# Scalable MatMul-free Language Modeling

**Rui-Jie Zhu**[1], **Yu Zhang**[2], **Steven Abreu**[3*], **Ethan Sifferman**[1], **Tyler Sheaves**[1],
**Yiqiao Wang**[4], **Dustin Richmond**[1], **Sumit Bam Shrestha**[5], **Peng Zhou**[14], **Jason K. Eshraghian**[1†]

[1] University of California, Santa Cruz, [2] Soochow University, [3] University of Groningen,
[4] LuxiTech, [5] Intel Labs

## Abstract

Large Language Models (LLMs) have fundamentally altered how we approach scaling in machine learning. However, these models pose substantial computational and memory challenges, primarily due to the reliance on matrix multiplication (MatMul) within their attention and feed-forward (FFN) layers. We demonstrate that MatMul operations can be eliminated from LLMs while maintaining strong performance, even at billion-parameter scales. Our MatMul-free models, tested on models up to 2.7B parameters, are comparable to state-of-the-art pre-trained Transformers, and the performance gap narrows as model size increases.

Our approach yields significant memory savings: a GPU-efficient implementation reduces memory consumption by up to 61% during training and over 10× during inference. When adapted for a multi-chip neuromorphic system, the model leverages asynchronous processing to achieve 4× higher throughput with 10× less energy than edge GPUs. These findings demonstrate a path toward dramatically simplified yet effective LLMs, advancing them toward brain-like efficiency and heralding a new generation of lightweight, high-performance language models. Our code implementation is available at https://github. com/ridgerchu/matmulfreellm.

## 1 Main

Matrix Multiplication (MatMul) is the dominant operation in most neural networks, where dense layers involve vector-matrix multiplication (VMM), convolutions can be implemented as block-sparse VMMs with shared weights, and self-attention relies on matrix-matrix multiplication (MMM). The prevalence of MatMul is primarily due to Graphics Processing Units (GPUs) being optimized for MatMul operations. By leveraging Compute Unified Device Architecture (CUDA) and highly optimized linear algebra libraries such as cuBLAS, the MatMul operation can be efficiently parallelized and accelerated. This optimization was a key factor in the victory of AlexNet in the ILSVRC2012 competition and a historic marker for the rise of deep learning [1]. AlexNet notably utilized GPUs to boost training speed beyond CPU capabilities, and as such, deep learning won the 'hardware lottery' [2]. It also helped that both training and inference rely on MatMul.

Despite its prevalence in deep learning, MatMul operations account for the dominant portion of computational expense, often consuming the majority of the execution time and memory access during both training and inference phases. Several works have replaced MatMul with simpler operations through two main strategies. The first strategy involves substituting MatMul with elementary operations, e.g., AdderNet replaces multiplication with signed addition in convolutional neural networks (CNNs) [3]. Given the focus on convolutions, AdderNet is intended for use in computer vision over language modeling. ShiftAddLLM enables re-parameterization of dense layers, though

---

attention layers still rely on dynamic MatMul operations [4], where I/O limits will quickly overtake arithmetic savings as the sequence length grows.

The second approach employs binary or ternary quantization, simplifying MatMul to operations where values are either flipped or zeroed out before accumulation. Quantization can be applied to either activations or weights: spiking neural networks (SNNs) use binarized activations [5, 6, 7], while binary and ternary neural networks (BNNs and TNNs) use quantized weights [8]. Both methods can also be combined [9, 10].

Recent advances in language modeling, like BitNet and Falcon-Edge [11, 12, 13], demonstrate quantization's scalability, replacing all dense layer weights with binary/ternary values to support up to 3B parameters. Despite replacing VMMs with accumulations in all dense layers, they retain the self-attention mechanism which relies on an expensive MMM. Dynamically computed matrices $Q$ (query) and $K$ (key) are multiplied to form the attention map. Since both $Q$ and $K$ matrices are dynamically computed from pre-activation values, achieving optimal hardware efficiency on GPUs requires custom optimizations, such as specialized kernels and advanced memory access patterns. Despite these efforts, such MatMul operations remain resource-intensive on GPUs, as they involve extensive data movement and synchronization which can significantly hinder computational throughput and efficiency [14]. In our experiments, ternary quantization of the attention matrices in BitNet causes a significant drop in performance and failure to reach model convergence (see Fig. 4). This raises the question: is it possible to completely eliminate MatMul from LLMs?

In this work, we develop the first scalable MatMul-free language model (Matmul-free LM) by using additive operations in dense layers and element-wise Hadamard products for self-attention-like functions. Specifically, ternary weights eliminate MatMul in dense layers, similar to BNNs. To remove MatMul from self-attention, we optimize the Gated Recurrent Unit (GRU) [15] to rely solely on element-wise products and show that this model competes with state-of-the-art Transformers while eliminating all MatMul operations.

To fully exploit the efficiency potential of our architecture, we developed microcode-level optimization to process this model on a neuromorphic cluster using Intel's Loihi 2 platform. The architecture of our MatMul-free LM naturally aligns with neuromorphic computing paradigms, allowing us to achieve remarkable efficiency gains. By mapping our model to Loihi 2's mesh of asynchronous neurocores, we surpass human-readable throughput by approximately $8\times$ at 4.2 W power consumption, representing a significant advancement over conventional hardware. For autoregressive generation, our model maintains consistent 59.4 tokens/second throughput at a highly efficient 70.8 mJ/token, far outperforming comparable Transformer models on embedded GPUs by at least $4\times$ in throughput and $10\times$ in energy efficiency, demonstrating how neuromorphic hardware can effectively harness the MatMul-free properties of our architecture. This implementation moves LLMs closer to brain-like efficiency and points towards a new generation of lightweight, high-performance language models.

## 2 Building the MatMul-free Language Model

A closer examination of the Transformer architecture reveals two fundamental components that heavily rely on MatMul operations. The first is the dense layer, which not only transforms input hidden states to create Query, Key, and Value matrices for attention computation but also serves as the core structure in Feed-Forward Network (FFN) layers. Indeed, dense layers account for nearly all parameters in Transformer models. The second MatMul-dependent component is the attention mechanism itself, which performs MatMul operations directly on activations to compute attention scores.

The challenge, therefore, is to eliminate MatMul operations from both these structures while maintaining comparable performance. Following the Metaformer [16] framework, we can decompose the Transformer architecture into two essential functions: token mixing and channel mixing. The token mixer operates on input sequences (implemented as attention in traditional Transformers), while the channel mixer processes information across embedding dimensions (implemented as FFN layers in Transformers).

Our solution addresses both components through a comprehensive architectural transformation. We replace the attention mechanism with an element-wise recurrent neural network (RNN) that provides similar token mixing capabilities but relies solely on element-wise operations. Simultaneously, we
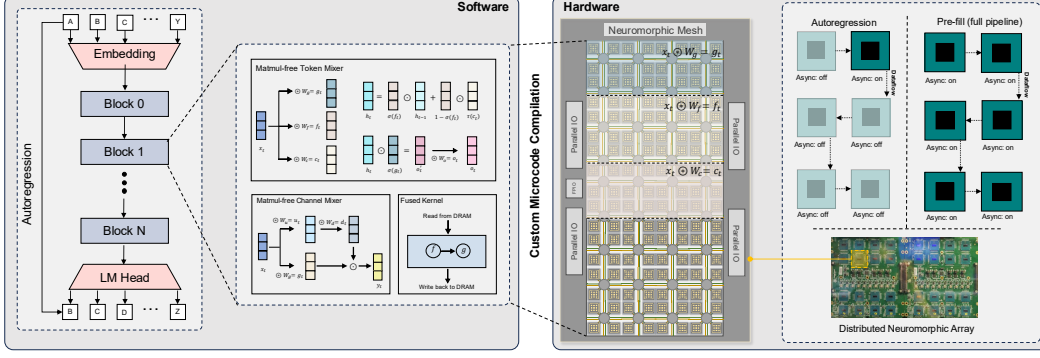
Figure 1: Overview of the MatMul-free LM. Left: general architecture of proposed model. Middle-right: Algorithm mapping of a single block across neurocores on a single Loihi2 chip. Top-Right: Multi-chip Dataflow for autoregression and pre-fill. During autoregression, only one chip consumes non-negligible dynamic power dissipation due to the clock-free system. Bottom-right: The MatMul-free LM is deployed on the Hala Point system which consists of 1,152 Loihi 2 chips.
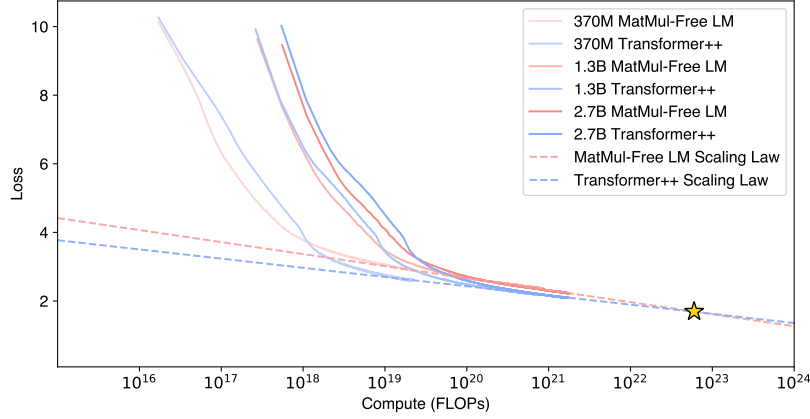


Figure 2: Scaling law comparison between MatMul-free LM and Transformer++ models, depicted through their loss curves. The red lines represent the loss trajectories of the MatMul-free LM, while the blue lines indicate the losses of the Transformer++ models. The star marks the intersection point of the scaling law projection for both model types. MatMul-free LM uses ternary parameters and BF16 activations, whereas Transformer++ uses BF16 parameters and activations.

transform all dense layers to use ternary weights, effectively eliminating MatMul operations throughout the entire architecture. This approach maintains the core functional properties of Transformers while dramatically reducing computational complexity. This architectural transformation achieves a complete elimination of MatMul operations through carefully designed structural modifications, while preserving the essential mixing capabilities that make Transformers effective. Specific implementation details are provided in 6. The following sections detail how each component is optimized and demonstrate the effectiveness of this MatMul-free approach.

## 3 Scaling Analysis

Neural scaling laws posit that model error decreases as a power function of training set size and model size. Such projections become important as training becomes increasingly expensive with larger models. A widely adopted best practice in LLM training is to first test scalability with smaller models, where scaling laws begin to take effect [17, 18, 19]. The GPT-4 technical report revealed that

Table 1: Zero-shot accuracy of MatMul-free LM and Transformer++ on benchmark datasets.

| Models | Size | ARCe | ARCc | HS | OQ | PQ | WGe | Avg. |
|---|---|---|---|---|---|---|---|---|
| *370M parameters with 15B training tokens, Layer=24, d=1024* | | | | | | | | |
| Transformer++ | 370M | 45.0 | 24.0 | 34.3 | 29.2 | 64.0 | 49.9 | 41.1 |
| MatMul-free RWKV-4 | 370M | 44.7 | 22.8 | 31.6 | 27.8 | 63.0 | 50.3 | 40.0 |
| **Ours** | 370M | 42.6 | 23.8 | 32.8 | 28.4 | 63.0 | 49.2 | 40.3 |
| *1.3B parameters with 100B training tokens, Layer=24, d=2048* | | | | | | | | |
| Transformer++ | 1.3B | 54.1 | 27.1 | 49.3 | 32.4 | 70.3 | 54.9 | 48.0 |
| MatMul-free RWKV-4 | 1.3B | 52.4 | 25.6 | 45.1 | 31.0 | 68.2 | 50.5 | 45.5 |
| **Ours** | 1.3B | 54.0 | 25.9 | 44.9 | 31.4 | 68.4 | 52.4 | 46.2 |
| *2.7B parameters with 100B training tokens, Layer=32, d=2560* | | | | | | | | |
| Transformer++ | 2.7B | 59.7 | 27.4 | 54.2 | 34.4 | 72.5 | 56.2 | 50.7 |
| **Ours** | 2.7B | 58.5 | 29.7 | 52.3 | 35.4 | 71.1 | 52.1 | 49.9 |

a prediction model just $1/10,000$ the size of the final model can still accurately forecast the full-sized model performance [20].

We evaluate how the scaling law fits to the 370M, 1.3B and 2.7B parameter models in both Transformer++ and MatMul-free LM, shown in Fig. 2. For a conservative comparison, each operation is treated identically between MatMul-free LM and Transformer++. But note that all weights and activations in Transformer++ are in BF16, while BitLinear layers in MatMul-free LM use ternary parameters, with BF16 activations. As such, an average operation in MatMul-free LM will be computationally cheaper than that of Transformer++.

