast, edge scenarios often involve user queries arriving one at a time, which limits the accelerator's performance due to the narrow memory bandwidth available on edge devices.

Fig. 1(a) illustrates the autoregressive inference paradigm, where the prefill stage involves Matrix-Matrix-Multiplication (MMM), while decode comprises mainly Matrix-Vector-Multiplication (MVM) workloads. Overall, the edge inference workloads consist of two typical scenarios of **1) Long Input Short Output (LISO)** (e.g., summarizing emails and extract-
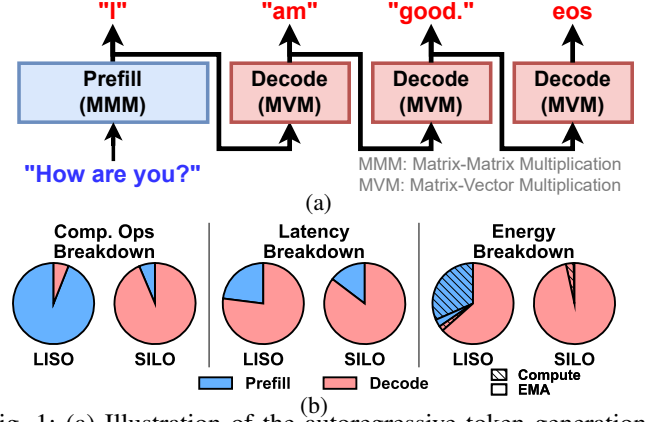
Fig. 1: (a) Illustration of the autoregressive token generation. (b) Breakdown of computational operations, latency, and energy for RetNet 1.3B model under LISO/SILO scenarios.

ing insights from long news articles) and **2) Short Input Long Output (SILO)** (e.g., text generation such as email templates).

For the LISO scenario, we assume input/output tokens are 750/50, while for the SILO scenario, they are 50/750. These numbers are based on the average length of an email or news article, which align with the average token statistics found in the news summarization dataset [16]. Although in practice many input and output token lengths vary, these two settings illustrate key computational differences between prefill and decode stages under long input/output scenarios, and how edge accelerators should be optimized accordingly.

We profile the computational operations, latency, and power consumption for both LISO and SILO scenarios. We run Ret-Net [23] 1.3B LLM model on Nvidia Jetson Orin Nano [17] as reference edge inference accelerator. As shown in Fig. 1(b), The commercial mobile GPU achieves peak performance of 40 TOPS, but the low memory bandwidth (68 GB/s with LPDDR5) limits the speed of decoding (only 1.7% of peak performance), which further reveals the following observations:

1) The decode stage dominates overall latency, contributing to over 80% of the runtime in both scenarios, even when input length is significantly larger than output.
2) EMA dominates energy consumption. Nonetheless, the computation energy consumption in the prefill stage cannot be neglected in LISO scenarios.

These observations are consistent with results obtained using other LLMs, such as Llama 2 [24]. Notably, edge accelerators operate under memory-bound conditions for most
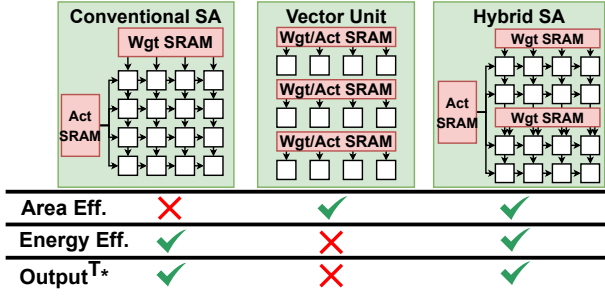
Fig. 2: Comparison of three MAC architectures, highlighting the **efficiency balance challenge**. The proposed hybrid SA combines the efficiency of both SA and vector unit architectures. *Output$^T$ indicates support for transpose output.



Fig. 3: Normalized latency/energy comparison of Llama-2 7B and RetNet 6.7B. Profiles made with Jetson Orin Nano as a reference edge accelerator. (Energy assumption: MAC=0.5pJ/Byte [17], DRAM access=32pJ/Byte [2], [18].)

of the runtime, resulting in low utilization (e.g., ~1.7% in Jetson Orin Nano's case). Additionally, despite the low latency of the prefill stage, its energy consumption cannot be neglected in scenarios involving long inputs. Here we conclude that edge accelerators should achieve high energy efficiency during the prefill stage, while reducing EMA and maintaining area efficiency during the decode stage.

Recent works [3], [8], [9], [19], [20] report high peak performance and energy efficiency, but exhibit inefficiency when running end-to-end workloads on edge devices due to low data reuse and low utilization. Specifically, [3], [8], [9] use vector unit architecture for MAC computation, which offers flexibility to support both MMM and MVM workloads, but it fails to reuse data during prefill and leads to low energy-efficiency. In contrast, [19], [20] employ conventional 2D systolic array (SA) architecture, which fails to support MVM workload during decode since the batch size is only one in the edge scenario, resulting in low utilization. In summary, prior architectures fail to balance energy efficiency during prefill and area efficiency during decode, as illustrated in Fig. 2.

Moreover, most of the aforementioned papers adopt INT8 as the weight format, which suffers from severe memory access bottlenecks. While [9] employs INT4 per-vector quantization, it fails to preserve accuracy in modern LLMs. Alternatively, [20] introduces weight partitioning to alleviate memory bottlenecks but incurs large power and area overhead.

