

# FORGE: COMPILING A UNIFIED ABSTRACTION INTO SCALABLE KERNELS FOR LINEAR ATTENTION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The quadratic complexity of softmax attention poses a major bottleneck for long-context modeling, motivating a surge of linear attention variants with linear complexity. Unlike softmax attention, which benefits from optimized kernels, linear attention lacks general-purpose, hardware-efficient support and scalable distributed implementations. We introduce Forge, a domain-specific compiler that automates the generation of high-performance, scalable kernels for a wide range of linear attention models directly from high-level PyTorch code. At its core, Forge employs an intuitive programming abstraction that decomposes any linear attention algorithm into three canonical phases: intra-chunk computation, inter-chunk state propagation, and output merging. This unified abstraction enables Forge to perform domain-specific optimizations, automatically generating kernels that fuse computation and communication at a fine-grained tile level and eliminating host synchronization. Our evaluation demonstrates that Forge combines programmability with performance: a wide range of linear attention variants can be implemented in just a few dozen lines of code, while the generated kernels deliver 1.01x-4.9x the performance of state-of-the-art expert-optimized library and scale with near-linear efficiency on scalar gated linear attention to 16 million tokens on 128 GPUs, surpassing the state-of-the-art distributed baseline by up to 7.2x.

## 1 INTRODUCTION

Transformer models rely on self-attention, which has quadratic time and memory complexity with respect to sequence length. As models handle increasingly long contexts, this quadratic bottleneck severely limits scalability. In response, many efficient-attention mechanisms have been proposed. In particular, linear attention methods remove the softmax nonlinearity and reorder computations to achieve linear computation complexity and constant-memory inference. This has led to a proliferation of innovative architectures, such as Mamba (Gu & Dao, 2023; Dao & Gu, 2024), RetNet (Sun et al., 2023), RWKV (Peng et al., 2023), GLA (Yang et al., 2023), HGRN (Qin et al., b) and Gated DeltaNet (Yang et al., 2024a). These architectures demonstrate capabilities competitive with, or even superior to, standard transformers.

Unlike softmax attention, which has benefited from highly optimized and now-standardized kernels like Flash-Attention (Dao et al., 2022) and Ring-Attention (Liu et al., 2023) that efficiently map its computation and communication to modern AI infrastructure, the landscape for linear attention is far more fragmented. Flash-Linear-Attention (FLA) (Yang & Zhang, 2024) provide a valuable collection of triton kernels for a series of linear attention variants. But it essentially relies on expert developers to provide manual implementation for each variant. The rapid evolution of linear attention variants means that a one-size-fits-all solution does not exist. This forces researchers into a costly and inefficient cycle of manual kernel development for each new variant, a process fraught with two major challenges.

First, the implementation of high-performance kernels is an arduous task requiring deep hardware expertise. While many linear attention models share a conceptual similarity, their specific state update rules and memory access patterns can differ substantially. Achieving hardware efficiency necessitates not only fusing the state update rule into a single kernel but also manually tuning hardware-specific parameters like pipeline schedules and tile sizes. Even with high-level DSLs like



Figure 1: Comparison between PyTorch and our proposed DSL in writing scalar gated linear attention (data type conversion is omitted for simplicity). And the performance comparison of different approaches. Test shape: 32 heads with 128 head dimension. Torch-Eager fails to parallelize the computation, while Torch-Compile also performs poorly. FLA provides expert optimized triton kernels, while program generated by Forge offers comparable performance to handwritten kernels with the scalability to distributed environments.

Triton (Tillet & Cox, 2019), developers must often delve into low-level hardware details, such as managing barriers or dealing with shared memory capacity limitation, to extract maximum performance for each variant. This creates a high barrier to entry and slows down the pace of innovation.

Second, existing solutions lack robust support for distributed execution, which is non-negotiable for scaling to contexts of hundreds of thousands or millions of tokens. FLA offers a bunch of single-device kernels but do not address the distributed scaling problem. When sequence lengths exceed the memory capacity of a single accelerator, distributed sequence parallelism becomes essential. However, enabling sequence-parallel linear attention is non-trivial: it typically requires custom communication schedules tailored to each variant’s state update rule. Existing sequence parallel schemes, such as LASP and LASP-2 (Sun et al., 2024a; 2025) are designed for specific architectures and employ generic communication primitives (e.g., All-Gather from NCCL). The mismatch of existing communication primitive and dataflow of distributed linear attention leading to significant network bandwidth underutilization (Chou et al., 2025).

**Can we provide a solution to bridge the gap between the rapid evolution of linear attention algorithms and the difficulty of developing scalable kernels?** We observe that many of these difficulties stem from not exploiting the common structure underlying linear-attention variants. Our central insight is that most linear-attention variants share a small set of canonical operations and data exchanges. Based on this, we introduce Forge, a compiler-driven framework that allows implementing the majority of linear attention variants in a few lines of idiomatic PyTorch code and scaling them to distributed system. As shown in Figure 1, our DSL expresses linear attention in three modular functions. The compiler translates the DSL into high performance kernels: reducing latency from 34.6 seconds (PyTorch eager) to 9.2 ms, even better than SOTA hand-written kernel from FLA. More importantly, the latency was further reduced to 2.7 ms when scaled to 4 GPUs distributed system.

The frontend of Forge ingests a user-defined computation logic for a linear attention variant. Its backend then intelligently maps this logic, along with potential communication operations, onto hardware accelerators and network interfaces, applying domain-specific optimizations to generate high-performance, distributed-aware kernels. Forge is built around three key principles: **1 A Linear-Attention-Specific Programming Abstraction.** We formulate our programming abstraction based on the canonical chunk-parallel representation of linear attention. This model decomposes the computation into three intuitive phases: a compute phase that processes local chunks of the sequence in parallel, an update phase that communicates and updates the inter-chunk states, and a merge phase that combines the global state with local chunk results. This abstraction aligns directly with the mathematical structure of chunk-wise parallel form, allowing researchers to translate their algorithm’s formulation into our framework with minimal effort by providing simple PyTorch callable. **2 Native Compute-Communication Fusion.** At compilation time, Forge lowers the three user-provided callables into Triton code. We leverage Triton-Distributed (Zheng et al., 2025a)

as our compiler backend, which extends Triton with native communication primitives. This allows our compiler to generate fine-grained, tile-level communication instructions that are fused directly with computation. By creating custom communication patterns tailored to the algorithm, we bypass the overhead and limitations of standard libraries like NCCL, enabling a more efficient use of the underlying network fabric and dramatically improving hardware utilization. **③ Targeted Optimization of System Bottlenecks.** Beyond kernel fusion, we identify and optimize other critical system-level bottlenecks that affect real-world system performance. Forge employs a suite of techniques, including Ahead-of-Time (AOT) compilation with static kernel dispatcher built on top of Triton to reduce runtime overhead and an adaptive parallelism scheduler that dynamically explores the optimal configuration of compute resources, further boosting the end-to-end efficiency of linear attention execution.

To demonstrate its flexibility, we implement a broad range of linear attention variants using Forge, each requiring only dozens of lines of code. This programmability does not come at the cost of performance. On a single GPU, our generated kernels achieve 1.01x to 4.9x the performance of the state-of-the-art FLA library of expert-tuned kernels. Furthermore, in distributed settings, Forge demonstrates near-linear scalability on up to 128 GPUs, outperforming the leading open-source baseline by up to 7.2x.

## 2 PRELIMINARY

### 2.1 LINEAR ATTENTION ARCHITECTURE

Given a sequence  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L]^\top \in \mathbb{R}^{L \times d}$ , the input of attention block:  $\mathbf{q}_i, \mathbf{k}_i, \mathbf{v}_i = \mathbf{W}_q \mathbf{x}_i, \mathbf{W}_k \mathbf{x}_i, \mathbf{W}_v \mathbf{x}_i$  where  $\mathbf{x}_i, \mathbf{q}_i, \mathbf{k}_i, \mathbf{v}_i, \mathbf{y}_i \in \mathbb{R}^d$  and the weights  $\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v \in \mathbb{R}^{d \times d}$ . Transformers employ softmax attention as a token mixer (Vaswani et al., 2017):

$$\mathbf{o}_i = \sum_{j=1}^i \frac{\exp(\mathbf{q}_i^\top \mathbf{k}_j)}{\sum_{p=1}^i \exp(\mathbf{q}_i^\top \mathbf{k}_p)} \mathbf{v}_j \quad (1) \quad \mathbf{O} = \text{softmax}(\mathbf{Q}\mathbf{K}^\top \odot \mathbf{M})\mathbf{V} \quad (2)$$

Equation 2 is the matrix form of Equation 1 where  $\mathbf{Q} := [\mathbf{q}_1, \dots, \mathbf{q}_L]^\top$ ,  $\mathbf{K} := [\mathbf{k}_1, \dots, \mathbf{k}_L]^\top$ ,  $\mathbf{V} := [\mathbf{v}_1, \dots, \mathbf{v}_L]^\top \in \mathbb{R}^{L \times d}$  and  $\mathbf{M} \in \{-\infty, 1\}^{L \times L}$  is a causal mask.