Interestingly, the scaling projection for the MatMul-free LM exhibits a steeper descent compared to that of Transformer++. This suggests that the MatMul-free LM is more efficient in leveraging additional compute resources to improve performance. As a result, the scaling curve of the MatMul-free LM is projected to intersect with the scaling curve of Transformer++ at approximately $10^{23}$ FLOPs. This compute scale is roughly equivalent to the training FLOPs required for Llama-3 8B (trained with 15 trillion tokens) and Llama-2 70B (trained with 2 trillion tokens), suggesting that MatMul-free LM not only outperforms in efficiency, but can also outperform in terms of loss when scaled up.

## 4    Performance on Downstream Tasks

In line with benchmarking in BitNet, we evaluated the zero-shot performance of these models on a range of language tasks, including ARC-Easy [21], ARC-Challenge [21], Hellaswag [22], Winogrande [23], PIQA [24], and OpenbookQA [25]. The results are shown in Tab. 1. All evaluations are performed using the LM evaluation harness [26]. The MatMul-free LM models achieve competitive performance compared to the Transformer++ baselines across all tasks, demonstrating its effectiveness in zero-shot learning despite the absence of MatMul operations, and the lower memory required from ternary weights. Notably, the 2.7B MatMul-free LM model outperforms its Transformer++ counterpart on ARC-Challenge and OpenbookQA, while maintaining comparable performance on the other tasks. As the model size increases, the performance gap between MatMul-free LM and Transformer++ narrows, which is consistent with the scaling law. These results highlight that MatMul-free architectures are capable of achieving strong zero-shot performance on a diverse set of language tasks, ranging from question answering and commonsense reasoning to physical understanding.

## 5    Deployment on Neuromorphic Hardware

### 5.1    Training Efficiency Optimization on GPU

We evaluate our proposed Fused BitLinear and Vanilla BitLinear implementations in terms of training time and memory usage, shown in Fig. 3(a-b). For each experiment, we set the input size and sequence length to 1024. All experiments are conducted using an NVIDIA A100 80GB GPU. Note

**(a) Computational Latency vs. Batch Size**: Comparative performance of Fused and Vanilla BitLinear implementations

**(b) Memory Utilization vs. Batch Size**: Memory efficiency comparison between Fused and Vanilla BitLinear architectures

**(c) Resource Efficiency Analysis**: GPU memory consumption and inference latency comparison between MatMul-free LM and Transformer++ across model size
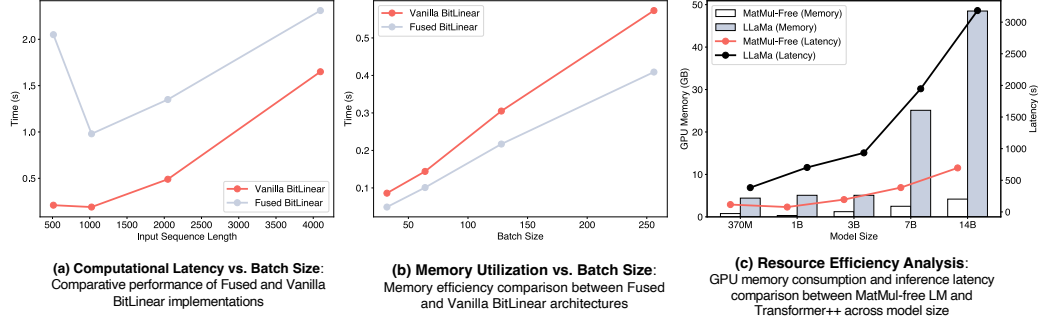
Figure 3: Performance comparison and analysis of different models and configurations. (a) and (b) show the training performance comparison between Vanilla BitLinear and Fused BitLinear in terms of time and memory consumption as a function of batch size. (c) compares the inference memory consumption and latency between MatMul-free LM and Transformer++ across various model sizes.

that during training, the sequence length and batch dimensions are flattened, making the effective batch size the product of these dimensions.

Our experiments show that our fused operator benefits from larger batch sizes in terms of faster training speeds and reduced memory consumption. When the batch size is $2^8$, the training speed of the 1.3B parameter model improves from 1.52s to 1.21s per iteration, a 25.6% speedup over the vanilla implementation. Additionally, memory usage decreases from 82GB to 32GB, a 61.0% reduction. The performance of the fused implementation improves significantly with larger batch sizes, allowing more samples to be processed simultaneously and reducing the total number of iterations.

## 5.2 Inference Efficiency Optimization on GPU

Fig. 3(c) presents a comparison of GPU inference memory consumption and latency between the proposed MatMul-free LM and Transformer++ for various model sizes. In the MatMul-free LM, we employ BitBLAS [27] for acceleration to further improve efficiency. The evaluation is conducted with a batch size of 1 and a sequence length of 2048. The MatMul-free LM consistently demonstrates lower memory usage and latency compared to Transformer++ across all model sizes. For a single layer, the MatMul-free LM requires only 0.12 GB of GPU memory and achieves a latency of 3.79 ms, while Transformer++ consumes 0.21 GB of memory and has a latency of 13.87 ms. As the model size increases, the memory and latency advantages of the MatMul-free LM become more pronounced. It is worth noting that for model sizes larger than 2.7B, the results are simulated using randomly initialized weights. For the largest model size of 13B parameters, the MatMul-free LM uses only 4.19 GB of GPU memory and has a latency of 695.48 ms, whereas Transformer++ requires 48.50 GB of memory and exhibits a latency of 3183.10 ms. These results highlight the efficiency gains achieved by the MatMul-free LM, making it a promising approach for large-scale language modeling tasks, particularly during inference.

## 5.3 Neuromorphic Computing with Intel Loihi 2

While GPUs excel at conventional deep learning workloads grounded in dense floating-point matrix multiplications, the MatMul-free LM we present is dominated by low-precision element-wise operations. The arithmetic intensity is so low that many CUDA cores remain idle during inference. Additionally, using ternary weights in matrix multiplications naturally induces unstructured sparsity of $\approx 35\%$ (for the 370M model) which cannot easily be accelerated on GPUs, resulting in redundant calculations and memory movement. Our MatMul-free LM maps naturally onto large-scale neuromorphic hardware that executes element-wise, recurrent, low-precision operations close to memory. Intel's Loihi 2, for instance, offers microcode-programmable neuron dynamics and on-chip support for sparse low-precision arithmetic, making it an ideal substrate for the MatMul-free model [28]. Processing is done locally within 120 fully asynchronous "neurocores" on each chip, with the option of connecting up to 1,152 chips together into a larger processing system, see Fig. 6 in Appendix C.

Each neurocore stores weights and recurrent states in its local SRAM, which minimizes memory movement and consequently improves both energy efficiency and latency.

Table 2: Quantization results for the 370M model, using PyTorch on GPU. All weight matrices are ternary. W8: using 8-bit element-wise weights. A8/A16: using 8-bit or 16-bit activations.

|  | ARCc | ARCe | HS | OQ | PQ | WGe | Avg. |
|---|---|---|---|---|---|---|---|
| Transformer++ | 24.0% | 45.0% | 34.3% | 29.2% | 64.0% | 49.9% | 41.1% |
| **Ours** | 23.8% | 42.6% | 32.8% | 28.4% | 63.0% | 49.2% | 40.3% |
| W8A8 | 28.3% | 26.8% | 26.1% | 27.0% | 52.7% | 51.5% | 35.4% |
| W8A16 | 23.0% | 42.4% | 32.4% | 27.8% | 63.0% | 50.1% | 39.8% |

To accommodate Loihi 2's low-precision fixed-point arithmetic, we studied the zero-shot accuracy of our 370M-parameter model under various quantization schemes (Table 2). Quantizing normalization parameters and activations to 8-bit weights and 16-bit activations (W8A16) preserves nearly the same performance as the original MatMul-free LM, whereas 8-bit activations (W8A8) reduce accuracy more noticeably, thus activations are kept in 16-bit precision. All quantized operators were implemented end-to-end in fixed-point arithmetic on Loihi 2, including a custom fixed-point RMSNorm layer (Appendix C.2).

Table 3: Throughput and energy efficiency for various transformer-based language models with at most 500M parameters running on the NVIDIA Jetson Orin Nano compared to our MatMul-free LM running on Intel's Loihi 2. Bolded metrics are based on single-chip metrics and inter-chip communication results[†].

|  | Throughput (tokens/sec) | | | | Efficiency (mJ/token) | | | |
|---|---|---|---|---|---|---|---|---|
| Generate | 500 | 1000 | 2000 | 4000 | 500 | 1000 | 2000 | 4000 |
| Alireo-400M | 14.3 | 14.9 | 15.0 | 14.7 | 723 | 719 | 751 | 853 |
| Qwen2-500M | 13.4 | 14.0 | 14.1 | 14.1 | 791 | 785 | 816 | 912 |
| Ours (370M, 1-chip*) | 71.3 | 71.3 | 71.3 | 71.3 | 59 | 59 | 59 | 59 |
| **Ours (370M, system[†])** | **59.4** | **59.4** | **59.4** | **59.4** | **70.8** | **70.8** | **70.8** | **70.8** |
| Prefill | 500 | 1000 | 2000 | 4000 | 500 | 1000 | 2000 | 4000 |
| Alireo-400M | 849.4 | 1620 | 2858 | 3153 | 11.7 | 7.8 | 5.8 | 6.8 |
| Qwen2-500M | 627 | 909 | 1514 | 2639 | 17.9 | 13.9 | 9.5 | 6.7 |
| Ours (370M, 1-chip*) | 13965 | 13965 | 13965 | 13965 | 2.8 | 2.8 | 2.8 | 2.8 |
| **Ours (370M, system[†])** | **11637** | **11637** | **11637** | **11637** | **3.4** | **3.4** | **3.4** | **3.4** |

* The MatMul-free LM on Loihi 2 was characterized on an Oheo Gulch single-chip Loihi 2 system (N3C1 silicon) running NxKernel v0.2.0 and NxCore v2.5.8 (only accessible to Intel Neuromorphic Research Community members). The 1-chip case neglects inter-chip communication.
[†] Inter-chip communication causes a derived $\approx$20% slowdown over the single chip case (Appendix C.4.1).
[‡] Transformer LMs characterized on NVIDIA Jetson Orin Nano 8GB using MAXN power mode running Jetpack 6.2, TensorRT 10.3.0, CUDA 12.4. Energy values include CPU_GPU_CV, SOC, and IO components as reported by jtop 4.3.0.
Performance results as of Jan 2025 and may not reflect all publicly available security updates. Results may vary.

Throughput is at least 4× higher for autoregressive generation and at least 3.6× higher for prefill. More details can be found in Appendix C.4.1. Table 3 contrasts the performance of the 370M MatMul-free model on Loihi 2 against Transformer baselines running on an NVIDIA Jetson Orin Nano, showing at least 1.7× less energy per token for prefill, and at least 10× less energy per token for auto-regressive generation. We compare against the 500M parameter Qwen2 model [29], and also against the 400M parameter Alireo model [30] which is representative of a small-scale Llama model [31]. We have included extended experimental results that compare against server-class GPUs (H100) in the Appendix to provide additional context (C.4), though omit them here due to the significant differences between memory and power draw.

# 6 Methods

This section first introduces MatMul-free BitLinear layers with ternary weights $(-1, 0, +1)$, replacing multiplications with additions and negations to cut compute and memory while preserving expressiveness (Sec. 6.1). It then presents a hardware-efficient fused BitLinear implementation (Sec. 6.2) and assembles the full MatMul-free LM architecture (Sec. 6.3): a MatMul-free token mixer built on the Linear Gated Recurrent Unit (MLGRU) to capture sequence dependencies and a channel mixer that employs a BitLinear-based Gated Linear Unit, so the entire model relies solely on additions and element-wise products. Training details are summarized in Sec. 6.4.