In this work, we present an area- and energy-efficient accelerator for edge inference, which is implemented in TSMC 28nm CMOS. We deploy RetNet 1.3B [23] as the target LLM to reduce memory access and eliminate the latency overhead of softmax operations. Our key contributions to address the shortcomings of prior works are as follows:

- **Hybrid Systolic Architecture (HSA)**: We introduce a hybrid systolic architecture that combines the strengths of conventional SA and vector units, supporting power-efficient prefilling while maintaining high hardware utilization for decoding, achieving high area efficiency.
- **Low-access High-utilization Dataflow**: During decode, we adopt MXINT4 [21] weight quantization scheme to reduce memory access. We propose a 100% utilization rate dataflow mapped to HSA under 500MHz by using
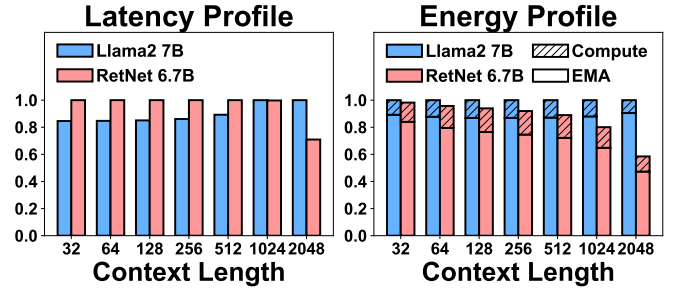
multiple PEs for dequantization.
- **Overhead-optimized RMSNorm and RoPE Units**: We introduce a layer-fused RMSNorm to eliminate latency overhead (5–10% of total) and remove a 32kB buffer. For the RoPE unit, we reuse embedding multipliers to compute sin/cos values online, reducing DRAM access.

Overall, our accelerator achieves 247/117 token/s/mm$^2$ in LISO/SILO scenarios, which represents over 2.45×/13.5× improvement compared to existing accelerator solutions while maintaining competitive energy efficiency.

## II. BACKGROUND

Retentive Network (RetNet) [23] is a recently developed, attention-free LLM. Compared to transformer-based models, RetNet replaces the softmax operation with a decaying causal mask $D$, transforming the attention block into a retention block. This modification transforms RetNet into a Structured State Space Duality model [5], providing computationally efficiency during training (similar to transformers) and enabling linear memory complexity during decoding.

RetNet achieves a desirable memory complexity of O(1) during decoding, making it particularly well-suited for edge scenarios where memory bandwidth and capacity are constrained, and users can tolerate accuracy trade-offs. The energy profile in Fig. 3 highlights RetNet's low memory footprint. In this profiling, we select RetNet 6.7B to match the model size of Llama 2 7B for a fair comparison.

While we chose RetNet as the target model in this work, we emphasize that all modern LLMs, including commonly used models like Llama [6], [24] or SLMs [1], employ auto-regressive token generation. Thus, the MMM (prefill) and MVM (decode) workloads remain consistent across models, all of which face the same challenges discussed in Sec. I.

## III. QUANTIZATION ALGORITHM

In this work, we quantize the target RetNet [23] and Llama [6], [24] models with 8-bit activation. The model weights are quantized to 4-bit with MXINT4 representation to reduce EMA by half during decode. We first adopt SmoothQuant [26] and compress the activation precision down to INT8. To save the bandwidth between the off-chip memory
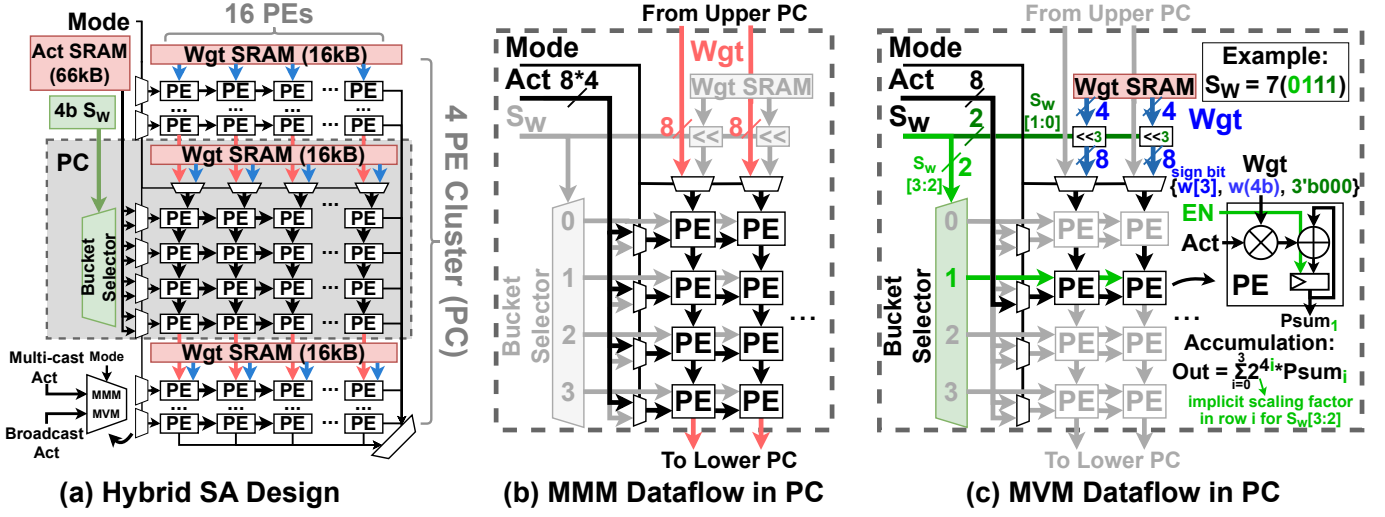
**(a) Hybrid SA Design**  **(b) MMM Dataflow in PC**  **(c) MVM Dataflow in PC**