Such matrix form is well suited to modern accelerators, which excel at large matrix multiplications, but it incurs  $\mathcal{O}(L^2 d)$  complexity. If we remove the softmax operation, the computation becomes associative:  $\mathbf{o}_i = \mathbf{q}_i (\mathbf{k}_i \mathbf{v}_i^\top)$  which reduces the complexity to  $\mathcal{O}(Ld^2)$ . The recurrence form is expressed as:

$$\mathbf{S}_t = \mathbf{S}_{t-1} + \mathbf{k}_t \mathbf{v}_t^\top, \quad \mathbf{o}_t = \mathbf{q}_t \mathbf{S}_t. \quad (3)$$

Here  $\mathbf{S} \in \mathbb{R}^{d \times d}$  is the state (or memory) updated in each time step. Equation 3 highlights the key idea of linear attention: replacing the exponential kernel in softmax attention with a linear recurrence. Although this formulation resembles RNNs (Hochreiter & Schmidhuber, 1997). The critical difference is that dependencies across time steps remain *linear*, which makes parallel training possible. Indeed, linear attention can be written in fully parallel form:

$$\mathbf{O} = (\mathbf{Q}\mathbf{K}^\top \odot \mathbf{M})\mathbf{V} \quad (4)$$

Modern linear attention often augments the recurrence with a decay or gating mechanism, e.g.,  $\mathbf{S}_t = \mathbf{G}_t \odot \mathbf{S}_{t-1} + \mathbf{k}_t \mathbf{v}_t^\top$  (Gu & Dao, 2023), or with more sophisticated update rules such as the delta rule (Yang et al., 2024b):  $\mathbf{S}_t = \mathbf{S}_{t-1}(\mathbf{I} - \beta_t \mathbf{k}_t \mathbf{k}_t^\top) + \beta_t \mathbf{v}_t \mathbf{k}_t^\top$ , which enhance memory utilization.

### 2.2 CHUNK-WISE PARALLEL FORM OF LINEAR ATTENTIONS

The fully parallel form in Equation 4 achieves maximum hardware utilization but retains quadratic complexity. Conversely, the recurrent form in Equation 3 has linear complexity but is inherently sequential and hardware-inefficient. In practice, linear attention strikes a balance by adopting chunk-wise parallelization (Hua et al., 2022; Sun et al., 2023; Yang et al., 2023).

Specifically, the sequence of length  $L$  is partitioned into  $\frac{L}{C}$  chunks of size  $C$ . Let  $\mathbf{Q}_{[i]}, \mathbf{K}_{[i]}, \mathbf{V}_{[i]} \in \mathbb{R}^{C \times d}$  denote the query, key, and value matrices of the  $i$ -th chunk, and let  $\mathbf{S}_{[i]} \in \mathbb{R}^{d \times d}$  be the state after processing chunk  $i$ . The chunk-wise formulation separates computation into intra-chunk and inter-chunk two parts, and then merge both together to get the final output as shown in Figure 2.

$$\mathbf{S}_{[i]} = \mathbf{S}_{[i-1]} + \sum_{j=(i-1)C}^{iC} \mathbf{k}_j^\top \mathbf{v}_j \quad (5)$$

Matrix:  $\mathbf{K}_{[i]}^\top \mathbf{V}_{[i]}$

$$\mathbf{O}_{[i]} = \underbrace{\mathbf{Q}_{[i]} \mathbf{S}_{[i-1]}}_{\text{inter}} + \underbrace{(\mathbf{Q}_{[i]} \mathbf{K}_{[i]}^\top \odot \mathbf{M}) \mathbf{V}_{[i]}}_{\text{intra}} \quad (6)$$

Figure 2: Chunk-wise parallel form demonstration for linear attention.

The inter chunk item can be viewed as readout memory from start of the sequence to the start of current chunk while the intra chunk item is processing the information in current chunk. When chunk size  $C$  is set to sequence length  $L$ , it becomes the fully parallel form as in Equation 4. The chunk size is a tradeoff between parallelism and FLOPs.

### 2.3 DSLs AND DOMAIN SPECIFIC COMPILER

A domain-specific language (DSL) trades generality for performance by restricting program expressivity to a particular domain, creating opportunities for targeted optimization. Deep learning compilers like TVM (Chen et al., 2018), ThunderKitten (Spector et al., 2024), TileLang (Wang et al., 2025) and `torch.compile` (Ansel et al., 2024) exemplify this approach: by operating on a constrained set of primitives (i.e. `tir` in TVM and aten operators in PyTorch), they can systematically explore a focused design space to apply transformations like operator fusion and loop tiling, thereby automatically generating high-performance code from high-level descriptions.

This paradigm is particularly effective for operations that possess rich computational structure but high implementation complexity. By exposing domain-relevant abstractions, a DSL allows developers to specify *what* to compute, while delegating the complex details of *how* to execute it efficiently to the compiler. This separation of concerns is key to enabling automated, domain-specific optimizations without burdening the user with low-level hardware details.

Applied to the domain of linear attention, an effective DSL must therefore provide abstractions that are expressive enough to capture the diverse patterns found in various state update rules and parallel scan formulations. Simultaneously, its compiler must be able to recognize these common patterns and systematically generate optimized kernels for them, bridging the gap between high-level algorithmic design and performance, hardware-aware code. Because chunked parallel forms are complex to implement and offer opportunities to fully exploit hardware, this work focuses on the prefill phase of linear attention (used in inference or the forward pass of training). The backward pass can be implemented in a similar manner (Qin et al., a).

## 3 FORGE

In this section, we propose a unified abstraction of diverse linear attention variants. This abstraction enables programmers to easily express linear attention semantics without worrying about implementation details and kernel performance.

### 3.1 PROGRAMMING ABSTRACTION

Despite the numerous linear attention variants designed by researchers, we unify these variants into three commonly shared phases based on chunk-wise parallel form introduced in subsection 2.2. **1** Intra-Chunk Computation. The first phase computes a local state within each chunk of the input

Table 1: A comparison of representative linear attention variants that can be easily mapped to our three-phase abstraction, including HGRN (Qin et al., 2023), RetNet (Sun et al., 2023), Mamba2 (Dao & Gu, 2024), GLA (Yang et al., 2023), and GDN (Yang et al., 2024a). Despite diverse state types (vector vs. matrix) and decay mechanisms (element-wise product vs. matrix multiplication).  $v_t, k_t, q_t$  are value, key and query projections;  $\alpha_t, \beta_t, r_t, i_t$  are gates;  $\odot$  is the Hadamard product.

Model	Update rule	Read-out	State + Decay type
HGRN	$h_t = \alpha_t \odot h_{t-1} + (1 - \alpha_t) \odot v_t$	$o_t = h_t \odot q_t$	vector + data-dependent vector
RetNet	$S_t = \gamma S_{t-1} + v_t k_t^\top$	$o_t = S_t q_t$	matrix + data-independent scalar
Mamba2	$S_t = \gamma_t S_{t-1} + v_t k_t^\top$	$o_t = S_t q_t$	matrix + data-dependent scalar
GLA	$S_t = S_{t-1} \odot (1\alpha_t^\top) + v_t k_t^\top$	$o_t = S_t q_t$	matrix + data-dependent vector
GDN	$S_t = \alpha_t S_{t-1} (I - \beta_t k_t k_t^\top) + \beta_t v_t k_t^\top$	$o_t = S_t q_t$	matrix + data-dependent matrix

sequence. In this stage, computation across different chunks is embarrassingly parallel, as there are no data dependencies between them. Each chunk is processed independently, transforming its sequence of inputs into a relative state summary. **② Inter-Chunk State Propagation.** The second phase addresses the dependencies between chunks. To compute the correct global state at the beginning of each chunk, the state summaries from all preceding chunks must be accumulated. For instance, in the case of vanilla linear attention, this propagation corresponds to a prefix sum (scan) operation, as shown in Equation 5. This phase is inherently sequential due to the temporal dependencies between chunk states. And cross-device communication happen in this phase. **③ Merging and Output Generation.** The final phase merges the results of the intra-chunk and inter-chunk computations. Here, operations are once again parallel across chunks. Each chunk utilizes the global state propagated from Phase 2 and its local inputs to compute its final output sequence.

Based on these insights, we designed our programming abstraction around three corresponding callable: `chunk_mode`, `decay_mode` and `merge_mode` correspond to three phases. This abstraction empowers users to implement a new linear attention variant by simply defining its chunk-wise parallel logic in idiomatic PyTorch code, decoupling algorithmic expression from system optimization. Furthermore, this three-phase decomposition also helps us optimize the program: we separate the parts of the entire program that can be executed in parallel and the parts that must be executed serially and may involve cross-device communication. With this information, Forge can perform more aggressive and accurate optimizations.

As shown in Table 1, prominent linear attention variants employ vastly different state representations and decay mechanisms. Nonetheless, all can be expressed using the chunk-wise parallel formulation. For example, the chunk-wise parallel form for Mamba2 can be specified as:

$$S_{[t]} = (\prod \alpha) \odot S_{[t-1]} + V K^\top \quad O = Q^\top K \odot M \odot G V^\top + S Q \quad (7)$$

Under our framework, a user only needs to implement these equations within our three-phase programming abstraction in native PyTorch code. Forge then automatically handles all subsequent code generation, performance tuning, and scaling to distributed systems.