## 6.1 MatMul-free Dense Layers with Ternary Weights

In a standard dense layer, the MatMul between the input $x \in \mathbb{R}^d$ and the weight matrix $W \in \mathbb{R}^{d \times m}$ can be expressed as:

$$y = xW = \sum_{j=1}^{d} x_j W_{ij} \quad \text{for } i = 1, 2, \ldots, m$$

where $y \in \mathbb{R}^m$ is the output. To avoid using standard MatMul-based dense layers, we adopt BitNet to replace dense layers containing MatMuls with BitLinear modules, which use ternary weights to transform MatMul operations into pure addition operation with accumulation, i.e., ternary accumulation. When using ternary weights, the elements from the weight matrix $W$ are constrained to values from the set $\{-1, 0, +1\}$. Let $\widetilde{\mathbf{W}}$ denote the ternary weight matrix. The MatMul with ternary weights can be expressed as:

$$\widetilde{\mathbf{Y}} = x \circledast \widetilde{\mathbf{W}} = \sum_{j=1}^{d} x_j \widetilde{\mathbf{W}}_{ij}, \quad \widetilde{\mathbf{W}}_{ij} \in \{-1, 0, +1\}, \quad \text{for } i = 1, 2, \ldots, m$$

where $\widetilde{\mathbf{Y}} \in \mathbb{R}^m$ is the output, and $\circledast$ represents a ternary MatMul, which can be simplified to accumulation. Since the ternary weights $\widetilde{\mathbf{W}}_{ij}$ can only take values from $\{-1, 0, +1\}$, the multiplication operation in the MatMul can be replaced by a simple addition or subtraction operation:

$$x_j \widetilde{\mathbf{W}}_{ij} = \begin{cases} x_j, & \text{if } \widetilde{\mathbf{W}}_{ij} = 1, \\ 0, & \text{if } \widetilde{\mathbf{W}}_{ij} = 0, \\ -x_j, & \text{if } \widetilde{\mathbf{W}}_{ij} = -1. \end{cases}$$

Therefore, ternary MatMul can be written as follows:

$$\widetilde{\mathbf{Y}}_i = \sum_{j=1}^{d} x_j \widetilde{\mathbf{W}}_{ij} = \sum_{j:\widetilde{\mathbf{W}}_{ij}=1} x_j - \sum_{j:\widetilde{\mathbf{W}}_{ij}=-1} x_j, \quad \text{for } i = 1, 2, \ldots, m$$

## 6.2 Hardware-efficient Fused BitLinear Layer

BitNet showed that stabilizing ternary layers requires an additional RMSNorm before the BitLinear input. However, the vanilla implementation of BitNet is not efficient. Modern GPUs feature a memory hierarchy with a large, global high-bandwidth memory (HBM) and smaller, faster shared memory (SRAM), and the implementation of BitNet introduced many I/O operations: reading the previous activation into SRAM for RMSNorm, writing it back for quantization, reading it again for quantization, storing it, and reading it once more for the Linear operation. To address this inefficiency, we read the activation only once and apply RMSNorm and quantization as fused operations in SRAM, which we present in Algorithm 1. Optimal utilization of SRAM to reduce HBM I/O costs can significantly speed up computations. Since the activations in this model have a larger memory footprint than ternary weights and the large number of element-wise operations, our optimization efforts focus on activations.

**Algorithm 1** Fused RMSNorm and BitLinear Algorithm with Quantization

**Define** FORWARDPASS($\mathbf{X}, \mathbf{W}, \boldsymbol{b}, \epsilon$)
$\quad\mathbf{X} \in \mathbb{R}^{M \times N}, \mathbf{W} \in \mathbb{R}^{N \times K}, \boldsymbol{b} \in \mathbb{R}^{K}$

$\quad$ **function** forward_pass($\mathbf{X}, \mathbf{W}, \boldsymbol{b}, \epsilon$)
$\quad\quad$ Load $\mathbf{X}, \mathbf{W}, \boldsymbol{b}, \epsilon$ from HBM
$\quad\quad$ On Chip: $\widetilde{\mathbf{Y}}, \mu, \sigma^2, r \leftarrow$ rms_norm_fwd($\mathbf{X}$)
$\quad\quad$ On Chip: $\widetilde{\mathbf{W}} \leftarrow$ weight_quant($\mathbf{W}$)
$\quad\quad$ On Chip: $\mathbf{O} \leftarrow \widetilde{\mathbf{Y}} \circledast \widetilde{\mathbf{W}} + \boldsymbol{b}$
$\quad\quad$ Store $\mathbf{O}, \mu, \sigma^2, r$ to HBM
$\quad\quad$ **return** $\mathbf{O}, \mu, \sigma^2, r$

$\quad$ **function** rms_norm_fwd($\mathbf{X}$)
$\quad\quad$ $\mu, \sigma^2 \leftarrow$ mean($\mathbf{X}$), variance($\mathbf{X}$)
$\quad\quad$ $r \leftarrow \frac{1}{\sqrt{\sigma^2 + \epsilon}}$
$\quad\quad$ $\widetilde{\mathbf{Y}} \leftarrow$ activation_quant($r(\mathbf{X} - \mu)$)
$\quad\quad$ **return** $\widetilde{\mathbf{Y}}, \mu, \sigma^2, r$

$\quad$ **function** activation_quant($\mathbf{X}$)
$\quad\quad$ $s \leftarrow \frac{127}{\max(|\mathbf{X}|)}$ $\quad\quad\triangleright \lfloor \cdot \rceil$ |. means round then clamp
$\quad\quad$ $\widetilde{X} \leftarrow \lfloor s\mathbf{X} \rceil |_{[-128,127]} \cdot \frac{1}{s}$
$\quad\quad$ **return** $\widetilde{X}$

$\quad$ **function** weight_quant($\mathbf{W}$)
$\quad\quad$ $s \leftarrow \frac{1}{\text{mean}(|\mathbf{W}|)}$
$\quad\quad$ $\widetilde{\mathbf{W}} \leftarrow \lfloor s\mathbf{X} \rceil |_{[-1,1]} \cdot \frac{1}{s}$
$\quad\quad$ **return** $\widetilde{\mathbf{W}}$
$\quad$ **return** $\mathbf{O}$

**Define** BACKWARDPASS($\mathbf{X}, \mathbf{W}, \boldsymbol{b}, \mathbf{O}, \mathrm{d}\mathbf{O}, \mu, \sigma^2, r$)
$\quad\mathbf{X} \in \mathbb{R}^{M \times N}, \mathbf{W} \in \mathbb{R}^{N \times K}, \boldsymbol{b} \in \mathbb{R}^{K}$
$\quad\mathbf{O} \in \mathbb{R}^{M \times K}, \mathrm{d}\mathbf{O} \in \mathbb{R}^{M \times K}$

$\quad$ **function** backward_pass($\mathbf{X}, \mathbf{W}, \boldsymbol{b}, \mathbf{O}, \mu, \sigma^2, r, \mathrm{d}\mathbf{O}$)
$\quad\quad$ Load $\mathbf{X}, \mathbf{W}, \boldsymbol{b}, \mathbf{O}, \mu, \sigma^2, r, \mathrm{d}\mathbf{O}$ from HBM
$\quad\quad$ On Chip: $\mathrm{d}\mathbf{Y} \leftarrow \mathrm{d}\mathbf{O} \times \mathbf{W}^\top$
$\quad\quad$ On Chip: $\mathrm{d}\mathbf{X}, \widetilde{\mathbf{Y}} \leftarrow$ rms_norm_bwd($\mathrm{d}\mathbf{Y}, \mathbf{X}, \mu, \sigma^2, r$)
$\quad\quad$ On Chip: $\mathrm{d}\mathbf{W} \leftarrow \mathrm{d}\mathbf{O}^\top \times \widehat{\mathbf{Y}}$
$\quad\quad$ On Chip: $\mathrm{d}\boldsymbol{b} \leftarrow$ sum($\mathrm{d}\mathbf{O}$)
$\quad\quad$ Store $\mathrm{d}\mathbf{X}, \mathrm{d}\mathbf{W}, \mathrm{d}\boldsymbol{b}$ to HBM
$\quad\quad$ **return** $\mathrm{d}\mathbf{X}, \mathrm{d}\mathbf{W}, \mathrm{d}\boldsymbol{b}$

$\quad$ **function** rms_norm_bwd($\mathrm{d}\mathbf{Y}, \mathbf{X}, \mu, \sigma^2, r$)
$\quad\quad$ $\widetilde{\mathbf{Y}} \leftarrow$ activation_quant($r(\mathbf{X} - \mu)$)
$\quad\quad$ $\mathrm{d}\sigma^2 \leftarrow$ sum($\mathrm{d}\mathbf{Y} \times (\mathbf{X} - \mu) \times -0.5 \times r^3$)
$\quad\quad$ $\mathrm{d}\mu \leftarrow$ sum($-r\mathrm{d}\mathbf{Y}$) + $\mathrm{d}\sigma^2 \times$ mean($\mathbf{X} - \mu$)
$\quad\quad$ $\mathrm{d}\mathbf{X} \leftarrow r\mathrm{d}\mathbf{Y} + 2\mathrm{d}\sigma^2(\mathbf{X} - \mu)/N + \mathrm{d}\mu/N$
$\quad\quad$ **return** $\mathrm{d}\mathbf{X}, \widetilde{\mathbf{Y}}$

The `forward_pass` function in Algorithm 1 first calls `rms_norm_fwd` to perform RMSNorm on input activations $\mathbf{X}$, loading normalized activations $\mathbf{Y}$, mean $\mu$, variance $\sigma^2$, and scaling factor $r$ from HBM. The normalized activations $\mathbf{Y}$ are then quantized to obtain $\widetilde{\mathbf{Y}}$ and the weights $\mathbf{W}$ are quantized using `weight_quant`, with both performed without off-chip data movement. Finally, the output $\mathbf{O}$ is computed on-chip by multiplying quantized activations $\widetilde{\mathbf{Y}}$ with the ternary weights $\widetilde{\mathbf{W}}$, adding the bias $\boldsymbol{b}$, and then storing the result back to HBM.

The `backward_pass` function first loads $\mathbf{X}, \mathbf{W}, \boldsymbol{b}, \mathbf{O}, \mu, \sigma^2, r$, and the output gradient $\mathrm{d}\mathbf{O}$ from HBM. The gradient $\mathrm{d}\mathbf{Y}$ is then computed on-chip by multiplying the output gradient $\mathrm{d}\mathbf{O}$ with the transposed weight matrix $\mathbf{W}^\top$. Next, it calls `rms_norm_bwd` on-chip to backpropagate through RMSNorm, computing the input gradient $\mathrm{d}\mathbf{X}$. The weight gradient $\mathrm{d}\mathbf{W}$ is calculated on-chip by multiplying the transposed output gradient $\mathrm{d}\mathbf{O}^\top$ with the quantized activations $\widehat{\mathbf{Y}}$, and the bias gradient $\mathrm{d}\boldsymbol{b}$ is obtained by summing
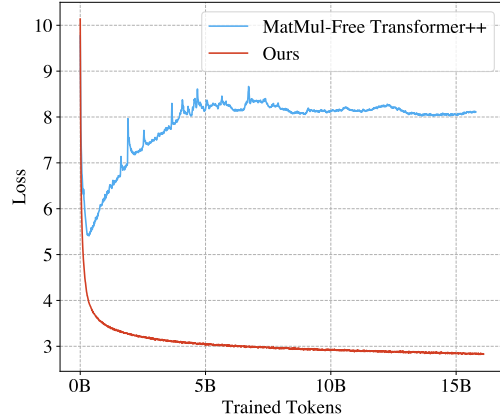


Figure 4: Training loss over steps for the MatMul-free Transformer++ and our proposed method in 370M. The MatMul-free Transformer++ fails to converge, while our method successfully converges under the MatMul-free setting.

$\mathrm{d}\mathbf{O}$. The computed gradients $\mathrm{d}\mathbf{X}, \mathrm{d}\mathbf{W}$, and $\mathrm{d}\boldsymbol{b}$ are then transferred back to HBM. Sec. 5.1 presents an experimental comparison between a vanilla BitLinear implementation and Fused BitLinear.

## 6.3 MatMul-free Language Model Architecture

We adopt the perspective from Metaformer [16], which suggests that Transformers consist of a token-mixer (for mixing temporal information, i.e., Self Attention [32], Mamba [33]) and a channel-mixer (for mixing embedding/spatial information, i.e., feed-forward network, Gated Linear Unit (GLU) [34, 35]). A high-level overview of the architecture is shown in Fig. 1.

### 6.3.1 MatMul-free Token Mixer

Self-attention is the most common token mixer in modern language models, relying on matrix multiplication between three matrices: $Q$, $K$, and $V$. To convert these operations into additions, we ternarize at least two of the matrices. Assuming all dense layer weights are ternary, we quantize $Q$ and $K$, resulting in a ternary attention map that eliminates multiplications in self-attention. However, as shown in Fig. 4, such a model fails to converge. One possible explanation is that activations contain outliers crucial for performance but difficult to quantize effectively [36, 37]. To address this challenge, we explore alternative methods for mixing tokens without relying on matrix multiplications.