Fig. 4: (a) **Left:** Proposed hybrid SA architecture. There are in total of 4 PE clusters arranged vertically with a unified activation SRAM. (b) **Middle:** MMM dataflow illustration, both Act and Wgt use INT8 format. (c) **Right:** MVM dataflow illustration, where each PC works independently in MVM since there is no weight reuse. MXINT4 format is used for Wgt.

and processing element, we follow the trend of recent state-of-the-art (SoTA) quantization algorithms [7], [12] by quantizing the model weights into 4-bit in a group-wise fashion. Recent SoTA methods [7], [12], [28] employ a floating-point 16-bit (FP16) scaling factor as the "dequantization scaling factor" between the stored weights (e.g., 4-bit) and the actual weights for computation (e.g., 8-bit). However, performing the high-precision FP16 dequantization on the fly is expensive due to the fine-grained (group-wise) mixed precision operation.

In this work, given the pre-trained weight $W_{FP16}$, we adopted the scaling format from micro-scaling (MX) [21] by characterizing the group-wise ($W_g$) difference via a shifting-based scaling ($S_g$). Given the output channel dimension $C$ and group size $g$ ($g < C$), the group-wise shifting factors are:

$$S_g = \text{floor}\big(\log_2(\max|W_g|)\big) \quad (1)$$

During post-training quantization (PTQ), the weights $W_{FP16}$ are quantized down to 4-bit after applying the shifting (MX-INT4). The amount of shifting is constrained between [-9, +5] to prevent overflow. On top of that, the quantization scaling factor $S_w$ remains tensor-wise, which can be fused together with the group-wise shifter ($S_g$). During prefilling, we employ INT8 weights and continue the subsequent decoding with MXINT4 format only. Unlike prior research works [7], [12], [28] where the shifting factors are separately stored as group-wise FP16 scaling factors, the shifting-based scaling is **naturally more compatible with hardware** compared to FP16 scaling [12], [28], while preserving minimal performance drop, as presented in Section V. Additionally, we further reduce shifting overhead by leveraging multiple PEs, taking advantage of the inherently low hardware utilization during the decode stage, as shown in Section IV-A2.

## IV. HARDWARE ACCELERATOR DESIGN

Our accelerator consists of two main components: the Hybrid Systolic Array (HSA) and Post-Processing Unit (PPU).

The HSA enables low-latency and energy-efficient computation, while PPU supports various post-processing functions to achieve end-to-end inference.

### A. Hybrid Systolic Array (HSA)

The hybrid systolic array (HSA), shown in Fig. 4(a), is designed to efficiently handle both MMM and MVM workloads. It comprises four PE clusters (PC), each with one weight SRAM, one bucket selector, 16 4-bit shifters, and 64 PEs arranged in four rows, totaling 256 PEs. All PCs share a single activation SRAM, sized to accommodate activation with a maximum dimension of 4,096 with batch size 16. While prior designs [3], [9], [20] incorporate more than 1,024 MAC units, we opt for 256 PEs in our accelerator due to the memory-bound nature of edge inference. As highlighted in [27], increasing compute capability significantly helps prefill but has minimal impact and degrades utilization ratio on decoding.

The 256 PEs are arranged as a 16×16 2-D array, enabling activation and weight reuse in MMM workloads, similar to conventional SA. Bucket selectors and shifters are designed to incorporate MXINT4 format with 4-bit weight scaling in MVM dataflow. Furthermore, HSA supports transpose tensor manipulation by draining output vertically or horizontally. Here, we elaborate on how MMM/MVM dataflows operate on HSA architecture.

*1) MMM Dataflow:* Fig. 4(b) shows the dataflow of MMM workload. As a compute-bound workload, only the weight SRAM in the first PC will be activated. Activations and weights are loaded and computed in INT8 format. Each PE adopts output stationary dataflow and is drained together at the end of operation. The draining phase consumes 16 cycles, which is negligible compared to the latency of computing large matrix multiplication. Output will be drained vertically when transposed output is required, otherwise horizontally.
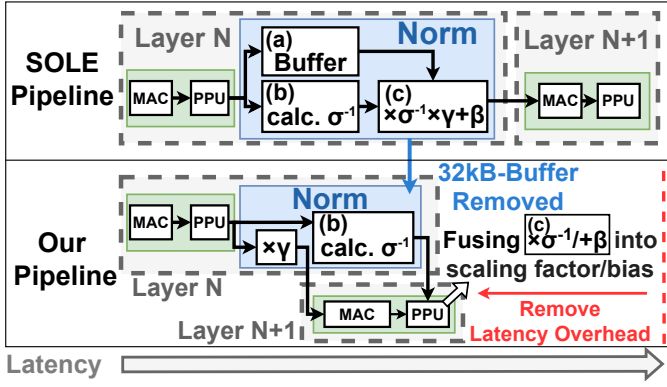
Fig. 5: Comparison of normalization pipeline between SOLE [25] and ours.