### 3.2 COMPILATION AND CODE GENERATION

Given the chunk-wise parallel form description in our DSL, Forge performs a series of graph-level and system-level transformations to produce optimized program as illustrated in Figure 3. The user-defined function is first captured as Torch.fx graph (Reed et al., 2022) via tracing. We choose fx graph as our intermediate representation (IR) because it is the most powerful tool in the PyTorch ecosystem, most torch operators can be directly captured, which makes our DSL expressive enough to implement most of linear attention variants. We first replace special op code in fx graph as custom instructions (e.g. `placeholder` is replaced with `load` instruction). Domain-specific optimization passes are then applied. To enable fusion of non-trivial operators (e.g. `lower_triangular_inverse` in GDN), we provide a bunch of custom triton kernels commonly used in linear attention domain and mapping them to corresponding torch operators. Once the compiler see these operators, for instance, `torch.inverse`, they will be marked and subsequently substituted with custom triton source code in code generation phase. Common optimization passes like transpose elimination are also applied in this phase. The IR is further rewritten with system-resource awareness. Forge take the hardware information to generate hardware-specific in-

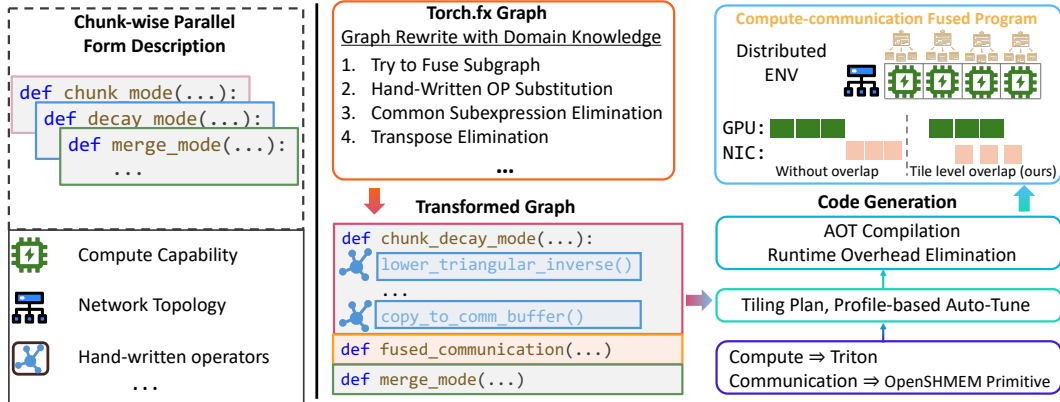


Figure 3: An overview of the Forge compilation pipeline. Our compiler ingests a high-level description of a linear attention algorithm, specified using our DSL (`chunk_mode`, `decay_mode`, `merge_mode`). This description is traced into a graph structure, which then undergoes a series of domain-specific optimization passes. The optimized graph is then compiled to Triton with native communication support, enabling fine-grained, tile-level overlap between computation and communication for better utilization of GPU and Network Interface Card (NIC).

structions. For instance, Tensor Memory Accelerator (TMA) availability is marked as an attribute of load instruction.

Recent studies have demonstrated that fine-grained compute–communication fusion can more effectively hide latency in distributed settings (Chang et al., 2024; Zheng et al., 2025b). We adopt this technique in Forge by automatically fusing computation and communication at the tile level, thereby reducing data dependency scope and eliminating the frequent GPU–host synchronizations inherent to traditional overlap strategies (Jangda et al., 2022). Within our programming abstraction, all cross-device communication is confined to the second phase, i.e., *inter-chunk state propagation*. Consequently, Forge first analyzes the data dependencies in this phase, then determines the corresponding computational tiling and communication tiling strategies based on the network topology, selects the appropriate communication mode, and generates computation and on-device communication instructions.

Finally, the IR is lowered into triton source code. Different with torch inductor (Ansel et al., 2024) targeting on official triton, Forge targets on Triton-Distributed (Zheng et al., 2025a), who additionally provide fine-grained communication control. On device computation is translated into computation primitives provided by Triton while the communication logic is mapped to the OpenShmem-style communication primitives provided exclusively by Triton-distributed, which are ultimately translated into GPU-initiated communication operations.

### 3.3 PERFORMANCE OPTIMIZATIONS

With domain-specific knowledge of linear attention, Forge can explore a compact yet effective optimization space. In particular, the choice of whether to fuse different phases introduces an important trade-off. For example, fusing `chunk_mode` with `decay_mode` avoids materializing intermediate states in global memory, thereby reducing memory traffic. However, such fusion also limits available parallelism, since computations can no longer be scheduled independently at the chunk level. There are many such trade-offs and Forge will handle all of these to get better performance. Forge employs a parallelism scheduling algorithm that dynamically chooses an optimal parallelization strategy based on input shapes and hardware information internally. The specifics of this parallelism scheduler are detailed in Appendix B.

Beyond the optimization for target program, we further incorporate a set of optimizations tailored to the compilation system itself. In practice, the Triton runtime introduces overheads on the order of hundreds of microseconds, often exceeding the actual most of linear attention kernel execution time at short to medium sequence length (e.g. 2K or even 4K). While approaches like CUDA-Graph can mitigate launch overhead for workloads with static input shapes, they often incur significant memory costs and are unsuitable for the dynamic workloads in inference scenarios. To address these limitations natively, we extend Triton compiler with a custom Ahead-of-Time (AOT) compilation



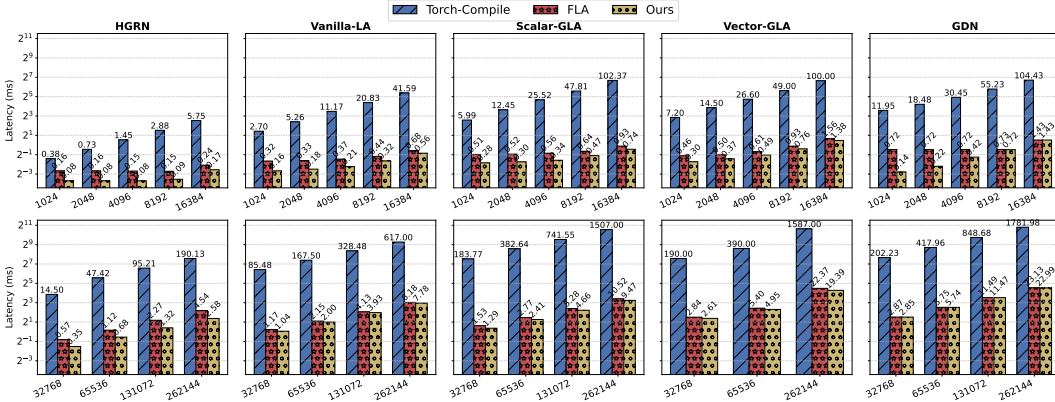


Figure 4: Latency comparison of different linear attention variants under varying sequence lengths on a single H100 GPU. Each subplot corresponds to one model, with the top row showing short to medium sequences (1K–16K) and the bottom row showing long sequences (32K–256K).

module. Specifically, our module compiles Triton source code into pre-linked dynamic libraries ahead of execution. At runtime, Forge employs a profile-guided static dispatcher that bypasses the Triton runtime entirely, invoking the optimal pre-compiled binary directly through the CUDA Driver API. The static dispatcher is automatically generated by Forge from an offline performance database, ensuring that the empirically best-performing kernel is selected for any given workload without incurring runtime overhead from hash lookups or dynamic compilation logic.

Another problem is redundant compilation. Since our system is specialized for linear attention, we observe that certain input tensor dimensions (e.g., head dimension and number of heads) remain relatively static across runs, while sequence length is typically dynamic. This property facilitates efficient AOT compilation: Forge allows users to specify constant dimensions and their admissible ranges via input metadata. Forge enumerates the Cartesian product of these ranges and generates all potentially required kernels in advance. Furthermore, by leveraging PyTorch’s symbolic tracing, our system supports tensors with symbolic shapes, ensuring that recompilation is unnecessary unless static dim change.

## 4 EXPERIMENTS

We implement several high-performance linear attention kernels using Forge, including HGRN, vanilla linear attention, scalar GLA, vector GLA, and Gated DeltaNet (Qin et al., 2023; Dao & Gu, 2024; Yang et al., 2023; 2024a). While many linear attention models differ in their parameterization, the computational patterns we implement cover over ten existing model designs (Appendix C).

### 4.1 SINGLE-DEVICE EVALUATION

Figure 4 reports the latency of kernels generated by Forge across sequence lengths ranging from 1K to 256K on a single H100 GPU. We compare with two baselines: Torch-Compile, representing a general-purpose compiler without domain-specific knowledge, and FLA (commit hash: 02766e71), the state-of-the-art library of providing expert-tuned Triton kernels for linear attention. For all variants except HGRN, we fix BatchSize = 1, NumHeads = 32, and HeadDim = 128, and vary sequence length. For HGRN, we follow its original single-head configuration. A batch size of one is a standard and reasonable choice, as a single long sequence in linear attention is computationally equivalent to a batch of packed shorter sequences.