By resorting to the use of ternary RNNs, which combine element-wise operations and accumulation, it becomes possible to construct a MatMul-free token mixer. Among various RNN architectures, the GRU is noted for its simplicity and efficiency, achieving similar performance to Long Short-Term Memory (LSTM) [38] cells while using fewer gates and having a simpler structure. Thus, we choose the GRU as the foundation for building a MatMul-free token mixer. We first revisit the standard GRU and then demonstrate, step by step, how we derive the MLGRU.

**Revisiting the Gated Recurrent Unit**  The GRU [15] can be formalized as follows:

$$\boldsymbol{r}_t = \sigma\left(\boldsymbol{x}_t\mathbf{W}_{xr} + \boldsymbol{h}_{t-1}\mathbf{W}_{hr} + \mathbf{b}_r\right) \in \mathbb{R}^d, \tag{1}$$

$$\boldsymbol{f}_t = \sigma\left(\boldsymbol{x}_t\mathbf{W}_{xf} + \boldsymbol{h}_{t-1}\mathbf{W}_{hf} + \mathbf{b}_f\right) \in \mathbb{R}^d, \tag{2}$$

$$\boldsymbol{c}_t = \tanh\left(\boldsymbol{x}_t\mathbf{W}_{xc} + (\boldsymbol{r}_t \odot \boldsymbol{h}_{t-1})\mathbf{W}_{cc} + \mathbf{b}_c\right) \in \mathbb{R}^d, \tag{3}$$

$$\boldsymbol{h}_t = \boldsymbol{f}_t \odot \boldsymbol{h}_{t-1} + (1 - \boldsymbol{f}_t) \odot \boldsymbol{c}_t \in \mathbb{R}^d, \tag{4}$$

$$\boldsymbol{o}_t = \boldsymbol{h}_t \tag{5}$$

where $\boldsymbol{x}_t \in \mathbb{R}^m$ is the input vector at time step $t$, $\boldsymbol{h}_{t-1} \in \mathbb{R}^d$ is the hidden state vector from the previous time step, $\boldsymbol{r}_t \in \mathbb{R}^d$ is the reset gate vector, $\boldsymbol{f}_t \in \mathbb{R}^d$ is the forget gate vector, $\boldsymbol{c}_t \in \mathbb{R}^d$ is the candidate hidden state, $\boldsymbol{h}_t \in \mathbb{R}^d$ is the final hidden state vector at time step $t$, $\boldsymbol{o}_t \in \mathbb{R}^d$ is the output vector at time step $t$, $\mathbf{W}(\cdot) \in \mathbb{R}^{m \times d}$ and $\mathbf{b}(\cdot) \in \mathbb{R}^d$ are learnable weight matrices and bias vectors, respectively, $\sigma(\cdot)$ is the sigmoid activation function, and $\odot$ denotes element-wise multiplication.

A key characteristic of the GRU is the coupling of the input gate vector $\boldsymbol{f}_t$ and the forget gate vector $(1 - \boldsymbol{f}_t)$, which together constitute the 'leakage' unit. This leakage unit decays the hidden state $\boldsymbol{h}_{t-1}$ and the candidate hidden state $\boldsymbol{c}_t$ through element-wise multiplication (Eq. 4). This operation allows the model to adaptively retain information from the previous hidden state $\boldsymbol{h}_{t-1}$ and incorporate new information from the candidate hidden state $\boldsymbol{c}_t$. Importantly, this operation relies solely on element-wise multiplication, avoiding the need for MatMul. We preserve this property of the GRU while introducing further modifications to create a MatMul-free variant of the model.

**MatMul-free Linear Gated Recurrent Unit**  We first remove hidden-state related weights $\mathbf{W}_{cc}$, $\mathbf{W}_{hr}$, $\mathbf{W}_{hf}$, and the activation between hidden states (tanh). This modification not only makes the model MatMul-free but also linearized the GRU through time, enabling training via a parallel scan [33, 39]. This approach is critical for improving computational efficiency, as transcendental functions are expensive to compute accurately, and non-diagonal transition matrices increase runtime complexity from $\mathcal{O}(\backslash)$ to $\mathcal{O}(\backslash^{\ni})$. This modification is a key feature of recent RNNs, such as the Linear Recurrent Unit [40], Hawk [41], and RWKV-4 [42]. We then add a data-dependent output gate between $\boldsymbol{h}_t$ and $\boldsymbol{o}_t$, inspired by the LSTM and widely adopted by recent RNN models:

$$\boldsymbol{g}_t = \boldsymbol{x}_t\mathbf{W}_g + \mathbf{b}_g \in \mathbb{R}^d,$$

$$\boldsymbol{o}_t' = \boldsymbol{g}_t \odot \sigma(\boldsymbol{h}_t) \in \mathbb{R}^d,$$

$$\boldsymbol{o}_t = \boldsymbol{o}_t'\mathbf{W}_o + \mathbf{b}_o \in \mathbb{R}^d.$$

Following the approach of HGRN [43], we further simplify the computation of the candidate hidden state by keeping it as a simple linear transform, rather than coupling it with the hidden state. This can be rewritten as a linear transformation of the input. Finally, we replace all remaining weight matrices with ternary weight matrices, completely removing the MatMul operations. The resulting MLGRU

9

architecture can be formalized as follows:

$$\boldsymbol{f}_t = \sigma\left(\boldsymbol{x}_t \circledast \mathbf{W}_f + \mathbf{b}_f\right) \in \mathbb{R}^d,$$

$$\boldsymbol{c}_t = \tau\left(\boldsymbol{x}_t \circledast \mathbf{W}_c + \mathbf{b}_c\right) \in \mathbb{R}^d,$$

$$\boldsymbol{h}_t = \boldsymbol{f}_t \odot \boldsymbol{h}_{t-1} + (1 - \boldsymbol{f}_t) \odot \boldsymbol{c}_t \in \mathbb{R}^d,$$

$$\boldsymbol{g}_t = \boldsymbol{x}_t \circledast \mathbf{W}_g + \mathbf{b}_g \in \mathbb{R}^d,$$

$$\boldsymbol{o}'_t = \boldsymbol{g}_t \odot \sigma(\boldsymbol{h}_t) \in \mathbb{R}^d,$$

$$\boldsymbol{o}_t = \boldsymbol{o}'_t \circledast \mathbf{W}_o + \mathbf{b}_o \in \mathbb{R}^d.$$

where $\mathbf{W}_c, \mathbf{W}_f, \mathbf{W}_o, \mathbf{W}_g \in \mathbb{R}^{d \times d}$ consist of ternary weights, $\boldsymbol{f}_t$ is the forget gate output, $\sigma$ is the Sigmoid activation function, $\boldsymbol{c}_t$ is the input vector, $\tau$ is the SiLU activation function, $\boldsymbol{h}_t$ is the hidden state, $\boldsymbol{g}_t$ is the output gate, $\boldsymbol{o}'_t$ is the intermediate output, and $\boldsymbol{o}_t$ is the final output at time step $t$. The initial hidden state $\boldsymbol{h}_0$ is set to $\mathbf{0}$. Similarly to HGRN, we also employ the cummax operation to bound the forget gate values in deeper layers closer to 1, though omit this above for brevity. The MLGRU can be viewed as a simplified variant of HGRN that omits complex-valued components and reduces the hidden state dimension from $2d$ to $d$. This simplification makes MLGRU more computationally efficient while preserving essential gating mechanisms and ternary weight quantization.

Recurrent LLMs are known to have limitations in long-context benchmarks and retrieval. Linear recurrent models partially address this by 1) avoiding non-linearities through time to improve gradient flow, 2) introducing data-dependent decay, and 3) introducing lower-bounds on the 'forget gate' [43]. MatMul-free LM adopts these best practices, and recent work shows that hybrid architectures with very few Transformer blocks can compensate for RNN performance across longer contexts [44].

### 6.3.2  MatMul-free Channel Mixer

For MatMul-free channel mixing, we use GLU, which is widely adopted in many modern LLMs, including Llama [45, 46, 31], Mistral [47] and RWKV [42], and a BitLinear-adapted version can be expressed as follows:

$$\boldsymbol{g}_t = \boldsymbol{x}_t \circledast \boldsymbol{W}_g \in \mathbb{R}^l,$$

$$\boldsymbol{u}_t = \boldsymbol{x}_t \circledast \boldsymbol{W}_u \in \mathbb{R}^l,$$

$$\boldsymbol{p}_t = \tau(\boldsymbol{g}_t) \odot \boldsymbol{u}_t \in \mathbb{R}^l,$$

$$\boldsymbol{d}_t = \boldsymbol{p}_t \circledast \mathbf{W}_d \in \mathbb{R}^d,$$

where $\tau$ denotes the SiLU activation function, $\circledast$ represents ternary accumulation, and $\odot$ represents the element-wise product.

The GLU consists of three main steps: 1) *upscaling* the $t$-step input $\boldsymbol{x}_t \in \mathbb{R}^d$ to $\boldsymbol{g}_t, \boldsymbol{u}_t \in \mathbb{R}^l$ using weight matrices $\mathbf{W}_g, \mathbf{W}_u \in \mathbb{R}^{d \times l}$ 2) *elementwise gating* $\boldsymbol{u}_t$ with $\boldsymbol{g}_t$ followed by a nonlinearity $f(\cdot)$, where we apply Swish [35]. 3) *Down-scaling* the gated representation $\boldsymbol{p}_t$ back to the original size through a linear transformation $\mathbf{W}_d$. Following Llama [45], we maintain the overall number of parameters of GLU at $8d^2$ by setting the upscaling factor to $\frac{8}{3}d$.

The channel mixer here only consists of dense layers, which are replaced with ternary accumulation operations. By using ternary weights in the BitLinear modules, we can eliminate the need for expensive MatMuls, making the channel mixer more computationally efficient while maintaining its effectiveness in mixing information across channels.

### 6.4  Training Details

**Surrogate Gradient**  To handle non-differentiable functions such as the Sign and Clip functions during backpropagation, we use the straight-through estimator (STE) [48] as a surrogate function for the gradient. STE allows gradients to flow through the network unaffected by these non-differentiable functions, enabling the training of our quantized model. This technique is widely adopted in BNNs and SNNs.

**Large Learning Rate**  When training a language model with ternary weights, using the same learning rate as regular models can lead to excessively small updates that have no impact on the

clipping operation. This prevents weights from being effectively updated and results in biased gradients and update estimates based on the ternary weights. To address this challenge, it is common practice to employ a larger learning rate when training binary or ternary weight language models, as it facilitates faster convergence [49, 50, 12]. In our experiments, we maintain consistent learning rates across both the 370M and 1.3B models, aligning with the approach described in Ref. [51]. Specifically, for the Transformer++ model, we use a learning rate of $3e-4$, while for the MatMul-free LM, we employ a learning rate of $4e-3, 2.5e-3, 1.5e-3$ in 370M, 1.5B and 2.7B, respectively. These learning rates are chosen based on the most effective hyperparameter sweeps for faster convergence during the training process.

**Learning Rate Scheduler**    When training conventional Transformers, it is common practice to employ a cosine learning rate scheduler and set a minimal learning rate, typically $0.1\times$ the initial learning rate. We follow this approach when training the full precision Transformer++ model. However, for the MatMul-free LM, the learning dynamics differ from those of conventional Transformer language models, necessitating a different learning strategy. We begin by maintaining the cosine learning rate scheduler and then reduce the learning rate by half midway through the training process. Interestingly, we observed that during the final training stage, when the network's learning rate approaches 0, the loss decreases significantly, exhibiting an *S*-shaped loss curve. This phenomenon has also been reported by [12, 49] when training binary/ternary language models.

## 6.5    Implementation Details on Intel Loihi 2

In order to run the MatMul-free LM on the Intel Loihi 2 chip, the original model must be quantized, implemented using fixed-point arithmetic using custom microcode on the Loihi 2 platform, and finally verified against the original model to make sure that the mismatch between the fixed-point model on Loihi 2 and the original model on GPU is limited. We present these steps in the following sections.

**Model Quantization**    As a first step, we quantize all weights and activations of the MatMul-free LM and verify the resulting performance of the quantized model on the downstream tasks presented previously. All weight matrices are already ternarized, thus not needing further quantization. However, every BitLinear layer with a ternary weight matrix $W \in \mathbb{R}^{d\times m}$ is preceded by a root mean square normalization (RMSNorm) operation, given by:

$$\text{RMSNorm}(x; w, \epsilon) = \frac{x}{\sqrt{\epsilon + \sum_i^d x_i^2}} \odot w \tag{6}$$

where $w \in \mathbb{R}^d$ is a learned vector of element-wise scaling factors, $\epsilon$ is a small constant set to $\epsilon = 10^{-6}$ to avoid zero division, and $x \in \mathbb{R}^d$ is the input activation vector. While the original MatMul-free model stores all normalization scales $w \in \mathbb{R}^d$ as 16-bit floating-point numbers, we quantize all normalization scales across the model to 8-bit integers using $w_q = \text{round}\,(w/s_w)$ where $s_w$ is the quantization scale, chosen as $s_w = 2^{\lfloor \log_2(\max(|w|)) \rceil}$ where $\lfloor \cdot \rceil$ denotes the rounding operation, and $\max(|x|)$ is the absolute maximum activation observed on the data across all activation channels. These quantization scales allow the rescaling between different quantization scales to be efficiently implemented on Loihi 2 using bit-shift operations. Quantizing all weights in this scheme results in a relative performance drop of 0.18% averages across all downstream tasks reported in Section 4.