*2) MVM Dataflow:* Fig. 4(c) shows the dataflow of MVM workload. During decode stage, as there is no reuse for weights, each PC works independently. To reduce memory access and improve hardware utilization, activations are still loaded in 8-bit while weights are loaded in 4-bit MXINT4 format. Inspired by [13], we map the dequantization of MX-INT4 onto the HSA. As mentioned in Section III, we design our 4-bit scaling factor $S_w$ for MXINT4 in output dimension, which allows row-wise reuse. The 4-bit weights will first be shifted by two LSBs of the scaling factor ($S_w[1:0]$), resulting in 8-bit weights, and will then be transferred vertically to PEs. Two MSBs of the scaling factor ($S_w[3:2]$) activate one row of the PEs and govern accumulation of the product of 8-bit activation and 8-bit weight by clock-gating with an enable signal. At the end, output will be drained vertically and four PEs' partial sum will be accumulated together. Note that in MVM dataflow, we reuse the inter-PE connections in MMM dataflow to reuse existing hardware resources. The only overhead added to support MVM dataflow is the bucket selector and the shifters that cost 0.4% of the total area.

### B. Post-Processing Unit (PPU)

Our post-processing unit (PPU) supports several non-linear functions, such as activation functions and quantization, after the MAC is conducted to enable end-to-end inference. We particularly optimized RMSNorm and RoPE units, which have larger area, latency, and energy overhead.

*1) Root Mean Square Normalization (RMSNorm):* LLM models adopt RMSNorm over standard normalization to eliminate the latency and area overhead caused by mean computation, which further simplifies the design of AI accelerators.

Prior works [25] implemented a dedicated module for standard normalization which contains (a) input buffer, (b) $\sigma^{-1}$ calculation, and (c) normalization with affine transform. Since our focus is solely on RMSNorm, we propose a fusion scheme that reduces area overhead by removing the need of the input buffer, while eliminating latency overhead by parallelly calculating $\sigma^{-1}$ and the next layer.

Eq. (2) describes a quantized matrix multiplication layer $n$ where $Y_n$, $X_n$, $W_n$, $S_n$ denote the output activation, input activation, weight, and quantization scaling factor respectively.

Eq. (3) describes the RMSNorm calculation, where $\sigma_{Y_n}$, $\gamma$, and $\beta$ are the standard deviation and learned parameters. Note that $\gamma$ and $\beta$ are vectors of the same dimension as $Y_n$.

$$Y_n = (X_n \cdot W_n) \cdot S_n \tag{2}$$

$$RMSNorm(Y_n) = \frac{Y_n}{\sigma_{Y_n}} * \gamma + \beta \tag{3}$$

Instead of performing the full RMSNorm on $Y_n$ in layer $n$, we apply only $*\gamma$ and fuse $*\sigma_{Y_n}^{-1}$ and $+\beta$ with the scaling factor $S_{n+1}$ and bias $B_{n+1}$ in layer $n+1$ as shown in Eq. (4). Bold terms are computed online. This fusion preserves mathematical equivalence and enables pipelining. $\beta$ is included for generality, although it is often omitted in modern LLMs.

$$\begin{aligned}
\mathbf{X_{n+1}} &= RMSNorm(\mathbf{Y_n}) \\
&= \mathbf{Y_n} * \sigma_{\mathbf{Y_n}}^{-1} * \gamma + \beta \\
\mathbf{Y_{n+1}} &= (W_{n+1} * \mathbf{X_{n+1}}) * S_{n+1} \\
&= W_{n+1} * (\mathbf{Y_n} * \sigma_{\mathbf{Y_n}}^{-1} * \gamma + \beta) * S_{n+1} \\
&= W_{n+1} * \mathbf{Y_n} * \gamma * \sigma_{\mathbf{Y_n}}^{-1} * S_{n+1} + W_{n+1} * \beta \\
&\quad * S_{n+1} \\
&= W_{n+1} * \mathbf{Y_n^*} * \mathbf{S_{n+1}^*} + B_{n+1} \\
\text{where } & \mathbf{Y_n^*} = \mathbf{Y_n} * \gamma, \quad \mathbf{S_{n+1}^*} = \sigma_{\mathbf{Y_n}}^{-1} * S_{n+1}, \\
& B_{n+1} = W_{n+1} * \beta * S_{n+1}
\end{aligned} \tag{4}$$

Fig. 5 illustrates the overhead optimization compared to SOLE. We effectively remove the (a) 32kB buffer (assuming embedding dimension of 4,096, batch size of 16) in RMSNorm unit and eliminate its latency overhead (5-10% of total latency) by pipelining its operations (b) with the subsequent layer's MAC operations. The (b) calculation of $\sigma^{-1}$ remains the same, which involves square accumulation and square root operation.

*2) Rotary Position Embedding (RoPE):* RoPE [22] is widely adopted to incorporates token position information in LLMs. In RoPE, for the $m^{th}$ token, its embeddings $x_n$ and $x_{n+1}$ are rotated by an angle of $m\theta_{n/2}$. The embedding computation is shown in Eq. (5), where $d$ represents the embedding dimension per head.

$$\begin{aligned}
RoPE&(x_n, x_{n+1}) = \\
&\begin{pmatrix} x_n \\ x_{n+1} \end{pmatrix} \otimes \begin{pmatrix} \cos m\theta_{n/2} \\ \cos m\theta_{n/2} \end{pmatrix} + \begin{pmatrix} -x_{n+1} \\ x_n \end{pmatrix} \otimes \begin{pmatrix} \sin m\theta_{n/2} \\ \sin m\theta_{n/2} \end{pmatrix} \\
&\theta_i = 10000^{-2(i-1)/d}, i \in [1, 2, ..., d/2]
\end{aligned} \tag{5}$$