Torch-Compile consistently exhibits poor performance, as it fails to apply the domain-specific fusion strategies required for linear attention. FLA achieves strong performance by carefully hand-tuning IO-aware tiling. Our work, Forge, automating this optimization process, consistently matches or outperforms FLA. On HGRN, Forge achieves  $1.64\text{--}2.02\times$  speedup, highlighting its ability to discover optimization opportunities beyond expert tuning (We check the generated code and find Forge use a more efficient threads allocation). On scalar and vector GLA as well as vanilla linear attention, Forge provides stable improvements ( $1.1\text{--}2.0\times$ ) while for GDN, performance converges to FLA at longer sequences. The speedup is most pronounced on short to medium sequence length. In this

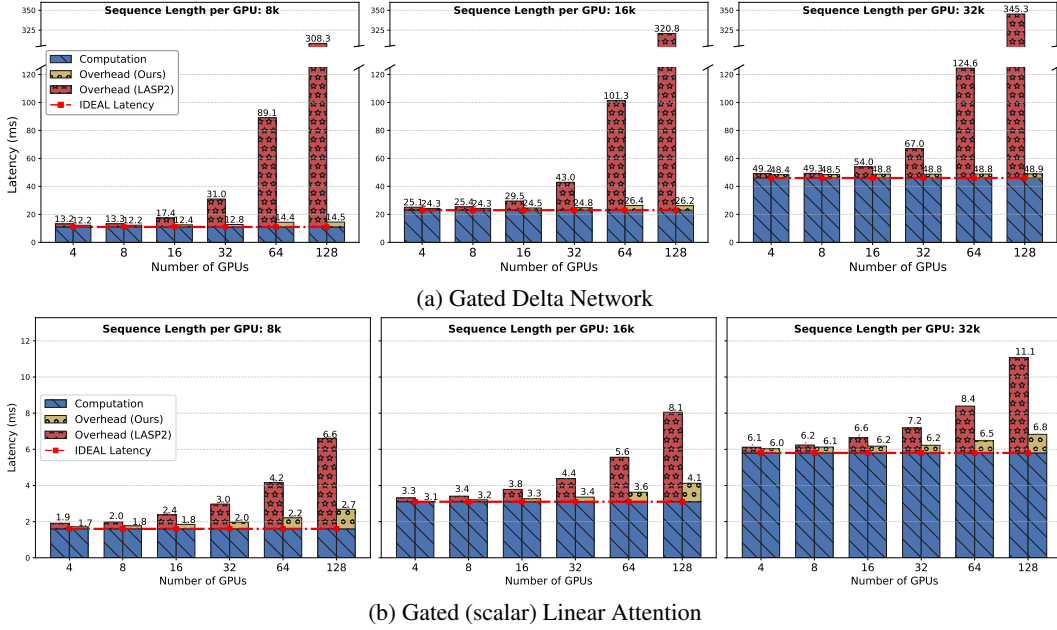


Figure 5: Weak scaling performance comparison for GDN and GLA models. Both figures show latency breakdown across a varying number of GPUs with a fixed sequence length per GPU.

regime, system-level overheads and the choice of parallelism strategy are the dominant factors, and Forge effectively eliminate these bottlenecks, as discussed in subsection 3.3.

## 4.2 SEQUENCE PARALLEL EVALUATION

We further evaluate weak scaling behavior under distributed training, comparing against LASP2 (Sun et al., 2025), the strongest open-source baseline at the time of writing. ZeCO (Chou et al., 2025) reports improved scaling via pipelined communication, but its implementation is not public. We benchmark two representative workloads: GDN and scalar GLA, which feature matrix-multiplication and element-wise decay mechanisms, respectively. Experiments were conducted on a cluster of up to 128 NVIDIA H20 GPUs, interconnected with NVSwitch within nodes and InfiniBand between nodes. We fix BatchSize = 4, NumHeads = 32, HeadDim = 128, and maintain a constant workload per GPU while scaling the total sequence length from 128K (on 4 GPUs) to 4 million tokens (on 128 GPUs). The ideal outcome for weak scaling is a constant execution time.

Figure 5 show that Forge exhibits near-ideal weak scaling for both workloads: latency remains flat as GPU count increases. This is attributable to two core features of our compiler. First, it generates communication patterns that avoid the data redundancy incurred by the All-Gather primitive used in LASP2. Second, its ability to fuse computation and communication effectively hides the latency of the local state update. This is particularly impactful for GDN, whose matrix-based update is time consuming. In contrast, the communication and computation redundancy of LASP2 is amplified as the number of nodes increases, causing its performance to degrade significantly (e.g., from 49.2ms on 4 GPUs to 345ms on 128 GPUs for GDN). This confirms that Forge eliminates redundant communication and achieves scalable performance on large GPU clusters.

## 4.3 ABLATION STUDY

**AOT compilation with static dispatcher.** We measure the end-to-end latency of the execution of Scalar GLA kernel including both kernel execution and to demonstrate Triton runtime overhead and the efficiency of our solution as shown in Figure 6. At a sequence length of 1024, this overhead (207 $\mu$ s) is over 4.4 times the actual kernel execution time (47 $\mu$ s). Our static dispatcher mitigates this issue by reducing the overhead by 46%, yielding a 1.6x end-to-end speedup on these latency-sensitive inputs. As the sequence length grows and execution time be-

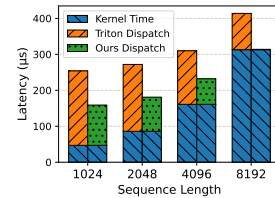


Figure 6: Execution time decomposition.



comes the dominant factor, our dispatcher consistently maintains a negligible overhead that is effectively eliminated at 8192 tokens (reducing from 101 $\mu$ s to just 1 $\mu$ s).

**Tile Level Compute-Communication Overlapping.** We conduct a targeted study to isolate and quantify the benefits of our tile-level compute-communication fusion. We set two baselines for the inter-rank state propagation: ① Serial communication, where each rank  $i$  waits to receive all state data from rank  $i - 1$  before performing its local update and sending to rank  $i + 1$ . ② Pipelined baseline, which chunks the state and uses standard NCCL send/recv operations to overlap the computation of one chunk with the communication of the next. Our method, in contrast, generates a single kernel that fuses the state update computation with communication primitives at the tile level. The results, presented in Table 2 show the standard Pipelined (PyTorch) baseline is slightly slower than the naive Serial implementation, demonstrating that host-managed pipelining can be counterproductive due to the significant overhead of launching numerous small operations and the required host-device synchronization. Conversely, our fused kernel substantially reduces the total time, achieving a 1.56x speedup over the serial baseline.

Table 2: State communication time on 8xH800 GPUs. State size is 67MB.

Method	Time (us)
Serial	873
Torch-Pipeline	902
Ours	560

## 5 RELATED WORK

There are a wide range of approaches to address the quadratic complexity of softmax attention. Sparse attention mechanisms (Zaheer et al., 2020; Xiao et al., 2023; Yuan et al., 2025) leverage structured or un-structured sparsity in attention to skip computation. Quantized attention (Shah et al., 2024; Zhang et al., 2024) use exploit low-precision arithmetic unit in modern hardware to get higher throughput. There are also various techniques to reduce key-value (KV) cache overhead. KIVI (Liu et al., 2024b) and SKVQ (Duanmu et al., 2024) directly compress KV cache using quantization while grouped-query attention (GQA) (Ainslie et al., 2023), and multi-head latent attention (MLA) (Liu et al., 2024a) alter the attention architecture to reduce memory overhead.

Another line of work proposes architectural alternatives with lower complexity, including linear attention variants (Katharopoulos et al., 2020; Dao & Gu, 2024; Peng et al., 2023; Yang et al., 2024a) as well as test-time-training approaches (Sun et al., 2024b; Behrouz et al., 2024). For linear attention sequence parallelism, LASP (Sun et al., 2024a) first extend linear attention to distributed environments with serial send-receive primitive. LASP2 (Sun et al., 2025) improves on this by leveraging collective communication primitives, yet both approaches still incur significant bandwidth under-utilization. ZeCO (Chou et al., 2025) introduces a pipelined send-receive scheme to hide send/receive latency but relies on manual chunk-size tuning and does not detail its implementation. Therefore, it was not included in our comparison.

AI compilers have been developed to optimize a broad range of workloads. `torch.compile` (Ansel et al., 2024), TVM (Chen et al., 2018), and TASO (Jia et al., 2019) are effective for common operators but their optimization spaces do not cover linear attention. Operator-level compilers such as Triton (Tillet & Cox, 2019), ThunderKitten Spector et al. (2024), TileLang (Wang et al., 2025), and Triton-Distributed (Zheng et al., 2025a) provide expressive abstractions for modern accelerators, but supporting the large and growing family of linear attention variants still requires substantial manual development effort. The most closely related work, FlexAttention (Dong et al., 2024), targets block-sparse softmax attention and is limited to single-device settings.

## 6 CONCLUSION

We presented Forge, a domain-specific compiler for linear attention that unifies diverse algorithmic variants under a common three-phase abstraction. By generating hardware-efficient kernels with native distributed execution support, Forge bridges the gap between rapidly evolving linear attention research and the complexity of hand-tuned implementations. Our evaluation demonstrates that Forge

achieves both high performance and broad applicability across modern linear attention models. We hope this work will accelerate the development of new architectures and inspire further research at the intersection of deep learning algorithms and domain-specific compilation.