Whereas activations of the MatMul-free LM are stored as 16-bit floating-point numbers, activations on Loihi 2 are sent between neurons as integer payloads. To test the performance degradation of the model when quantization activations to lower bit precision, the model was quantized using dynamic quantization in PyTorch, where quantization scales for the activations $s_a$ are computed for each batch individually. When quantizing activations in the RMSNorm layer, the $\epsilon$ value underflows in most layers, potentially leading to division by zero. We therefore set $\epsilon = 10^{-3}$ and report results with this modified value. Table 2 shows that for the 370M model, performance drops by less than 1.5% relative to the original model when using 16-bit activation quantization. Performance drops more significantly, by over 12% relative to the original model, when using 8-bit activations. All further results will assume 8-bit weights and 16-bit activations. We note that these results are for naive post-training quantization and may be significantly improved through quantization-aware training or fine-tuning the model with specific quantization requirements.

**Fixed-Point Arithmetic** Once all quantization scales are computed, the model can be run in static quantization with integer weights and activations. However, some operations used in the MatMul-free LM are not defined on integers, namely the sigmoid activation function $\sigma$ and the inverse-square-root in the RMSNorm.

We employ a look-up-table (LUT) for a fixed-point approximation of the logistic sigmoid function, $\sigma(x) = 1/(1 + e^{-x})$. Specifically, we scale the floating-point input $x$ by $2^{x_{\exp}}$ where $x_{\exp} = 6$, quantize it to an integer domain, and store precomputed $\lfloor \sigma\left(\frac{x}{2^{x_{\exp}}}\right) \cdot 2^{y_{\exp}} \rfloor$ values in a LUT for positive inputs. For negative inputs, we exploit $\sigma(-x) = 1 - \sigma(x)$, thus requiring only half-range values. During inference, a piecewise linear interpolation between adjacent LUT entries refines the output. This design offers efficient computation and controllable approximation error.

For the inverse-square-root in the RMSNorm layer, we adapted a well-known "fast inverse square root" algorithm `Fast InvSqrt` to operate entirely in fixed-point arithmetic. Our method treats the input $\tilde{x}$ as an integer paired with a fixed exponent, and uses a LUT with 24 values to produce an initial guess for $\sqrt{\tilde{x}}$. This estimate is then refined using five iterations of the Newton-Raphson method, all in a fixed-point format. See Appendix C.2 for more details.

**Mapping the MatMul-free LM onto Loihi** The model on Loihi 2 is implemented as a set of neurons that are connected through synapses. Each neuron asynchronously executes a simple microcode program and sends its output to a collection of other neurons through synapses. Because every neuron contains only information pertaining to its own dynamics, aggregate operations–such as the sum of squares over an activation vectors–must be implemented through an additional neuron that receives inputs from all neurons in the relevant layer. Fig. 9 in Appendix C.4.2 shows the partitioning of the MatMul-free model into Loihi-compatible neurons.

# References

[1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.

[2] Sara Hooker. The hardware lottery. *Communications of the ACM*, 64(12):58–65, 2021.

[3] Hanting Chen, Yunhe Wang, Chunjing Xu, Zhaohui Yang, Chuanjian Liu, Boxin Shi, Chao Xu, Chunfeng Xu, and Qi Tian. The addernet: Do we really need multiplications in deep learning? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1468–1477, 2020.

[4] Haoran You, Yipin Guo, Yichao Fu, Wei Zhou, Huihong Shi, Xiaofan Zhang, Souvik Kundu, Amir Yazdanbakhsh, and Yingyan Celine Lin. Shiftaddllm: Accelerating pretrained llms via post-training multiplication-less reparameterization. *Advances in Neural Information Processing Systems*, 37:24822–24848, 2024.

[5] Wolfgang Maass. Networks of spiking neurons: the third generation of neural network models. *Neural networks*, 10(9):1659–1671, 1997.

[6] Jason K Eshraghian, Max Ward, Emre O Neftci, Xinxin Wang, Gregor Lenz, Girish Dwivedi, Mohammed Bennamoun, Doo Seok Jeong, and Wei D Lu. Training spiking neural networks using lessons from deep learning. *Proceedings of the IEEE*, 2023.

[7] Rui-Jie Zhu, Qihang Zhao, Guoqi Li, and Jason K Eshraghian. SpikeGPT: Generative pre-trained language model with spiking neural networks. *arXiv preprint arXiv:2302.13939*, 2023.

[8] Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1. *arXiv preprint arXiv:1602.02830*, 2016.

[9] Sreyes Venkatesh, Razvan Marinescu, and Jason K Eshraghian. Squat: Stateful quantization-aware training in recurrent spiking neural networks. *arXiv preprint arXiv:2404.19668*, 2024.

[10] Jason K Eshraghian, Xinxin Wang, and Wei D Lu. Memristor-based binarized spiking neural networks: Challenges and applications. *IEEE Nanotechnology Magazine*, 16(2):14–23, 2022.

[11] Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Huaijie Wang, Lingxiao Ma, Fan Yang, Ruiping Wang, Yi Wu, and Furu Wei. Bitnet: Scaling 1-bit transformers for large language models. *arXiv preprint arXiv:2310.11453*, 2023.

[12] Shuming Ma, Hongyu Wang, Lingxiao Ma, Lei Wang, Wenhui Wang, Shaohan Huang, Li Dong, Ruiping Wang, Jilong Xue, and Furu Wei. The era of 1-bit llms: All large language models are in 1.58 bits. *arXiv preprint arXiv:2402.17764*, 2024.

[13] Falcon-LLM Team. Falcon-e, a series of powerful, universal and fine-tunable 1.58bit language models., April 2025.

[14] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359, 2022.

[15] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.

[16] Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what you need for vision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10819–10829, 2022.

[17] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.

[18] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.

[19] Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, et al. Baichuan 2: Open large-scale language models. *arXiv preprint arXiv:2309.10305*, 2023.

[20] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[21] Vikas Yadav, Steven Bethard, and Mihai Surdeanu. Quick and (not so) dirty: Unsupervised selection of justification sentences for multi-hop question answering. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *EMNLP-IJCNLP*, 2019.

[22] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: can a machine really finish your sentence? In *Proceedings of the 57th Conference of the Association for Computational Linguistics*, pages 4791–4800, 2019.

[23] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. WinoGrande: an adversarial winograd schema challenge at scale. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence*, pages 8732–8740, 2020.

[24] Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. PIQA: reasoning about physical commonsense in natural language. *CoRR*, abs/1911.11641, 2019.

[25] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? A new dataset for open book question answering. *CoRR*, abs/1809.02789, 2018.

[26] Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, 12 2023.

[27] Lei Wang, Lingxiao Ma, Shijie Cao, Quanlu Zhang, Jilong Xue, Yining Shi, Ningxin Zheng, Ziming Miao, Fan Yang, Ting Cao, Yuqing Yang, and Mao Yang. Ladder: Enabling efficient low-precision deep learning computing through hardware-aware tensor transformation. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*, 2024.

[28] Mike Davies, Andreas Wild, Garrick Orchard, Yulia Sandamirskaya, Gabriel A. Fonseca Guerra, Prasad Joshi, Philipp Plank, and Sumedh R. Risbud. Advancing Neuromorphic Computing With Loihi: A Survey of Results and Outlook. *Proceedings of the IEEE*, 109(5):911–934, 2021. Publisher: Institute of Electrical and Electronics Engineers (IEEE).

[29] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report, 2024.

[30] Michele Montebovi. Alireo-400m: A lightweight italian language model, 2024.

[31] AI@Meta. Llama 3 model card. 2024.

[32] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[33] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.

[34] Yann N Dauphin, Angela Fan, Michael Auli, and David Grangier. Language modeling with gated convolutional networks. In *International conference on machine learning*, pages 933–941. PMLR, 2017.

[35] Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.

[36] Jiayi Pan, Chengcan Wang, Kaifu Zheng, Yangguang Li, Zhenyu Wang, and Bin Feng. Smoothquant+: Accurate and efficient 4-bit post-training weightquantization for llm. *arXiv preprint arXiv:2312.03788*, 2023.

[37] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, pages 38087–38099. PMLR, 2023.

[38] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[39] Aaron Voelker, Ivana Kajić, and Chris Eliasmith. Legendre memory units: Continuous-time representation in recurrent neural networks. *Advances in neural information processing systems*, 32, 2019.

[40] Antonio Orvieto, Samuel L Smith, Albert Gu, Anushan Fernando, Caglar Gulcehre, Razvan Pascanu, and Soham De. Resurrecting recurrent neural networks for long sequences. In *International Conference on Machine Learning*, pages 26670–26698. PMLR, 2023.

[41] Soham De, Samuel L Smith, Anushan Fernando, Aleksandar Botev, George Cristian-Muraru, Albert Gu, Ruba Haroun, Leonard Berrada, Yutian Chen, Srivatsan Srinivasan, et al. Griffin: Mixing gated linear recurrences with local attention for efficient language models. *arXiv preprint arXiv:2402.19427*, 2024.

[42] Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, Kranthi Kiran GV, et al. Rwkv: Reinventing rnns for the transformer era. *arXiv preprint arXiv:2305.13048*, 2023.

[43] Zhen Qin, Songlin Yang, and Yiran Zhong. Hierarchically gated recurrent neural network for sequence modeling. *Advances in Neural Information Processing Systems*, 36, 2024.

[44] Dustin Wang, Rui-Jie Zhu, Steven Abreu, Yong Shan, Taylor Kergan, Yuqi Pan, Yuhong Chou, Zheng Li, Ge Zhang, Wenhao Huang, and Jason Eshraghian. A systematic analysis of hybrid linear attention. *arXiv preprint arXiv:2507.06457*, 2025.

[45] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[46] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[47] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.

[48] Yoshua Bengio, Nicholas Léonard, and Aaron C. Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *CoRR*, abs/1308.3432, 2013.

[49] Yichi Zhang, Ankush Garg, Yuan Cao, Lukasz Lew, Behrooz Ghorbani, Zhiru Zhang, and Orhan Firat. Binarized neural machine translation. *Advances in Neural Information Processing Systems*, 36, 2024.

[50] Zechun Liu, Barlas Oguz, Aasish Pappu, Yangyang Shi, and Raghuraman Krishnamoorthi. Binary and ternary natural language generation. *arXiv preprint arXiv:2306.01841*, 2023.

[51] Zhen Qin, Dong Li, Weigao Sun, Weixuan Sun, Xuyang Shen, Xiaodong Han, Yunshen Wei, Baohong Lv, Fei Yuan, Xiao Luo, et al. Scaling transnormer to 175 billion parameters. *arXiv preprint arXiv:2307.14995*, 2023.

[52] Garrick Orchard, E Paxon Frady, Daniel Ben Dayan Rubin, Sophia Sanborn, Sumit Bam Shrestha, Friedrich T Sommer, and Mike Davies. Efficient neuromorphic signal processing with loihi 2. In *2021 IEEE Workshop on Signal Processing Systems (SiPS)*, pages 254–259. IEEE, 2021.

[53] Sumit Bam Shrestha, Jonathan Timcheck, Paxon Frady, Leobardo Campos-Macias, and Mike Davies. Efficient Video and Audio Processing with Loihi 2. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 13481–13485, April 2024. ISSN: 2379-190X.

# A  FPGA Implementation and Results

Neither GPUs nor Intel Loihi 2 support 2-bit weights, so we additionally did an RTL-only implementation of the MatMul-free LM on a D5005 Stratix 10.