Naïve implementations of RoPE either involve loading pre-computed sine/cosine values as weights, or compute them on the accelerator. Consider the memory-bound nature of LLM inference and the high cost of on-chip CORDIC module, our RoPE unit reuses the existing embedding multipliers and adders to compute sine and cosine values for the next token with trigonometric identities (Eq. (6)). This approach reuses the hardware resources and also eliminates the need for repeated DRAM access.

$$\begin{aligned}
sin((m+1)\theta_i) &= sin(m\theta_i)cos(\theta_i) + sin(\theta_i)cos(m\theta_i) \\
cos((m+1)\theta_i) &= cos(m\theta_i)cos(\theta_i) - sin(m\theta_i)sin(\theta_i)
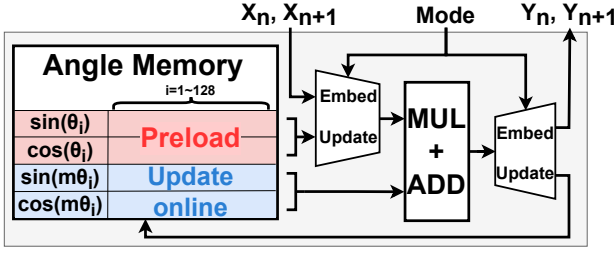\end{aligned} \tag{6}$$

Fig. 6: Proposed RoPE unit. Angle memory stores sin, cos values across head dimension $d$. RoPE has two modes of operation: **Embed** mode and **Update** mode.
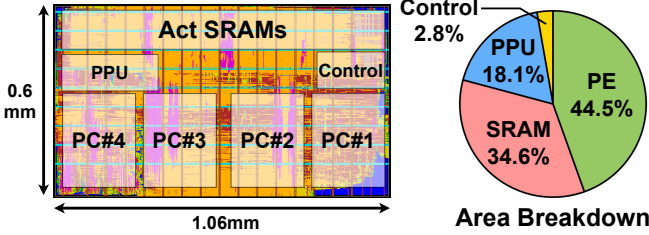


Fig. 7: 28nm accelerator layout and area breakdown.

Fig. 6 illustrates the details of the RoPE unit. While performing rotary position embedding for token $m$, the unit operates in **Embed** mode, where token embeddings $x_n, x_{n+1}$ are embedded with $sin(m\theta_i)$ and $cos(m\theta_i)$. Subsequently, the unit switches to **Update** mode to calculate $sin((m+1)\theta_i)$ and $cos((m+1)\theta_i)$ using trigonometric identities, preparing for the the next token $m+1$. The RoPE unit acquires $\sim$4% of the total area of our accelerator.

## V. EXPERIMENTAL RESULTS

Our accelerator is prototyped in TSMC 28nm CMOS and is evaluated with post-layout simulation results. The design is implemented in SystemC code following the high level synthesis (HLS) design flow [10] with Siemens Catapult. For gate-level design, Synopsys Design Compiler is used to synthesize RTL code from Catapult, where clock gating is added to the design. Cadence Innovus is used as the place-and-route tool. To obtain the power and latency results with post-layout simulation, we use Synopsys VCS and Primetime. Fig. 7 shows the 28nm accelerator layout and area breakdown.

### A. Architecture Comparison

To highlight HSA's advantages over conventional SA and vector unit, we compare their end-to-end energy and latency while running RetNet 1.3B under LISO and SILO scenarios. HSA is evaluated using post-layout simulation, while [9] scaled to 28nm as the vector unit baseline. Conventional SA performance is estimated based on HSA, as it introduces minimal overhead. Note that as both prefill and decode stages are evaluated, "token" refers to total tokens, including prompt and output, in token-related metrics like tokens/s.

Fig. 8 presents energy and latency results for the prefill and decode stages. We report only prefill energy under LISO scenario and decode latency under SILO scenario for two
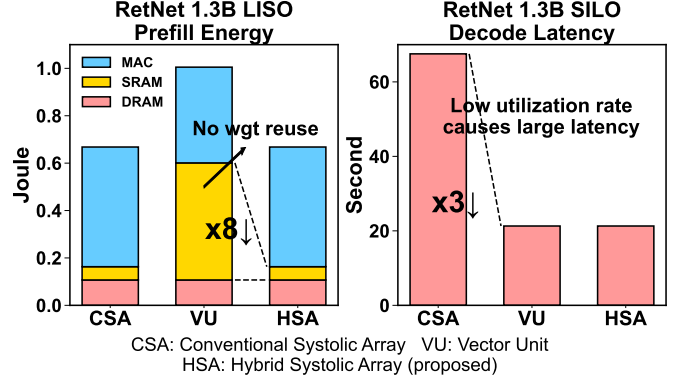


Fig. 8: Energy and latency in prefill and decoding.

TABLE I: Performance comparison of different architectures.

| Architecture | | Conv. SA[a] | Vector Unit[b] | HSA[c] |
|---|---|---|---|---|
| Tokens/s | LISO[d] | 90.2 | 138.3 | **138.3** |
| | SILO[e] | 11.8 | 37.6 | **37.6** |
| Tokens/J | LISO[d] | 1060.7 | 719.1 | **1060.7** |
| | SILO[e] | 21.83 | 21.6 | **21.83** |
| Low-bit quantization support | | None | VSQ W4A4 | **MXINT4** **W4A8** |

Hardware Settings: (1) 256 PEs (2) 500MHz (3) Same tiling strategy
LLM Settings: End-to-end RetNet 1.3B
[a]Estimated based on HSA by considering the dataflow differences
[b]The result of [9] is scaled to 28nm CMOS as vector unit baseline
[c]Post-layout simulation results
[d,e]Prompt/output tokens is 750/50 for LISO, 50/750 for SILO

reasons: 1) Decode energy remains consistent due to excessive DRAM access, while LISO's large prefill computation highlights architectural differences. 2) During prefill, accelerators are compute-bound, making latency differences negligible for designs with same PE counts. In contrast, The SILO scenario better captures variations in decode latency. Our evaluation shows that vector unit consume 36% more energy than the other two architectures due to excessive SRAM access, as there is no weight reuse (activations can be reused by register buffers). In terms of latency, conventional SA suffers from low utilization rate, leading to poor decoding latency.