## 7 REPRODUCIBILITY STATEMENT

We are committed to the reproducibility of our work. Most of algorithms, implementation details, and experimental setups are described in the paper. Due to organizational policies requiring an internal review prior to public release, the source code is not included with the submission. However, we are committed to open-sourcing the code and will provide it to reviewers upon request for the purpose of evaluation.

## REFERENCES

- Joshua Ainslie, James Lee-Thorp, Michiel De Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. Gqa: Training generalized multi-query transformer models from multi-head check-points. *arXiv preprint arXiv:2305.13245*, 2023.
- Jason Ansel, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznesensky, Bin Bao, Peter Bell, David Berard, Evgeni Burovski, et al. Pytorch 2: Faster machine learning through dynamic python bytecode transformation and graph compilation. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*, pp. 929–947, 2024.
- Ali Behrouz, Peilin Zhong, and Vahab Mirrokni. Titans: Learning to memorize at test time. *arXiv preprint arXiv:2501.00663*, 2024.
- Li-Wen Chang, Wenlei Bao, Qi Hou, Chengquan Jiang, Ningxin Zheng, Yinmin Zhong, Xuanrun Zhang, Zuquan Song, Chengji Yao, Ziheng Jiang, et al. Flux: Fast software-based communication overlap on gpus through kernel fusion. *arXiv preprint arXiv:2406.06858*, 2024.
- Feiyang Chen, Yu Cheng, Lei Wang, Yuqing Xia, Ziming Miao, Lingxiao Ma, Fan Yang, Jilong Xue, Zhi Yang, Mao Yang, et al. Attentionengine: A versatile framework for efficient attention mechanisms on diverse hardware platforms. *arXiv preprint arXiv:2502.15349*, 2025.
- Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Meghan Cowan, Leyuan Wang, Yuwei Hu, Luis Ceze, et al. {TVM}: An automated {End-to-End} optimizing compiler for deep learning. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, pp. 578–594, 2018.
- Yuhong Chou, Zehao Liu, Ruijie Zhu, Xinyi Wan, Tianjian Li, Congying Chu, Qian Liu, Jibin Wu, and Zejun Ma. Zeco: Zero communication overhead sequence parallelism for linear attention. *arXiv preprint arXiv:2507.01004*, 2025.
- Tri Dao and Albert Gu. Transformers are SSMS: Generalized models and efficient algorithms through structured state space duality. In *International Conference on Machine Learning (ICML)*, 2024.
- 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.
- 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.
- Juechu Dong, Boyuan Feng, Driss Guessous, Yanbo Liang, and Horace He. Flex attention: A programming model for generating optimized attention kernels. *arXiv preprint arXiv:2412.05496*, 2024.
- Haojie Duanmu, Zhihang Yuan, Xiuhong Li, Jiangfei Duan, Xingcheng Zhang, and Dahua Lin. Skvq: Sliding-window key and value cache quantization for large language models. *arXiv preprint arXiv:2405.06219*, 2024.
- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.

- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.
- William Hu, Drew Wadsworth, Sean Siddens, Stanley Winata, Daniel Y. Fu, Ryann Swann, Muhammad Osama, Christopher Ré, and Simran Arora. Hipkittens: Fast and furious amd kernels, 2025. URL <https://arxiv.org/abs/2511.08083>.
- Weizhe Hua, Zihang Dai, Hanxiao Liu, and Quoc Le. Transformer quality in linear time. In *International conference on machine learning*, pp. 9099–9117. PMLR, 2022.
- Abhinav Jangda, Jun Huang, Guodong Liu, Amir Hossein Nodehi Sabet, Saeed Maleki, Youshan Miao, Madanlal Musuvathi, Todd Mytkowicz, and Olli Saarikivi. Breaking the computation and communication abstraction barrier in distributed machine learning workloads. In *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, pp. 402–416, 2022.
- Zhihao Jia, Oded Padon, James Thomas, Todd Warszawski, Matei Zaharia, and Alex Aiken. Taso: optimizing deep learning computation with automatic generation of graph substitutions. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*, pp. 47–62, 2019.
- Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International conference on machine learning*, pp. 5156–5165. PMLR, 2020.
- Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Daya Guo, et al. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *arXiv preprint arXiv:2405.04434*, 2024a.
- Hao Liu, Matei Zaharia, and Pieter Abbeel. Ring attention with blockwise transformers for near-infinite context. *arXiv preprint arXiv:2310.01889*, 2023.
- Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhuo Xu, Vladimir Braverman, Beidi Chen, and Xia Hu. Kivi: A tuning-free asymmetric 2bit quantization for kv cache. *arXiv preprint arXiv:2402.02750*, 2024b.
- Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, et al. RwkV: Reinventing rnns for the transformer era. *arXiv preprint arXiv:2305.13048*, 2023.
- Zhen Qin, Xuyang Shen, Dong Li, and Yiran Zhong. Accelerating linear attention design by unifying forward & backward propagation. In *ES-FoMo III: 3rd Workshop on Efficient Systems for Foundation Models*, a.
- Zhen Qin, Songlin Yang, Weixuan Sun, Xuyang Shen, Dong Li, Weigao Sun, and Yiran Zhong. Hgrn2: Gated linear rnns with state expansion. In *First Conference on Language Modeling*, b.
- Zhen Qin, Songlin Yang, and Yiran Zhong. Hierarchically gated recurrent neural network for sequence modeling. *Advances in Neural Information Processing Systems*, 36:33202–33221, 2023.
- Zhen Qin, Weigao Sun, Dong Li, Xuyang Shen, Weixuan Sun, and Yiran Zhong. Lightning attention-2: A free lunch for handling unlimited sequence lengths in large language models. *arXiv preprint arXiv:2401.04658*, 2024.
- James Reed, Zachary DeVito, Horace He, Ansley Ussery, and Jason Ansel. torch.fx: Practical program capture and transformation for deep learning in python. *Proceedings of Machine Learning and Systems*, 4:638–651, 2022.
- Jay Shah, Ganesh Bikshandi, Ying Zhang, Vijay Thakkar, Pradeep Ramani, and Tri Dao. Flashattention-3: Fast and accurate attention with asynchrony and low-precision. *Advances in Neural Information Processing Systems*, 37:68658–68685, 2024.
- Julien Siems, Timur Carstensen, Arber Zela, Frank Hutter, Massimiliano Pontil, and Riccardo Grazi. Deltaproduct: Improving state-tracking in linear rnns via householder products. *arXiv preprint arXiv:2502.10297*, 2025.

- Benjamin F Spector, Simran Arora, Aaryan Singhal, Daniel Y Fu, and Christopher Ré. Thunderkittens: Simple, fast, and adorable ai kernels. *arXiv preprint arXiv:2410.20399*, 2024.
- Weigao Sun, Zhen Qin, Dong Li, Xuyang Shen, Yu Qiao, and Yiran Zhong. Linear attention sequence parallelism. *arXiv preprint arXiv:2404.02882*, 2024a.
- Weigao Sun, Disen Lan, Yiran Zhong, Xiaoye Qu, and Yu Cheng. Lasp-2: Rethinking sequence parallelism for linear attention and its hybrid. *arXiv preprint arXiv:2502.07563*, 2025.
- Yu Sun, Xinhao Li, Karan Dalal, Jiarui Xu, Arjun Vikram, Genghan Zhang, Yann Dubois, Xinlei Chen, Xiaolong Wang, Sanmi Koyejo, et al. Learning to (learn at test time): Rnns with expressive hidden states. *arXiv preprint arXiv:2407.04620*, 2024b.
- Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and Furu Wei. Retentive network: A successor to transformer for large language models. *arXiv preprint arXiv:2307.08621*, 2023.
- Philippe Tillet and David Cox. Triton: an intermediate language and compiler for tiled neural network computations. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*, pp. 10–19, 2019.
- 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.
- Lei Wang, Yu Cheng, Yining Shi, Zhengju Tang, Zhiwen Mo, Wenhao Xie, Lingxiao Ma, Yuqing Xia, Jilong Xue, Fan Yang, et al. Tilelang: A composable tiled programming model for ai systems. *arXiv preprint arXiv:2504.17577*, 2025.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. *arXiv preprint arXiv:2309.17453*, 2023.
- Songlin Yang and Yu Zhang. Fla: A triton-based library for hardware-efficient implementations of linear attention mechanism, January 2024. URL <https://github.com/fla-org/flash-linear-attention>.
- Songlin Yang, Bailin Wang, Yikang Shen, Rameswar Panda, and Yoon Kim. Gated linear attention transformers with hardware-efficient training. *arXiv preprint arXiv:2312.06635*, 2023.
- Songlin Yang, Jan Kautz, and Ali Hatamizadeh. Gated delta networks: Improving mamba2 with delta rule. *arXiv preprint arXiv:2412.06464*, 2024a.
- Songlin Yang, Bailin Wang, Yu Zhang, Yikang Shen, and Yoon Kim. Parallelizing linear transformers with the delta rule over sequence length. *Advances in neural information processing systems*, 37:115491–115522, 2024b.
- Jingyang Yuan, Huazuo Gao, Damai Dai, Junyu Luo, Liang Zhao, Zhengyan Zhang, Zhenda Xie, YX Wei, Lean Wang, Zhiping Xiao, et al. Native sparse attention: Hardware-aligned and natively trainable sparse attention. *arXiv preprint arXiv:2502.11089*, 2025.
- Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for longer sequences. *Advances in neural information processing systems*, 33:17283–17297, 2020.
- Jintao Zhang, Jia Wei, Haofeng Huang, Pengl Zhang, Jun Zhu, and Jianfei Chen. Sageattention: Accurate 8-bit attention for plug-and-play inference acceleration. *arXiv preprint arXiv:2410.02367*, 2024.
- Size Zheng, Wenlei Bao, Qi Hou, Xuegui Zheng, Jin Fang, Chenhui Huang, Tianqi Li, Haojie Duanmu, Renze Chen, Ruifan Xu, et al. Triton-distributed: Programming overlapping kernels on distributed ai systems with the triton compiler. *arXiv preprint arXiv:2504.19442*, 2025a.
- Size Zheng, Jin Fang, Xuegui Zheng, Qi Hou, Wenlei Bao, Ningxin Zheng, Ziheng Jiang, Dongyang Wang, Jianxi Ye, Haibin Lin, et al. Tilelink: Generating efficient compute-communication overlapping kernels using tile-centric primitives. *arXiv preprint arXiv:2503.20313*, 2025b.