## A.1  Implementation

To test the power usage and effectiveness of the MatMul-free LM on custom hardware that can better exploit ternary operations, we created an FPGA accelerator in SystemVerilog. The overview is shown in Fig. 5.
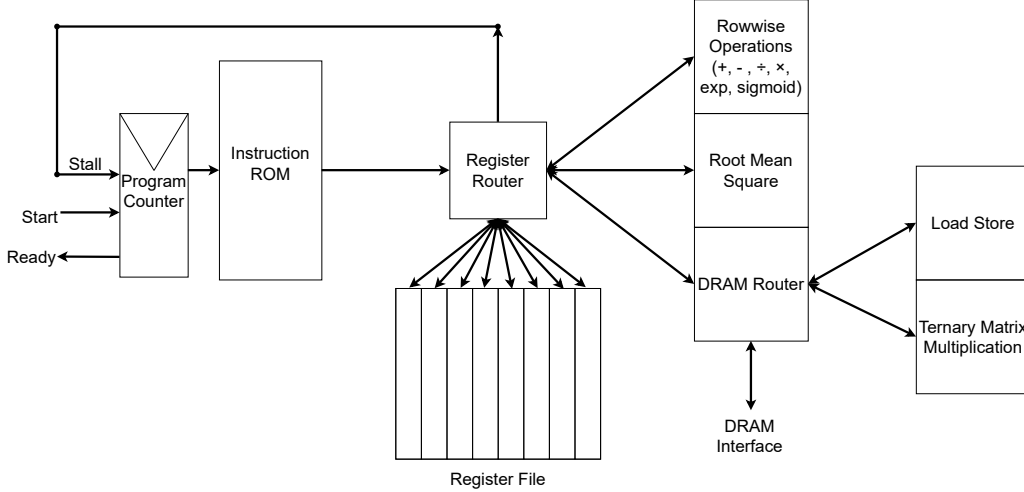


Figure 5: RTL implementation for running MatMul-free token generation

There are 4 functional units in this design: "Row-wise Operation," "Root Mean Square," "Load Store," and "Ternary Matrix Multiplication," and they each allow for simple out-of-order execution. We wrote a custom assembler for our custom instruction set, which was used to convert assembly files into an instruction ROM. The custom instruction set is given below:

- LDV: LoaD Vector from memory
- SV: Store Vector to memory
- ADD: row-wise ADDition
- SUB: row-wise SUBtraction
- MUL: row-wise MULtiplication
- DIV: row-wise DIVision
- EXP: row-wise EXPonential function
- SIG: row-wise SIGmoid
- NORM: NORMalization with root-mean-square
- TMATMUL: Ternary MATrix MULtiplication

**The register router** delegates incoming instructions to available registers. The register file consists of 8 registers, each storing 1 vector in a separate SRAM array. Each register SRAM array has a read and write port that are delegated to at most one instruction at a time. If an instruction requests access to a functional unit or a register that is busy, the program counter will stall until the functional unit or register has been freed. If two instructions do not block each other, they execute simultaneously.

Table 4: MatMul-free token generation FPGA core resource utilization and performance metrics. The ternary matrix multiplication operation dominates latency for the current implementation and there is not an observed bottleneck in the local DDR4 bridge. In future implementations, this functional unit will be optimized and the DDR interface will likely become the primary bottleneck.

| | %ALMs | | %M20Ks | | Avg Power (W) | | Latency (ms) | |
| Core Count | Core | Total | Core | Total | Core Active | Core Idle | Core | DDR4 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 1 | 2.9 | 9 | 0.01 | 2.87 | 13.67 | 13.68 | 46.36 | 0.09 |
| 8 | 23.21 | 26.9 | 0.08 | 3.06 | 39.78 | 39.94 | 46.36 | 0.18 |
| 16 | 46.43 | 50.1 | 0.15 | 5.13 | 75.25 | 73.97 | 46.36 | 0.72 |
| 26 | 75.45 | 100 | 0.25 | 22.64 | 166.30 | 149.66 | 46.36 | 5.76 |

**The "Root Mean Square" functional unit** uses a specialized hardware algorithm to preserve precision, and runs in 3 stages. Stage 1 will copy the target vector to an internal-temporary register, and perform a square on each element using a lookup-table. Stage 2 will divide-and-conquer to average neighboring vector elements, generating the Root-Mean-Square result. Stage 3 will perform normalization by dividing each element in the original vector by the Root-Mean-Square result. By using divide-and-conquer for averaging, instead of a typical rolling sum then large divide, rounding errors are significantly reduced.

**The "Ternary Matrix Multiplication" functional unit** takes in a DRAM address for a ternary matrix, then performs a TMATMUL on the specified vector. Our architecture entirely places the ternary matrices in DRAM. While running a TMATMUL instruction, an SRAM FIFO is simultaneously filled with sequential DRAM fetch results, and emptied by a power-efficient ternary-add operation. At the moment, the three required TMATMUL instructions take up nearly all of the total execution time. In future work, we will introduce parallelism and caching to improve TMATMUL execution time.

## A.2 Results

The RTL implementation of the MatMul-free token generation core is deployed on a D5005 Stratix 10 programmable acceleration card (PAC) in the Intel FPGA Devcloud. The core completes a forward-pass of a block in 43ms at $d = 512$ and achieves a clock rate of 60MHz. The resource utilization, power and performance of the single-core implementation of a single block ($N = 1$) are shown in Tab. 4. '% ALM Core' refers to the percentage of the total adaptive logic modules used by the core logic, and '%ALM Total' includes the core, the additional interconnect/arbitration logic, and "shell" logic for the FPGA Interface Manager. 'M20K' refers to the utilization of the memory blocks, and indicates that the number of cores are constrained by ALMs, and not on-chip memory (for this DDR implementation). We implement a single token generation core, and estimate the total number of cores that could fit on the platform and the corresponding power, performance and area impact. This is the simplest case where the core only receives 8 bits at a time from memory.

**The single core implementation** exhibits extremely low dynamic power that is hardly distinguishable from power measured while the core is inactive. Each core requires access to a DDR4 interface and MMIO bridges for host control. In this implementation, the majority of resources are dedicated to the provided shell logic and only 0.4% of programmable logic resources are dedicated to logic for core interconnect and arbitration to DDR4 interfaces/MMIO. As described above, the core latency is primarily due to the larger execution time of the ternary matrix multiply functional unit.

By instead using the full 512-bit DDR4 interface and parallelizing the TMATMUL functional unit, which dominates 99% of core processing time, a further speed-up of approximately 64× is projected, while maintaining the same clock rate without additional optimizations or pipelining, as shown in Table 5. Given the 370M parameter model where $L = 24$, $d = 512$, the total projected runtime is 16.08ms, and a throughput of approximately 62 tokens per second. The 1.3B parameter model, where $L = 24$ and $d = 2048$, has a projected runtime of 42ms, and a throughput of 23.8 tokens per second. Despite not being optimized for 2-b precision weights, Intel Loihi 2 results outperform the

Table 5: FPGA Performance Metrics for Different Embedding Dimensions ($d$)

| $d$ | Runtime (ms) | Projected Runtime w/Bursting (ms) | Power (W) Idle | Active | ALM Utilization (%) Core | Total | Clock (MHz) |
|---|---|---|---|---|---|---|---|
| 512 | 43 | 0.67 | 13.36 | 13.39 | 2.8 | 9 | 60 |
| 1024 | 112 | 1.75 | 13.64 | 13.65 | 5.7 | 11 | 54 |
| 2048 | 456 | 7.13 | 13.92 | 13.93 | 11 | 16 | 52 |

FPGA implementation in terms of both throughput and power likely due to a mix of efficient packet routing, deep optimizations for recurrent state updates, unused digital signal processing (DSP) units on the FPGA, and process node deviations. This is for the case of a single core with a single batch of data, and can be scaled up significantly through batch processing by pipelining the single core with a negligible increase in average power, or alternatively, by increasing the core count with an increase in power (Table 4).

**Estimates of multi-core implementation latencies** are generated by scaling the overheads of the single core implementation and factoring in the growth of logic to accommodate contention on the DDR4 channels. Each core connects to one of four DDR4 channels, and each additional core connected to a channel will double the required arbitration and buffering logic for that channel. As both the host and core share DDR4 channels, this overhead will scale proportional to the number of cores attached to the channel. To mitigate this, future work could bring additional caching optimizations to the core and functional units. Core latency is the compute time of the core from start to ready and DDR4 latency is the required time to transfer input vectors from the host to the PAC local DDR4.

**Estimates of multi-core implementation power** are calculated by scaling the measured power of a single-core implementation. Idle power is estimated by scaling the total estimated resource overhead of all additional logic added to a constant estimate of idle power consumed by the platform shell. The single-core active power is scaled by the additional arbitration, interconnect and core overhead. We assume a constant clock rate for all implementations.

We note that the FPGA implementation is done in RTL from top to bottom, and there are many optimizations that could be added. For example, we are not using any vendor-provided IPs, and we are not bursting DDR transactions, both of which would significantly accelerate operation. This approach is to achieve the most generic and cross-platform evaluation possible.

## B   Proof of Linear Scaling in the Proposed MatMul-free Linear Gated Recurrent Unit

The MatMul-free Linear Gated Recurrent Unit (MLGRU) presents a more efficient structure compared to the conventional Key-Value (KV) Cache mechanism used in Transformer models. To demonstrate its advantages, we examine both models in terms of computational complexity and memory requirements with respect to the sequence length $T$.

In Transformers, the KV Cache mechanism stores all keys and values from previous time steps for autoregressive generation. At each time step $t$, the model requires three steps: computing the query, key, and value vectors for the input $\boldsymbol{x}_t$, followed by calculating the attention weights by comparing the query $\boldsymbol{q}_t$ with all previous keys, and finally constructing the context vector by a weighted sum of the previous values. Formally, the query, key, and value are calculated as

$$\boldsymbol{q}_t = \boldsymbol{x}_t \mathbf{W}_q, \quad \boldsymbol{k}_t = \boldsymbol{x}_t \mathbf{W}_k, \quad \boldsymbol{v}_t = \boldsymbol{x}_t \mathbf{W}_v,$$

where $\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v \in \mathbb{R}^{d_{\text{model}} \times d_{\text{attn}}}$ are weight matrices. The attention weights, calculated as

$$\alpha_{ti} = \frac{\exp\left(\frac{\boldsymbol{q}_t^\top \boldsymbol{k}_i}{\sqrt{d_{\text{attn}}}}\right)}{\sum_{j=1}^{t} \exp\left(\frac{\boldsymbol{q}_t^\top \boldsymbol{k}_j}{\sqrt{d_{\text{attn}}}}\right)}, \quad \text{for } i = 1, \ldots, t,$$

lead to the context vector

$$\boldsymbol{z}_t = \sum_{i=1}^{t} \alpha_{ti} \boldsymbol{v}_i.$$

These operations grow linearly in complexity with $t$, producing $O(t \cdot d_{\text{attn}})$ computational steps per time step.

The computational complexity of the KV Cache over a sequence of length $T$ therefore follows a quadratic scaling pattern:

$$\sum_{t=1}^{T} O(t \cdot d_{\text{attn}}) = O(d_{\text{attn}} \cdot T^2).$$

The memory requirement, dictated by the storage of all previous keys and values, also scales linearly with $T$, resulting in $O(T \cdot d_{\text{attn}})$.

In contrast, the proposed MLGRU architecture achieves simplicity and efficiency by eliminating matrix multiplications and relying on element-wise ternary-weight operations, making it a MatMul-free structure. This architecture uses a recurrent update rule for the hidden state $\boldsymbol{h}_t$, leveraging the previous hidden state $\boldsymbol{h}_{t-1}$ and the current input $\boldsymbol{x}_t$ without retaining past sequence information. The key equations include the forget gate

$$\boldsymbol{f}_t = \sigma\left(\boldsymbol{x}_t \circledast \mathbf{W}_f + \boldsymbol{b}_f\right),$$

where $\sigma$ represents the sigmoid function and $\circledast$ denotes element-wise multiplication with ternary weights $\mathbf{W}_f \in \{-1, 0, 1\}^{d \times d}$ and bias $\boldsymbol{b}_f \in \mathbb{R}^d$. The candidate hidden state is given by

$$\boldsymbol{c}_t = \tau\left(\boldsymbol{x}_t \circledast \mathbf{W}_c + \boldsymbol{b}_c\right),$$

where $\tau$ is the SiLU activation. The hidden state update, simplified through element-wise operations, is defined as

$$\boldsymbol{h}_t = \boldsymbol{f}_t \odot \boldsymbol{h}_{t-1} + (1 - \boldsymbol{f}_t) \odot \boldsymbol{c}_t.$$

The output computations, finalized by applying a sigmoid activation to the input and combining it with the hidden state, result in

$$\boldsymbol{g}_t = \boldsymbol{x}_t \circledast \mathbf{W}_g + \boldsymbol{b}_g, \quad \boldsymbol{o}_t = (\boldsymbol{g}_t \odot \sigma(\boldsymbol{h}_t)) \circledast \mathbf{W}_o + \boldsymbol{b}_o.$$

Each step requires only constant-time element-wise operations, establishing $O(d)$ computational complexity per step.