Table I presents the end-to-end performance comparison of the three architectures. HSA reducing power consumption by reusing activations and weights, while maintaining high utilization. These results demonstrate its effectiveness of addressing efficiency balance challenge mentioned in Fig. 2.

### B. Comparison with Prior Works

Table II shows the comparison with other works, including Keller et al. [9], FACT [19], MECLA [20], and other state-of-the-art ASIC accelerators [11], [15]. While these SoTA works report high 'peak' performance (TOPS) and energy efficiency (TOPS/W), memory bandwidth limitation in edge scenarios often constrain the accelerator from achieving such peak during inference. To reflect real-world edge conditions, we evaluate our design using end-to-end RetNet 1.3B inference with DDR5 memory bandwidth limitations (51.2GB/s).

Our accelerator achieves 2.45-10.83$\times$ and 13.5-61.99$\times$ improvements in area efficiency for token generation

TABLE II: Comparison of proposed accelerator and SoTA works.

| | Keller et al.[a] JSSC'23 [9] | MECLA[b] ISCA'24 [20] | FACT[c] ISCA'23 [19] | Kim et al.[d] ISSCC'24 [11] | Moon et al.[d] ISSCC'23 [15] | This work |
|---|---|---|---|---|---|---|
| Process Tech (nm) | 5 (28) | 28 | 28 | 28 | 28 | 28 |
| Dataflow | MVM | MMM | MMM | MMM | MVM | **Hybrid** |
| Area (mm$^2$) | 0.153 (4.8) | 22.02 | 6.03 | 20.25 | 7.29 | **0.636** |
| Frequency (MHz) | 152-1760 | 1000 | 500 | 50-200 | 20-400 | 500 |
| Power (W) | Not Reported | 2.87 | 0.337 | Not Reported | 0.002-0.237 | **0.108** |
| Post Proccesing Unit | YES | YES | YES | NO | NO | **YES** |
| Data Format | INT4/INT8 | INT8 | INT8 | INT8 | AQ 1-8b | **MXINT4/INT8** |
| Peak Performance (TOPS) | 3.6/1.8 | 14 | 1.02 | 3.41 | 0.52 | 0.256 |
| Peak Energy Efficiency (TOPS/W) | 91.1(2.9)/39.1(1.25) @0.47V,152MHz | 7.08 | 4.388 | 22.9 @0.5V,50MHz | 8.94 @0.46V,20MHz | 2.37 @0.81V |
| LISO[e] Area Effi. (token/s/mm$^2$)[f] | 395.18 (82.36)[g] | 100.82 | 50.83 | 22.84 | 30.83 | **247.38** |
| SILO[e] Area Effi. (token/s/mm$^2$)[f] | 38.00(7.92)[g] | 8.63 | 6.29 | 1.88 | 5.19 | **116.55** |
| Prefill Energy w/ EMA (mJ/token)[f] | 0.246[g] | 0.449 | 0.685 | 0.187 | 0.148 | 0.773 |
| Decode Energy w/ EMA (mJ/token)[f] | 48.59[g] | 25.10 | 54.54 | 47.96 | 45.15 | **24.06** |

[a]Scaled results show in parentheses for fair comparison across different technology node.
[b]TOPS reported is optimized by fine-tuning weight compression algorithm with accuracy loss, using SSMP standard setting
[c]TOPS reported is optimized by exploiting sparsity with accuracy loss    [d]Results without applying sparsity
[e]Running on RetNet 1.3B with prompt/output tokens of 750/50 for LISO and 50/750 for SILO
[f]Using DDR5 51.2GB/s bandwidth to simulate accelerator's behavior during edge decoding. DRAM access energy consumption 32pJ/Byte
[g]INT4-VSQ datapath cannot maintain acceptable model accuracy based on our experiment, results shown are using INT8 datapath

TABLE III: WikiText2 [14] perplexity↓ with 2,048 text length.

| Method | Precision | Scaling | RetNet-3B | Llama 3.2-3B | Llama 2-7B |
|---|---|---|---|---|---|
| Baseline | FP16 | - | 16.58 | 11.66 | 5.47 |
| SmoothQ. [26] | W8A8 | FP16 | 17.97 | 11.79 | 5.54 |
| Atom [28] | W4A4 | FP16 | - | - | 5.91 |
| QServe [12] | W4A8 | FP16 | - | - | 5.70 |
| VSQ [4] | W8A8 | INT8 | - | 11.736 | 5.44 |
| | W4A8 | INT8 | - | 1e35 | 107329 |
| **This work** | W4A8 | **4b Shift** | 18.22 | 13.11 | 5.81 |

TABLE IV: Generative reasoning accuracy↑ of GSM8K.

| Method | Precision | Scaling | Llama 3.1-8B-Instruct | Llama 3.2-3B-Instruct |
|---|---|---|---|---|
| Baseline | FP16 | - | 85.45 | 75.22 |
| SmoothQ. [26] | W8A8 | FP16 | 83.11 | 71.42 |
| SmoothQ. [26] | W4A8 | FP16 | 2.54 | 10.56 |
| **This work** | W4A8 | **4b Shift** | 82.54 | 69.22 |

TABLE V: Dequantization hardware overhead.

| Dequant. Scaling | 4-bit Shift | INT8 | FP16 |
|---|---|---|---|
| Normalized Area | **1.0×** | 10.33× | 15.96× |
| Normalized Power | **1.0×** | 7.19× | 9.60× |

when deploying other commonly used models like Llama and in different token length scenarios, as modern LLMs share the same core computations—MMM and MVM.