## A USAGE OF LLMs

For the preparation of this manuscript, we utilized LLM, as writing assistant to enhance the quality of the prose. Our process was interactive: the authors provided initial drafts and specific sentences to the LLM and used its suggestions to refine the text. Majority of prompts aimed at refining sentences to improve conciseness, in addition to correcting grammar and improving overall clarity. The core scientific ideas, methodology, and experimental results were developed exclusively by the human authors, who bear full responsibility for all claims and content within this paper.

## B PARALLISIM SCHEDULER

Forge automatically explore different parallel schemes for linear attention during compilation. With the domain knowledge of linear attention, we can shrink the search space into two mainly used schemes, which is partially explored in FLA (Yang & Zhang, 2024).

In this section, we provide a theoretical analysis of two parallel execution strategies for linear attention using vanilla linear attention on GPU as an example. Then we derive a heuristic scheduling algorithm 1 that Forge used to build a parallelism scheduler, selecting the optimal strategy based on the input tensor shapes and hardware characteristics.

### B.1 NOTATION

The vanilla linear attention is defined by the recurrence:

$$\mathbf{S}_t = \mathbf{S}_{t-1} + \mathbf{k}_t \mathbf{v}_t^T \quad (8)$$

$$\mathbf{o}_t = \mathbf{q}_t \mathbf{S}_t \quad (9)$$

where  $\mathbf{k}_t, \mathbf{q}_t \in \mathbb{R}^{D_k}$  and  $\mathbf{v}_t \in \mathbb{R}^{D_v}$  are the key, query, and value vectors at timestep  $t$ , and  $\mathbf{S}_t \in \mathbb{R}^{D_k \times D_v}$  is the state matrix. We consider batched inputs  $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{B \times T \times H \times D_k}$  and  $\mathbf{V} \in \mathbb{R}^{B \times T \times H \times D_v}$ . The sequence of length  $T$  is divided into  $N_C$  chunks of size  $C$ , such that  $T = N_C \times C$ . Let  $N_{SM}$  be the number of Streaming Multiprocessors (SMs) on the target GPU, representing its parallel execution capacity.

### B.2 STRATEGY 1: DECOUPLED THREE-PHASE EXECUTION

This strategy directly maps the chunk-wise parallel algorithm onto our three-phase programming model. Each phase is a separate kernel launch.

**Parallelism Analysis.** The degree of parallelism varies significantly between phases.

- **Phase 1 (Intra-Chunk State Computation):** Computations for each chunk are independent. The total number of parallel tasks is  $B \times H \times N_C$ . This phase exhibits the highest degree of parallelism.
- **Phase 2 (Inter-Chunk State Propagation):** The prefix sum (scan) is sequential across the  $N_C$  dimension. Parallelism is limited to the batch ( $B$ ) and head ( $H$ ) dimensions. Further parallelism, which we denote as  $P_{state}$ , can be achieved by partitioning the state matrix  $\mathbf{S} \in \mathbb{R}^{D_k \times D_v}$  and processing its partitions on different thread blocks. The total number of parallel tasks is  $B \times H \times P_{state}$ .
- **Phase 3 (Merge):** Similar to Phase 1, the merge operation is independent across chunks, offering a high degree of parallelism with  $B \times H \times N_C$  tasks.

**Memory Access Analysis.** The defining characteristic of this strategy is the materialization of intermediate states in global memory (GMEM) between phases.

- **Phase 1:** Reads  $\mathbf{K}$  and  $\mathbf{V}$  from GMEM. Writes the intermediate, chunk-local states  $\mathbf{S}' \in \mathbb{R}^{B \times H \times N_C \times D_k \times D_v}$  back to GMEM.

$$\text{GMEM Traffic}_{P1} = \underbrace{BHT(D_k + D_v)}_{\text{Reads}} + \underbrace{BHN_C D_k D_v}_{\text{Writes}} \quad (10)$$



- **Phase 2:** Reads the  $B \times H \times N_C$  chunk-local states from GMEM, performs the scan, and writes the updated global states  $\mathbf{S}_{global} \in \mathbb{R}^{B \times H \times N_C \times D_k \times D_v}$  back to GMEM.

$$\text{GMEM Traffic}_{P2} = \underbrace{BHN_C D_k D_v}_{\text{Reads}} + \underbrace{BHN_C D_k D_v}_{\text{Writes}} = 2BHN_C D_k D_v \quad (11)$$

The total GMEM traffic introduced between Phase 1 and Phase 2 is  $3 \times B \times H \times N_C \times D_k D_v$ .

### B.3 STRATEGY 2: FUSED INTRA- AND INTER-CHUNK EXECUTION

This strategy fuses Phase 1 and Phase 2, computing the local state of a chunk and immediately uses it to update the global state, all within on-chip memory.

**Parallelism Analysis.** By fusing the phases, the execution is constrained by the most sequential part, which is the inter-chunk scan. Therefore, the maximum number of parallel tasks is identical to that of Phase 2 in the decoupled strategy:  $B \times H \times P_{state}$ . The high parallelism across the chunk dimension ( $N_C$ ) is sacrificed.

**Memory Access Analysis.** The primary advantage of this strategy is the elimination of the intermediate GMEM traffic.

- **Fused Phase (1+2):** Each of the  $B \times H \times P_{state}$  parallel tasks loads its corresponding slice of  $\mathbf{K}$  and  $\mathbf{V}$  from GMEM chunk by chunk. The intermediate, chunk-local states are generated and accumulated entirely within SMEM. The only state written to GMEM is the final updated global state for each chunk, which is required by the Merge phase.

$$\text{GMEM Traffic}_{P1+P2} = \underbrace{BHT(D_k + D_v)}_{\text{Reads}} + \underbrace{BHN_C D_k D_v}_{\text{Writes}} \quad (12)$$

Compared to Strategy 1, this approach saves  $2 \times B \times H \times N_C \times D_k D_v$  worth of GMEM traffic, which is the cost of one full read and one full write of the intermediate state tensor.

### B.4 HEURISTIC SCHEDULING ALGORITHM

The choice between the Decoupled and Fused strategies presents a classic trade-off between parallelism and memory locality.

- **Strategy 1 (Decoupled)** is favored when the GPU has a high degree of parallelism ( $N_{SM}$  is large) that is not saturated by the task parallelism of Strategy 2 ( $B \times H \times P_{state}$ ). The performance gain from launching more parallel tasks in Phase 1 and 3 must outweigh the latency incurred by the extra GMEM I/O.
- **Strategy 2 (Fused)** is favored when the task parallelism of the scan ( $B \times H \times P_{state}$ ) is already sufficient to saturate the GPU’s SMs, or when the cost of reading/writing the intermediate states is the dominant performance bottleneck.

We can formulate a simple heuristic to guide this choice as shown in Algo 1.

In practice, we set  $T_{\text{comp\_chunk}}$  to 0 to reduce the runtime overhead and omit the micro-benchmark effort. The insight is that most linear attention variants is memory bound instead of compute-intensive.

## C COMPUTATIONAL PATTERNS AND MODEL COVERAGE

This appendix elaborates on the claim that our implemented kernels for a few representative models can support a much broader range of linear attention variants. The mapping from our representative implementations to the computational patterns and the models they cover is detailed below and summarized in Table 3.