Over a sequence of length $T$, MLGRU achieves total computational complexity of

$$\sum_{t=1}^{T} O(d) = O(d \cdot T).$$

Memory requirements remain constant at $O(d)$ since only the current hidden state is maintained. Thus, MLGRU's efficiency is clearly shown by its linear computational complexity and constant memory demand, contrasting the quadratic growth and linear memory of the KV Cache.

The recursion in MLGRU allows the hidden state $\boldsymbol{h}_t$ to encapsulate all historical information without explicit storage, as shown by

$$\boldsymbol{h}_t = \Phi(\boldsymbol{h}_{t-1}, \boldsymbol{x}_t),$$

where $\Phi$ represents the recursive update rule using element-wise ternary operations. This recursive nature ensures $\boldsymbol{h}_t$ implicitly reflects the sequence $\{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_t\}$ while avoiding the need for an accumulating memory footprint.

## C    Loihi 2: Detailed Implementation and Experimental Results

This section of the appendix provides a comprehensive overview of the methodologies used for implementing the MatMul-free Language Model on Intel's Loihi 2 neuromorphic research chip, along with a detailed presentation of the experimental results. We elaborate on the relevant aspects of the Loihi 2 hardware architecture, the specifics of fixed-point conversion and custom function implementation critical for on-chip execution, and the performance benchmarks achieved, including throughput and energy efficiency. The subsequent subsections delve into these areas in greater detail.

## C.1 Loihi 2 Hardware Architecture

Intel's second-generation neuromorphic research processor, Loihi 2, was purpose-built for sparse, event-based neural networks [52]. Each chip comprises 120 massively parallel compute units called *neuro-cores* that can be scaled to systems containing up to 1,152 chips (Figure 6). These neuro-cores operate asynchronously while maintaining global algorithmic time steps via barrier synchronization. Co-located memory enables efficient local state updates, supporting up to 8,192 stateful neurons per core. Users can program each neuron's temporal dynamics through assembly code, with flexible memory allocation achieved by trading neuron count for memory per core. External communication supports up to 160 million 32-bit integer messages per second [53]. Fully-digital stacked systems like Hala Point can scale to 1 billion neurons and 128 billion synapses (Figure 6). Networks must use



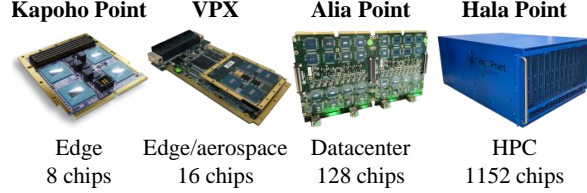| Kapoho Point | VPX | Alia Point | Hala Point |
|---|---|---|---|
| Edge | Edge/aerospace | Datacenter | HPC |
| 8 chips | 16 chips | 128 chips | 1152 chips |

Figure 6: Different Loihi 2 systems are available to cover a wide range of applications from the edge to HPC with up to 1 billion neurons.

fixed-point arithmetic: 8-bit synaptic weights[3] and up to 32-bit messages[4]. Unlike GPUs, Loihi 2 prioritizes local neuronal computations—a defining characteristic of neuromorphic processors. This design enables: (1) fast, efficient recurrent state updates with minimal data movement through neuro-core-local memory; (2) efficient processing of unstructured sparse weight matrices via asynchronous event-driven architecture; and (3) exploitation of activation sparsity, as asynchronous communication transfers only non-zero messages.

### C.1.1 Execution Modes on Loihi 2

Loihi 2's asynchronous architecture enables dynamic throughput-latency optimization (Fig. 7). *Pipelined mode* maximizes throughput by injecting new inputs every time step and forwarding them through neuronal layers. *Fall-through mode* minimizes latency by introducing new inputs only after complete processing of previous inputs. LLM deployment naturally maps to these modes: prefill processing of long sequences utilizes pipelined mode for optimal throughput, while autoregressive generation employs fall-through mode since token generation at time $t$ must complete before processing token $t + 1$.

## C.2 Fixed-Point Implementation Details

### C.2.1 Fixed-Point Implementation of the Sigmoid Function

We implement the logistic sigmoid function $\sigma(x) = 1/(1 + e^{-x})$ using fixed-point arithmetic with an interpolated look-up table (LUT). The floating-point input $x$ is scaled by $2^{x_{\exp}}$ (where $x_{\exp} = 6$) and quantized to integer domain: $x_{\text{fxp}} = \lfloor x2^{x_{\exp}} \rceil$. The LUT stores precomputed values $\left( x_{\text{int}}^{\text{lut,i}} \mapsto y_{\text{int}}^{\text{lut,i}} \right)_{i \in \{0, \dots N_\sigma\}}$ with $N_\sigma = 8$ entries:

$$y_{\text{int}}^{\text{lut,i}} = \lfloor \sigma(\frac{x_{\text{int}}^{\text{lut,i}}}{2^{x_{\exp}}}) \cdot 2^{y_{\exp}} \rfloor \tag{7}$$

Only positive inputs are stored, as negative values are computed via $\sigma(-x) = 1 - \sigma(x)$, halving memory requirements. Piecewise linear interpolation between adjacent entries enhances output precision. Parameters $x_{\exp}$ and $N_\sigma$ control the computation-accuracy tradeoff.

---

[3]This is not a hard limit, as an $8n$-bit synapse can be implemented through $n$ separate 8-bit synapses that are added together with different fixed-point exponents.

[4]Local states are not restricted in precision, and one may also transmit messages with more than 32 bits in an analogous way to what is described above for synaptic weights.
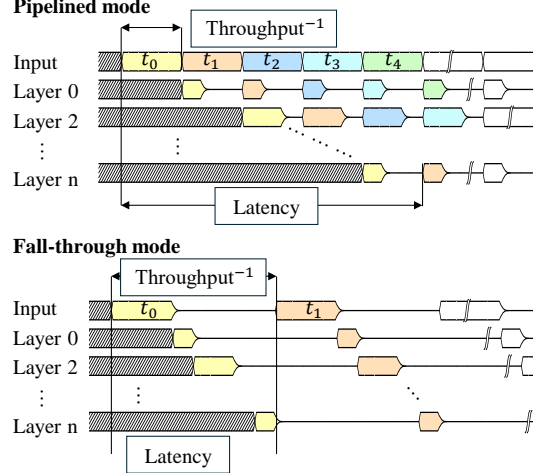
Figure 7: Different execution modes on Loihi 2 that either optimize throughput or latency. In the *pipelined mode*, a new data point is inserted in each time step, to use all processing cores and maximize the throughput–at the expense of latency because equal time bins $t_0 = t_1 = \ldots$ are enforced. In the *fall-through mode*, a new data points is only provided once the last data point has been fully processed with minimum latency. Only a single neuronal layer is active at any step as data travels through the network. The time per step is thus minimized as traffic is reduced and potentially more complex neuronal layers are not updated.

### C.2.2 Fixed point implementation of the inverse square root

For the inverse-square-root in the RMSNorm layer, we adapted a well-known "fast inverse square root" algorithm `FastInvSqrt` to operate entirely in fixed-point arithmetic. Our method treats the input $\tilde{x}$ as an integer paired with a fixed exponent, and uses a LUT with 24 values to produce an initial guess for $\sqrt{\tilde{x}}$. This estimate is then refined using five iterations of the Newton-Raphson method, all in a fixed-point format.

### C.3 Double RMSNorm Derivation

Let $x \in \mathbb{R}^d$ and $y = \text{RMSNorm}(x)$ and $z = \text{RMSNorm}(y)$, in expanded form:

$$y = \frac{x}{\sqrt{\epsilon + \sum_i^d x_i^2}} \odot g_1 \tag{8}$$

$$z = \frac{y}{\sqrt{\epsilon + \sum_i^d y_i^2}} \odot g_2 \tag{9}$$

We wish to derive an equation for $z = \text{DoubleRMSNorm}(x) = \text{RMSNorm}(\text{RMSNorm}(x))$.

First, we express $\mu_{\mathbf{y}}$ in terms of $\mu_{\mathbf{x}}$:

$$\mu_{\mathbf{y}} = \frac{1}{D} \sum_{i=1}^{D} y_i^2 = \frac{1}{D} \sum_{i=1}^{D} \left( x_i \cdot \frac{g_1}{\sqrt{\mu_{\mathbf{x}} + \varepsilon}} \right)^2 \tag{10}$$

$$= \left( \frac{g_1^2}{\mu_{\mathbf{x}} + \varepsilon} \right) \cdot \left( \frac{1}{D} \sum_{i=1}^{D} x_i^2 \right) = \frac{g_1^2 \mu_{\mathbf{x}}}{\mu_{\mathbf{x}} + \varepsilon}. \tag{11}$$

We then express $\mathbf{z}$ in terms of $\mathbf{x}$ by plugging in the equation for $\mathbf{y}$:

$$\mathbf{z} = \mathbf{y} \cdot \frac{g_2}{\sqrt{\mu_{\mathbf{y}} + \varepsilon}} = \left( \mathbf{x} \cdot \frac{g_1}{\sqrt{\mu_{\mathbf{x}} + \varepsilon}} \right) \cdot \frac{g_2}{\sqrt{\mu_{\mathbf{y}} + \varepsilon}} \tag{12}$$

$$= \mathbf{x} \cdot \frac{g_1 g_2}{\sqrt{(\mu_{\mathbf{x}} + \varepsilon)(\mu_{\mathbf{y}} + \varepsilon)}}. \tag{13}$$

We simplify the denominator:

$$\sqrt{(\mu_{\mathbf{x}} + \varepsilon)(\mu_{\mathbf{y}} + \varepsilon)} = \sqrt{(\mu_{\mathbf{x}} + \varepsilon) \cdot \frac{g_1^2 \mu_{\mathbf{x}}}{\mu_{\mathbf{x}} + \varepsilon} + \varepsilon} \tag{14}$$

$$= \sqrt{(\mu_{\mathbf{x}} + \varepsilon) \cdot \frac{\mu_{\mathbf{x}}(g_1^2 + \varepsilon) + \varepsilon^2}{\mu_{\mathbf{x}} + \varepsilon}} \tag{15}$$

$$= \sqrt{\mu_{\mathbf{x}}(g_1^2 + \varepsilon) + \varepsilon^2}. \tag{16}$$

We then derive the final expression for $\mathbf{z}$:

$$\mathbf{z} = \mathbf{x} \cdot \frac{g_1 g_2}{\sqrt{\mu_{\mathbf{x}}(g_1^2 + \varepsilon) + \varepsilon^2}}. \tag{17}$$

This provides the combined RMSNorm operation with different scaling parameters $g_1$ and $g_2$. By combining two RMSNorm operations with different scaling parameters, we arrive at a single normalization step:

$$\mathbf{z} = \mathbf{x} \cdot \frac{g_{\text{combined}}}{\sqrt{\mu_{\mathbf{x}} + \varepsilon_{\text{combined}}}}, \tag{18}$$

where:

$$g_{\text{combined}} = \frac{g_1 g_2}{\sqrt{g_1^2 + \varepsilon}}, \quad \varepsilon_{\text{combined}} = \frac{\varepsilon^2}{g_1^2 + \varepsilon}. \tag{19}$$

Alternatively, since the denominator depends on $\mu_{\mathbf{x}}$, it may not be possible to express $\varepsilon_{\text{combined}}$ independently without further approximations.

## C.4    Detailed Hardware Results

Table 6 presents comprehensive performance comparisons between the MatMul-free LLM on Loihi 2 and transformer-based LLMs on NVIDIA Jetson Orin Nano, additionally including H100 GPU results for both MatMul-free LLM and Transformer++ baseline. The metrics for Loihi 2 are based on experiments on multi-chip systems. See Appendix C.4.1 for more details. The MatMul-free LLM on GPU demonstrates improved throughput and efficiency at longer sequences, attributed to linear token mixing scaling and enhanced GPU utilization. Conversely, the Transformer++ baseline exhibits initial throughput gains at shorter sequences through improved hardware utilization, followed by performance degradation at longer sequences due to self-attention's quadratic scaling. The MatMul-free LLM on Loihi 2 shows constant scaling in both throughput and efficiency for different sequence lengths. This is due to full utilization of the chip at sequences as short as 500, and constant scaling of the MatMul-free token mixer.