### C. Accuracy with Proposed Quantization Scheme

We evaluate the proposed hardware-friendly quantization scheme (Section III) for both RetNet and transformer-based Llama LLMs. We choose the weight group size to 16 (along the output channel) to match the capacity of each PE of the accelerator. Table III and Table IV show the comparison of perplexity results on WikiText2 dataset [14] and GSM8K dataset with SoTA quantization schemes [12], [26]. With 4-bit weight and 8-bit activation (W4A8), the proposed quantization method achieves performance on par with recent SoTA quantization methods with 8-bit weight and 8-bit activation (W8A8). For the reasoning-based GSM8K dataset, we evaluate our proposed method with instruction-tuned Llama model only since the RetNet is not designed for reasoning.

### VI. CONCLUSION

This work presents an efficient LLM accelerator featuring the hybrid systolic array architecture, which addresses the efficiency balance challenge in edge devices. We further reduce EMA by adopting MXINT4 format on weights, which is mapped onto our HSA architecture with high hardware utilization. Evaluated on post-layout simulation results with 500 MHz operating frequency, our accelerator achieves >2.45×/13.5× higher area efficiency across LISO/SILO scenarios for end-to-end token generation, while maintaining high energy efficiency in token generation.

(token/s/mm$^2$) under LISO and SILO, respectively. This gain primarily stems from the high utilization of HSA and the integration of low-bit quantization to mitigate memory bottlenecks. Our accelerator maintains superior area efficiency across various token length settings, with LISO and SILO serving as representative evaluation scenarios for long context.

For energy efficiency, reducing EMA is crucial, as shown in Fig. 1. While [4], [9] supports 4-bit per-vector quantization, our results (Table III) show that it fails to sustain performance under W4A8. [20] proposes a weight partitioning technique that successfully reduces EMA by 80% but introduces significant area and operating power overhead–both critical concerns for edge devices–and also requires fine-tuning on models.

In contrast, our approach reduce EMA by adopting MXINT4 weight quantization, maintaining high accuracy while being plug-and-play solution for any LLM. Our accelerator consumes 24.06 mJ/token for RetNet 1.3B token generation. The energy consumption reported in [11] was evaluated on Wikitext2 [14], which only goes though prefill stage and generates one token, whereas real-world applications decode numerous tokens and decoding dominates overall energy cost.

Although our evaluation focuses on the RetNet model with certain token length use cases, the high utilization rate and improvements in area and energy efficiency can be retained

## REFERENCES

[1] M. Abdin, J. Aneja, H. Awadalla, A. Awadallah, A. A. Awan, N. Bach, A. Bahree, A. Bakhtiari, J. Bao, H. Behl *et al.*, "Phi-3 technical report: A highly capable language model locally on your phone," *arXiv preprint arXiv:2404.14219*, 2024.

[2] N. Chatterjee, M. O'Connor, D. Lee, D. R. Johnson, S. W. Keckler, M. Rhu, and W. J. Dally, "Architecting an energy-efficient dram system for gpus," in *2017 IEEE International Symposium on High Performance Computer Architecture (HPCA)*. IEEE, 2017, pp. 73–84.

[3] H. Cho, D. Kim, S.-E. Hwang, and J. Park, "Sparc: Token similarity-aware sparse attention transformer accelerator via row-wise clustering," in *Proceedings of the 61st ACM/IEEE Design Automation Conference*, ser. DAC '24. New York, NY, USA: Association for Computing Machinery, 2024. [Online]. Available: https://doi.org/10.1145/3649329.3655936

[4] S. Dai, R. Venkatesan, M. Ren, B. Zimmer, W. Dally, and B. Khailany, "Vs-quant: Per-vector scaled quantization for accurate low-precision neural network inference," *Proceedings of Machine Learning and Systems*, vol. 3, pp. 873–884, 2021.

[5] T. Dao and A. Gu, "Transformers are ssms: generalized models and efficient algorithms through structured state space duality," in *Proceedings of the 41st International Conference on Machine Learning*, ser. ICML'24. JMLR.org, 2024.

[6] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan *et al.*, "The llama 3 herd of models," *arXiv preprint arXiv:2407.21783*, 2024.

[7] E. Frantar, S. Ashkboos, T. Hoefler, and D. Alistarh, "Gptq: Accurate post-training quantization for generative pre-trained transformers," *arXiv preprint arXiv:2210.17323*, 2022.

[8] G. Islamoglu, M. Scherer, G. Paulin, T. Fischer, V. J. Jung, A. Garofalo, and L. Benini, "Ita: An energy-efficient attention and softmax accelerator for quantized transformers," in *2023 IEEE/ACM International Symposium on Low Power Electronics and Design (ISLPED)*. IEEE, 2023, pp. 1–6.