**Algorithm 1** Parallelism Scheduler)

---

**Require:** Shapes  $B, T, H, d_k, d_v$ , chunk size  $C$ , tile sizes  $t_k, t_v$   
**Require:** Device params: memory bandwidth  $BW$  (bytes/s),  $P_{\max}$  (hardware concurrency cap)  
**Require:** Cost param:  $T_{\text{comp\_chunk}}$  (measured per-chunk compute time in seconds)  
**Require:** Thresholds:  $P_{\min}$  (saturation threshold, tasks), element size  $s$  (bytes)  
**Ensure:** Return strategy  $\in \{\text{FUSED}, \text{DECOUPLED}\}$

```

1:  $N_c \leftarrow \lceil T/C \rceil$ 
2:  $P_{\text{state}} \leftarrow \lceil d_k/t_k \rceil \cdot \lceil d_v/t_v \rceil$ 
3:  $P_{\text{dec}} \leftarrow B \cdot H \cdot N_c \cdot P_{\text{state}}$ 
4:  $P_{\text{fused}} \leftarrow B \cdot H \cdot P_{\text{state}}$ 
5:  $\text{GMEM}_{\text{dec}} \leftarrow (2C(d_k + d_v) + 5d_k d_v) \cdot s$ 
6:  $\text{GMEM}_{\text{fused}} \leftarrow (2C(d_k + d_v) + 4d_k d_v) \cdot s$ 
    $\triangleright$  Fast path: if fused already provides enough parallelism, prefer fused to save GMEM IO
7: if  $P_{\text{fused}} \geq P_{\min}$  then
8:   return FUSED
9: end if
    $\triangleright$  Otherwise estimate per-chunk runtime under each mapping
10:  $\tilde{P}_{\text{dec}} \leftarrow \min(P_{\text{dec}}, P_{\max})$ 
11:  $\tilde{P}_{\text{fused}} \leftarrow \min(P_{\text{fused}}, P_{\max})$ 
12:  $\hat{T}_{\text{dec}} \leftarrow \frac{\text{GMEM}_{\text{dec}}}{BW \cdot \tilde{P}_{\text{dec}}} + \frac{T_{\text{comp\_chunk}}}{\tilde{P}_{\text{dec}}}$ 
13:  $\hat{T}_{\text{fused}} \leftarrow \frac{\text{GMEM}_{\text{fused}}}{BW \cdot \tilde{P}_{\text{fused}}} + \frac{T_{\text{comp\_chunk}}}{\tilde{P}_{\text{fused}}}$ 
14: if  $\hat{T}_{\text{fused}} \leq \hat{T}_{\text{dec}}$  then
15:   return FUSED
16: else
17:   return DECOUPLED
18: end if

```

---

**Vector-State Linear RNNs (represented by HGRN)** Our HGRN kernel embodies the computational pattern of linear recurrent networks where the state is a vector, updated via element-wise operations with gated inputs. This is a common pattern for models aiming for high efficiency with a compact state. Models sharing this fundamental structure include the original HGRN (Qin et al., 2023) and Hawk (RG-LRU) (De et al., 2024).

**Matrix-State with Scalar Decay (represented by scalar-GLA)** The scalar-GLA kernel represents the widely adopted matrix-state linear attention pattern. In this formulation, the state is a matrix updated via an outer product of key and value vectors, combined with a simple data-dependent or data-independent scalar decay. This pattern is foundational to many prominent and powerful models such as RetNet (Sun et al., 2023), Mamba-2 (Dao & Gu, 2024), and Lightning Attention (Qin et al., 2024).

**Matrix-State with Vector Decay (represented by vector-GLA)** Our vector-GLA implementation captures the pattern of matrix-state models that employ a more expressive, data-dependent vector decay. This allows for per-feature state transition dynamics, offering a richer representation than a single scalar decay. Models in this category include the original Gated Linear Attention (GLA) (Yang et al., 2023), HGRN-2 (Qin et al., b), and RWKV-6 (Peng et al., 2023).

**Delta-Rule Updates (represented by GDN)** Finally, our Gated DeltaNet (GDN) kernel is representative of a family of models based on the delta rule for state updates. This update mechanism can be interpreted as applying a series of Householder transformations to the state matrix, enabling more complex state transitions. This pattern is central to the family of DeltaNets, including the original DeltaNet (Yang et al., 2024b), GatedDeltaNet (Yang et al., 2024a), and DeltaProduct (Siems et al., 2025).

Table 3: Mapping of representative kernels implemented in subsection 4.1 to the broader set of models they cover. This demonstrates the generality of our compiler.

Representative Kernel	Covered Models (Examples)
HGRN	HGRN (Qin et al., 2023), Hawk (RG-LRU) (De et al., 2024)
Scalar-GLA	RetNet (Sun et al., 2023), Mamba-2 (Dao & Gu, 2024), Lightning Attention (Qin et al., 2024)
Vector-GLA	GLA (Yang et al., 2023), HGRN-2 (Qin et al., b), RWKV-6 (Peng et al., 2023)
GDN	DeltaNet (Yang et al., 2024b), GatedDeltaNet (Yang et al., 2024a), DeltaProduct (Siems et al., 2025)

## D DISCUSSION ON EXTENSIBILITY AND SUSTAINABILITY

The longevity and utility of a Domain-Specific Language (DSL) depend heavily on whether its abstraction captures the invariant properties of the target domain rather than transient heuristics. In this section, we discuss the future-proofing of Forge from two perspectives: algorithmic expressiveness and hardware sustainability.

### D.1 ALGORITHMIC EXPRESSIVENESS: GROUNDED IN ASSOCIATIVITY

The universality of Forge’s three-phase abstraction (Intra-Chunk, Inter-Chunk, and Merge) derives from the mathematical foundation of efficient sequence modeling: the **associative property**.

**Mathematical Invariance.** Virtually all Linear Attention and State Space Duality (SSD) models aim to achieve  $O(N)$  computational complexity by formulating the attention mechanism as a recurrence or a parallel prefix scan. Mathematically, any algorithm that can be decomposed into a chunk parallel form fits the Forge abstraction. This is not an ad-hoc design choice but a direct mapping of the underlying algebraic structure required for parallelization.

**Handling Complex Dependencies.** Forge is expressive enough to handle complex state updates found in modern architectures, provided they satisfy chunk-wise associativity. For instance:

- **The Delta Rule:** Despite involving data-dependent updates (e.g.,  $h_t = h_{t-1} + \beta_t(v_t - h_{t-1})$ ), the Delta Rule preserves the associative structure within local chunks, allowing Forge to effectively fuse operations.
- **Element-wise Decay:** Varying decay rates (as seen in Vector-Gated GLA) are fully supported, as the element-wise multiplication distributes over addition, maintaining the scan property.

**Theoretical Limitations and Edge Cases.** The boundary of Forge’s applicability is strictly defined by **non-associative recurrences**. If an architecture introduces a dependency where the state update  $h_t = f(h_{t-1}, x_t)$  involves a non-linear function  $f$  that prevents parallel(e.g., passing the hidden state through a complex MLP at every step before the next update), it cannot be expressed in Forge.

A notable example is the Test-Time Training (TTT) layer (Sun et al., 2024b). While TTT can be viewed through the lens of linear attention or fast weights, its gradient-based updates involve non-linearities that break associativity. Consequently, TTT cannot be accelerated via the parallel prefix scans used in Forge. However, it is worth noting that this limitation is mutual: by abandoning associativity, such models fundamentally forego the massively parallel efficiency that characterizes the linear attention regime targeted by Forge.

### D.2 HARDWARE SUSTAINABILITY: HIERARCHICAL DECOUPLING

To ensure sustainability amidst rapid hardware evolution, Forge employs a strict hierarchical decoupling between the algorithmic description and hardware-specific instructions.

**Layered Compilation.** Forge operates as a high-level graph compiler rather than a low-level assembler. It does not directly emit hardware-specific assembly (ISA). Instead, it lowers the PyTorch-based algorithmic description into an intermediate kernel DSL. In our current implementation, we target `Triton`, effectively delegating low-level complexities such as register allocation, instruction scheduling, and tensor core management (e.g., `WGMMA` on Hopper) to the Triton compiler.

**Backend Agnosticism.** This layered design makes Forge inherently adaptable. While the current backend is Triton, the architecture allows for retargeting the mapping layer to other emerging DSLs (e.g., ThunderKitten (Spector et al., 2024) or TileLang (Wang et al., 2025)). These DSLs provide better performance with the cost of hardware overfit (e.g. ThunderKittens target on NVidia GPU, HipKittens (Hu et al., 2025) target solely on AMD GPU while Triton support multiple hardware backend). Triton allows us to minimize the use of inline assembly (e.g. PTX on NVGPU), reserving it only for specific extensions not yet exposed by the intermediate representation.

**Forward Compatibility.** Consequently, as hardware architectures evolve (e.g., the transition to NVIDIA Blackwell), Forge benefits automatically from updates made by the community to the intermediate compiler stack. This ensures that user-level code remains stable and performant without requiring modification, solving the maintenance bottleneck often associated with hand-written kernels.

## E DISCUSSION ON GENERALIZABILITY TO SOFTMAX ATTENTION

While Forge is explicitly designed for Linear Attention, the underlying mathematical associativity, which enables our optimization is shared by Softmax attention (exploited via the online softmax trick (Dao et al., 2022)). In this section, we analyze the theoretical feasibility of extending Forge to support Softmax attention, and the design rationale behind our decision to specialize in the Linear Attention domain.

Extending the Forge abstraction to support Softmax-based mechanisms is theoretically feasible but would necessitate two primary modifications to the current three-phase abstraction:

- **Generalizing the Inter-Chunk State.** In Linear Attention, the state propagated between chunks is typically a feature map (e.g.,  $S \in \mathbb{R}^{d \times d}$ ) governed by a linear recurrence. In contrast, Softmax attention requires the propagation of normalization statistics, specifically the running maximum  $m$  and the running sum  $l$  to stabilize the computation of exponentials. Extending Forge would require modifying the `decay_mode` phase to support the propagation of these scalar or vector statistics alongside, or instead of, the matrix state.
- **Introducing a Rescaling Primitive.** The `merge_mode` in Forge is currently designed for linear combinations (accumulation). Softmax attention, however, mandates a renormalization step when merging partial results. Specifically, the output of a preceding block  $O_1$  must be scaled down based on the difference between its local maximum  $m_1$  and the new global maximum  $m_{new}$ :

$$O_{new} = O_1 \times e^{m_1 - m_{new}} + O_2 \times e^{m_2 - m_{new}} \quad (13)$$

Supporting this would require introducing a native `rescale(output, old_stats, new_stats)` primitive into the Forge DSL.