### C.4.1    Detailed Loihi 2 Results

**Single chip experiments** As detailed in Section 5.3, Loihi 2 energy and throughput metrics were derived from experiments deployed on a single MatMul-free LM block on an Oheo Gulch single-chip Loihi 2 system. We measured average time per step (TPS, $T_{\text{TPS}}$), representing single execution timestep duration. Total model latency (time-to-first-token, $T_{\text{ttft}}$) was calculated as $T_{\text{ttft}} = N_{\text{blocks}} \times N_{\text{steps/block}} \times T_{\text{TPS}}$, where $N_{\text{blocks}} = 24$ for the 370M MatMul-free language model.

Prefill operations employ pipelined mode with constant TPS due to enforced equal time bins (Appendix C.1.1). Prefill throughput follows $f_{\text{prefill}} = T_{\text{TPS}}^{-1}$, as tokens are processed at interval $T_{\text{TPS}}$. Single-chip power measurements comprise static and dynamic components: $P^{\text{1-chip}} = \tilde{P}^{\text{1-chip}} + \bar{P}^{\text{1-chip}}$, where $\bar{P}$ and $\tilde{P}$ denote static and dynamic power respectively. Total prefill power estimates yield $\hat{P}_{\text{prefill}} = 24 \times P^{\text{1-chip}}$, with energy per token calculated as $\hat{E}_{\text{prefill}} = \hat{P}_{\text{prefill}} * T_{\text{TPS}}$. Fig. 8 illustrates single-chip static and dynamic power measurements.

Generation utilizes fallthrough mode where TPS varies temporally, reflecting current operation latency (Appendix C.1.1). Average TPS across >1000 timesteps determines $T_{\text{TPS}}$. Generation throughput equals $f_{\text{generate}} = T_{\text{ttft}}^{-1}$, typically lower than $f_{\text{prefill}}$. Power measurements follow prefill methodology,

Table 6: Throughput and energy efficiency for two transformer-based language models running on the NVIDIA Jetson Orin Nano compared to our MatMul-free LM running on Intel's Loihi 2, across different sequence lengths for prefill and generation. The best-performing sequence length for each model and metric is <u>underlined</u>.

| | | Throughput (↑ tokens/sec) | | | | | Efficiency (↓ mJ/token) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Sequence length | 500 | 1000 | 4000 | 8000 | 16000 | 500 | 1000 | 4000 | 8000 | 16000 |
| **Generate** **MMF (370M)** | **Loihi 2†** | **59.4** | **59.4** | **59.4** | **59.4** | **59.4** | **70.8** | **70.8** | **70.8** | **70.8** | **70.8** |
| MMF (370M) | Loihi 2* | 71.3 | 71.3 | 71.3 | 71.3 | 71.3 | 59 | 59 | 59 | 59 | 59 |
| MMF (370M) | H100 | 13.4 | 13.3 | <u>13.5</u> | 13.2 | <u>13.5</u> | 10.1k | 10.1k | 10.0k | 9.9k | <u>9.8k</u> |
| TF++ (370M) | H100 | 22.4 | <u>22.9</u> | 21.7 | 21.3 | 20.9 | <u>5.5k</u> | 5.6k | 6.2k | 6.8k | 8.2k |
| Alireo (400M) | Jetson | 14.3 | 14.9 | 14.7 | <u>15.2</u> | 12.8 | 723 | <u>719</u> | 853 | 812 | 1.2k |
| Qwen2 (500M) | Jetson | 13.4 | 14.0 | 14.1 | <u>15.4</u> | 12.6 | 791 | <u>785</u> | 912 | 839 | 1.2k |
| **Prefill** **Ours (370M)** | **Loihi 2†** | **11637** | **11637** | **11637** | **11637** | **11637** | **3.4** | **3.4** | **3.4** | **3.4** | **3.4** |
| MMF (370M) | Loihi 2* | 13965 | 13965 | 13965 | 13965 | 13965 | 2.8 | 2.8 | 2.8 | 2.8 | 2.8 |
| MMF (370M) | H100 | 11.4k | 13.1k | 30.6k | 51.6k | <u>84.6k</u> | 6.1 | 5.3 | 2.5 | 1.4 | <u>0.9</u> |
| TF++ (370M) | H100 | 21.6k | 32.7k | 44.3k | 55.4k | <u>60.5k</u> | 11.3 | 7.3 | 5.4 | 4.3 | <u>3.8</u> |
| Alireo (400M) | Jetson | 849 | 1620 | <u>3153</u> | 2258 | 1440 | 11.7 | 7.8 | <u>6.8</u> | 7.6 | 11.5 |
| Qwen2 (500M) | Jetson | 627 | 909 | 2639 | <u>3861</u> | 3617 | 17.9 | 13.9 | 6.7 | <u>4.4</u> | 5.3 |

\* The MatMul-free LM on Loihi 2 was characterized on an Oheo Gulch single-chip Loihi 2 system (N3C1 silicon) running NxKernel v0.2.0 and NxCore v2.5.8 (only accessible to Intel Neuromorphic Research Community members). The 1-chip case neglects inter-chip communication.

† Inter-chip communication causes a derived ≈20% slowdown over the single chip case (Appendix C.4.1).

‡ Transformer LMs characterized on NVIDIA Jetson Orin Nano 8GB using MAXN power mode running Jetpack 6.2, TensorRT 10.3.0, CUDA 12.4. Energy values include CPU_GPU_CV, SOC, and IO components as reported by jtop 4.3.0.

Performance results as of Jan 2025 and may not reflect all publicly available security updates. Results may vary.
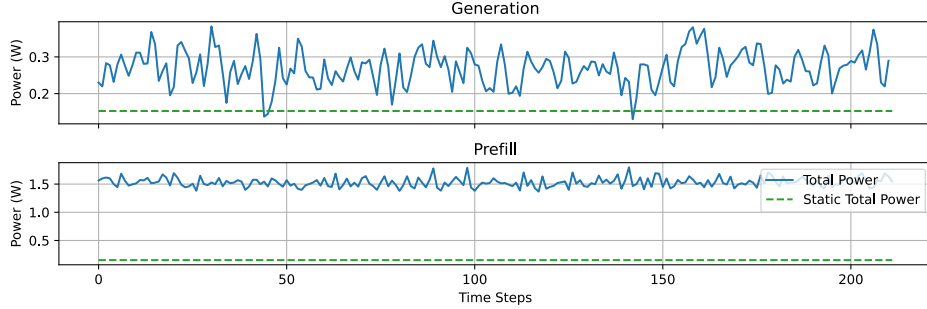


Figure 8: Power of one MatMul-free block on a single-chip Loihi 2 system.

but full-model power estimation uses $\hat{P}_{\text{generate}} = \tilde{P}^{\text{1-chip}} + 24 \times \bar{P}^{\text{1-chip}}$, as only one chip processes information at any timestep while others remain idle, drawing only static power. Energy per token estimates follow $\hat{E}_{\text{generate}} = \hat{P}_{\text{generate}} * T_{\text{ttft}}$.

**Multi-chip experiments** Single-chip estimates inherently exclude inter-chip communication overhead. To address this limitation, we deployed all 24 MatMul-free LM blocks on the Alia Point system (Figure 6), utilizing 32 of 128 available chips. Each block maps to one chip, occupying 24 of 32 chips. Throughput calculations follow single-chip methodology, while power and energy metrics are directly measured. We measured the time-per-step (TPS) with increasing chip count on the Alia Point Loihi 2 system. The TPS remains stable for systems using up to four chips, then increases by ≈ 20% and remains stable up to the full 24-chip system. Energy per token remains stable across all chip counts. Accordingly, inter-chip communication is derived to cause an approximately 20% slowdown over the single chip performance results (Table 3).

This scaling behavior indicates favorable performance for larger MatMul-free models on Loihi 2 systems.

**Embedding and un-embedding layers** Current experiments exclude embedding and un-embedding layers mapping between $\mathbb{R}^{1024}$ embedding space and vocabulary size $V = 32,000$. Our experiments mapped un-embedding to 7 chips, totaling 31 chips for the complete 370M model. The unembedding layer is partitioned and heavily parallelized such that throughput remains stable.

### C.4.2 Detailed NVIDIA Jetson Results

We evaluated state-of-the-art transformer language models on NVIDIA Jetson Orin Nano 8GB (MAXN power mode, Jetpack 6.2, TensorRT 10.3.0, CUDA 12.4). Power measurements included GPU, SOC, and I/O subsystems via integrated profiling tools. Two inference modes were evaluated: **prefill mode** (parallel input prompt processing) and **generate mode** (sequential autoregressive token generation). Generate mode's inherent token dependencies preclude sequence-level parallelization, yielding lower throughput than prefill mode. Performance metrics—throughput (tokens/second), average power (Watts), and energy per token (Joules)—were derived from 30-iteration averaged measurements using `jtop` utility.

Prefill mode achieved thousands of tokens/second through parallel token processing and effective pipelining. Generate mode demonstrated lower throughput due to sequential processing constraints. Power consumption varied between modes: prefill sustained higher continuous power draw, while generate exhibited variable profiles alternating between active processing and idle states.

Energy analysis revealed prefill mode's high throughput coincided with increased energy per token. Generate mode's extended latency elevated total energy per token despite lower instantaneous power draw.

These baseline measurements enable direct comparison with Loihi 2 neuromorphic implementations. Transformer models on Jetson Orin Nano achieve substantial throughput but at higher energy costs compared to MatMul-free LM estimates on Loihi 2, see Table 6.
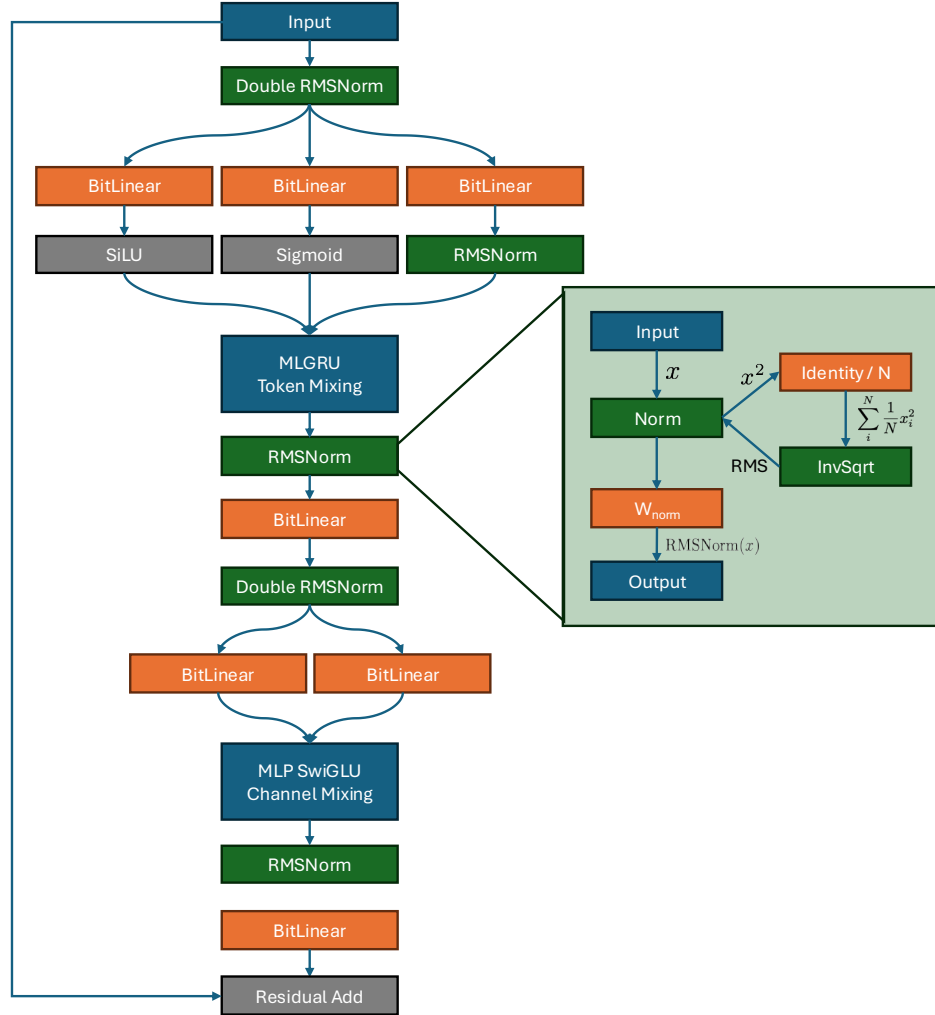
Figure 9: **Left**: Computational graph of a single MatMul-free LM layer, simplified from the actual computational graph that is mapped on the Loihi 2 chip. The RMSNorm is visualized as a single node. **Right**: Computational graph of the RMSNorm layer implemented on the Loihi 2 chip. See Sec. 6.5 for corresponding explanation.