[9] B. Keller, R. Venkatesan, S. Dai, S. G. Tell, B. Zimmer, C. Sakr, W. J. Dally, C. T. Gray, and B. Khailany, "A 95.6-tops/w deep learning inference accelerator with per-vector scaled 4-bit quantization in 5 nm," *IEEE Journal of Solid-State Circuits*, vol. 58, no. 4, pp. 1129–1141, 2023.

[10] B. Khailany, E. Khmer, R. Venkatesan, J. Clemons, J. S. Emer, M. Fojtik, A. Klinefelter, M. Pellauer, N. Pinckney, Y. S. Shao *et al.*, "A modular digital vlsi flow for high-productivity soc design," in *Proceedings of the 55th Annual Design Automation Conference*, 2018, pp. 1–6.

[11] S. Kim, S. Kim, W. Jo, S. Kim, S. Hong, and H.-J. Yoo, "20.5 c-transformer: A 2.6-18.1 $\mu$j/token homogeneous dnn-transformer/spiking-transformer processor with big-little network and implicit weight generation for large language models," in *2024 IEEE International Solid-State Circuits Conference (ISSCC)*, vol. 67. IEEE, 2024, pp. 368–370.

[12] Y. Lin, H. Tang, S. Yang, Z. Zhang, G. Xiao, C. Gan, and S. Han, "Qserve: W4a8kv4 quantization and system co-design for efficient llm serving," *arXiv preprint arXiv:2405.04532*, 2024.

[13] Y.-C. Lo and R.-S. Liu, "Bucket getter: A bucket-based processing engine for low-bit block floating point (bfp) dnns," in *Proceedings of the 56th Annual IEEE/ACM International Symposium on Microarchitecture*, 2023, pp. 1002–1015.

[14] S. Merity, C. Xiong, J. Bradbury, and R. Socher, "Pointer sentinel mixture models," 2016.

[15] S. Moon, H.-G. Mun, H. Son, and J.-Y. Sim, "A 127.8 tops/w arbitrarily quantized 1-to-8b scalable-precision accelerator for general-purpose deep learning with reduction of storage, logic and latency waste," in *2023 IEEE International Solid-State Circuits Conference (ISSCC)*. IEEE, 2023, pp. 21–23.

[16] R. Nallapati, B. Zhou, C. Gulcehre, B. Xiang *et al.*, "Abstractive text summarization using sequence-to-sequence rnns and beyond," *arXiv preprint arXiv:1602.06023*, 2016.

[17] Nvidia, "Nvidia jetson orin nano," https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/jetson-orin/, accessed: 19 November 2024.

[18] M. O'Connor, N. Chatterjee, D. Lee, J. Wilson, A. Agrawal, S. W. Keckler, and W. J. Dally, "Fine-grained dram: Energy-efficient dram for extreme bandwidth systems," in *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*, 2017, pp. 41–54.

[19] Y. Qin, Y. Wang, D. Deng, Z. Zhao, X. Yang, L. Liu, S. Wei, Y. Hu, and S. Yin, "Fact: Ffn-attention co-optimized transformer architecture with eager correlation prediction," in *Proceedings of the 50th Annual International Symposium on Computer Architecture*, 2023, pp. 1–14.

[20] Y. Qin, Y. Wang, Z. Zhao, X. Yang, Y. Zhou, S. Wei, Y. Hu, and S. Yin, "Mecla: Memory-compute-efficient llm accelerator with scaling sub-matrix partition," in *2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 2024, pp. 1032–1047.

[21] B. D. Rouhani, R. Zhao, A. More, M. Hall, A. Khodamoradi, S. Deng, D. Choudhary, M. Cornea, E. Dellinger, K. Denolf *et al.*, "Microscaling data formats for deep learning," *arXiv preprint arXiv:2310.10537*, 2023.

[22] J. Su, M. Ahmed, Y. Lu, S. Pan, W. Bo, and Y. Liu, "Roformer: Enhanced transformer with rotary position embedding," *Neurocomputing*, vol. 568, p. 127063, 2024.

[23] Y. Sun, L. Dong, S. Huang, S. Ma, Y. Xia, J. Xue, J. Wang, and F. Wei, "Retentive network: A successor to transformer for large language models," *arXiv preprint arXiv:2307.08621*, 2023.

[24] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale *et al.*, "Llama 2: Open foundation and fine-tuned chat models," *arXiv preprint arXiv:2307.09288*, 2023.

[25] W. Wang, S. Zhou, W. Sun, P. Sun, and Y. Liu, "Sole: Hardware-software co-design of softmax and layernorm for efficient transformer inference," in *2023 IEEE/ACM International Conference on Computer Aided Design (ICCAD)*. IEEE, 2023, pp. 1–9.

[26] G. Xiao, J. Lin, M. Seznec, H. Wu, J. Demouth, and S. Han, "Smoothquant: Accurate and efficient post-training quantization for large language models," in *International Conference on Machine Learning*. PMLR, 2023, pp. 38 087–38 099.

[27] H. Zhang, A. Ning, R. B. Prabhakar, and D. Wentzlaff, "Llmcompass: Enabling efficient hardware design for large language model inference," in *2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 2024, pp. 1080–1096.

[28] Y. Zhao, C.-Y. Lin, K. Zhu, Z. Ye, L. Chen, S. Zheng, L. Ceze, A. Krishnamurthy, T. Chen, and B. Kasikci, "Atom: Low-bit quantization for efficient and accurate llm serving," *Proceedings of Machine Learning and Systems (MLSys)*, 2024.