**Despite the feasibility of these extensions, we deliberately narrowed the scope of Forge to Linear Attention to maximize expressiveness.**

**The Tension between Universality and Specificity.** Constructing a "universal" DSL often necessitates a rigid abstraction that compromises the ability to model complex, domain-specific patterns. A recent framework, AttentionEngine (Chen et al., 2025), attempts to unify both Softmax and Linear Attention. However, to maintain this broad generality, its abstraction struggles to express advanced Linear Attention variants that feature intricate dependencies, such as the **Delta Rule** or **Vector-Gated GLA**. By specializing in Linear Attention, Forge avoids these constraints, supporting these complex patterns effortlessly.

From an ecosystem perspective, Softmax attention is already well-served by highly optimized libraries and compilers (e.g., FlexAttention (Dong et al., 2024)). In contrast, the Linear Attention landscape is characterized by rapid algorithmic fragmentation, where new variants are proposed frequently but lack efficient kernel support. Forge addresses this specific "N-to-1" compiler challenge, solving a critical bottleneck that is currently more pressing for the research community than further optimizing Softmax kernels.

Finally, while the abstraction of Forge is specialized, the *system-level optimizations* are broadly applicable. For example, our Ahead-of-Time (AOT) compilation pipeline and static dispatcher, designed to eliminate runtime Python overheads, are universal optimizations. These techniques can be directly adopted by Softmax-focused DSLs (such as FlexAttention) to significantly improve performance, particularly in short-sequence regimes where dispatch latency is a dominant factor.

## F EXPERIMENTS: DISTRIBUTED EVALUATION ON H100 CLUSTER

To verify the robustness of our distributed scaling strategy, we conducted additional experiments on an **8x H100 node** connected via NVLink.

Table 4 presents the latency comparison between Ring-Attention (a standard distributed softmax attention baseline) (Liu et al., 2023) with FlashAttention-3 which exploit advanced hardware feature on Hopper GPU, LASP-2, and Forge (both LASP2 and Forge using the Scalar GLA variant). The experiments were conducted with a **batch size of 1**, NumHeads = 32, and HeadDim = 128. The **SeqLen** in table means sequence length per GPU(global sequence length is equal to **SeqLen**  $\times$  **NumberOfGPUs**).

The result shows Forge consistently outperforms the LASP-2 baseline across all sequence lengths. For instance, at a sequence length of 16k on 4 GPUs, Forge achieves a **1.17 $\times$  speedup** over LASP-2 (0.96ms vs. 1.13ms). Even with the high bandwidth of NVLink on H100s, Forge maintains its efficiency advantage. This confirms that our fine-grained compute-communication overlap strategy is effective not only on bandwidth-constrained hardware (like the H20 used in the main result) but also on high-performance flagship clusters. As expected, Linear Attention methods (both LASP-2 and Forge) are orders of magnitude faster than Ring-Attention, especially at longer sequences (e.g., at 512k global sequence length, Forge is over **160 $\times$  faster** than Ring-Attention).

These results validate that the performance gains reported in the main paper are not artifacts of the H20 hardware or larger batch sizes, but rather stem from the fundamental efficiency of Forge’s generated kernels and scheduling logic.

Table 4: **Distributed Latency Comparison on H100 GPUs (Batch Size = 1)**. We compare the end-to-end latency of Forge (Scalar GLA) against Ring-Attention (Softmax) and LASP-2 (State-of-the-art Linear Attention Sequence Parallelism). Forge consistently achieves the lowest latency across all sequence lengths and GPU configurations.

# GPUs	SeqLen Per GPU	Latency (ms)		
		Ring-Attn	LASP-2	Forge (Ours)
4	8192	8.40	1.02	<b>0.79</b>
	16384	17.77	1.13	<b>0.96</b>
	32768	61.65	1.66	<b>1.57</b>
	65536	229.43	2.90	<b>2.76</b>
8	8192	19.32	1.00	<b>0.77</b>
	16384	38.09	1.15	<b>0.97</b>
	32768	133.14	1.67	<b>1.62</b>
	65536	462.30	2.92	<b>2.84</b>

## G CODE EXAMPLE

In this section we provide code examples to implement different linear attention variants and the code generate by Forge. Firt we show the Scalar GLA in Listing 1 and the generated code in List-

ing 2. Forge enable pre-compiled kernel (AOT) and static dispatch at runtime. In this example, we tune the kernel on  $H=[1,4,8,16,32]$  and  $D=[64, 128]$ , so there are a lot of auto-generated code encoded pre-tuned information for dispatcher (our AOT optimization). Using Forge to implement this kernel only involves 50+ lines of code (LOC), which the generated triton code and host launcher is around 1200 LOC. Even without including the static dispatcher and host launcher in count, the generated Triton kernel alone has over 400 LOC. And this is only partially generated content, because Forge also generates code for other parallel strategies.

Then we demonstrate the implementation of DeltaNet using Forge in Listing 3. Note that DeltaNet involves complex data dependencies where intermediate results computed in the chunk phase (specifically  $U$  and  $W$ ) are required in the merge phase. Forge introduces a primitive `forge.cache_result` to explicit mark these tensors. The compiler then automatically handles the memory layout and data movement to ensure these values are efficiently reused across phases without redundant re-computation or manual memory management.

### G.1 IMPLEMENTATION AND CODE GENERATION FOR SCALAR GLA

```

1042 1 def chunk_mode_scalar_gla(k: Tensor, v: Tensor, g: Tensor) -> Tensor:
1043 2     """k: [C, K], v: [C, V], g: [C]"""
1044 3     g_cumsum = g.cumsum(dim=0)
1045 4     g_cumsum_last = g.sum(dim=0)
1046 5     g_cumsum = (g_cumsum_last - g_cumsum).unsqueeze(0).exp()
1047 6     k_fp32 = k.permute([1, 0]).to(g.dtype)
1048 7     k_decay = (k_fp32 * g_cumsum).to(v.dtype)
1049 8     chunk_state = k_decay @ v
1050 9     return chunk_state
1051 10
1051 11 def decay_mode_scalar_gla(prev_s: Tensor, chunk_state: Tensor, g: Tensor)
1052 12     -> Tensor:
1053 13     """g: [C]"""
1054 14     g_sum = g.sum(dim=0).exp()
1055 15     return prev_s * g_sum + chunk_state
1056 16
1056 17 def merge_mode_scalar_gla(
1057 18     q: Tensor,
1058 19     k: Tensor,
1059 20     v: Tensor,
1060 21     g: Tensor,
1061 22     chunk_state: Tensor,
1062 23     scale: Tensor,
1063 24 ) -> Tensor:
1064 25     """q: [C, K], k: [C, K], v: [C, V], chunk_state: [K, V]"""
1065 26     g_cumsum = g.cumsum(0)
1066 27     chunk_state = chunk_state.to(q.dtype)
1067 28     p = (q @ k.T).tril(0)
1068 29     p = (p * (g_cumsum[..., None] - g_cumsum[None, ...]).exp()).to(v.
1069 30     dtype)
1070 31     return (p @ v + q @ chunk_state * g_cumsum.exp()[..., None]) * scale
1071 32
1072 33 CONST_H = ConstExpr("H", H)
1073 34 CONST_K = ConstExpr("K", K)
1074 35 CONST_V = ConstExpr("V", V)
1075 36 meta = {
1076 37     "q": SymbTensor(["T", CONST_H, CONST_K], dtype=dtype),
1077 38     "k": SymbTensor(["T", CONST_H, CONST_K], dtype=dtype),
1078 39     "v": SymbTensor(["T", CONST_H, CONST_V], dtype=dtype),
1079 40     "g": SymbTensor(["T", CONST_H], dtype=torch.float32),
1080 41     "prev_s": SymbTensor(["NS", CONST_H, CONST_K, CONST_V], dtype=torch.
1081 42     float32),
1082 43     "chunk_state": SymbTensor(["NC", CONST_H, CONST_K, CONST_V], dtype=
1083 44     dtype),
1084 45     "scale": 1 / math.sqrt(K),
1085 46 }

```



```

1080 43 LinearAttention(
1081 44     input_meta=meta,
1082 45     sp_group=pg if args.sp else None,
1083 46     enable_aot=args.aot,
1084 47     code_dir=args.dir,
1085 48     chunk_mode=chunk_mode_scalar_gla,
1086 49     decay_mode=decay_mode_scalar_gla,
1087 50     merge_mode=merge_mode_scalar_gla,
1088 51 )

```

Listing 1: Implementation of Scalar GLA in Forge

```

1090 1 fuse_chunk_decay_kernel_signature = (
1091 2     "*bf16:16, "
1092 3     "*bf16:16, "
1093 4     "*fp32, "
1094 5     "*bf16:16, "
1095 6     "*bf16:16, "
1096 7     "*fp32, "
1097 8     "*i32, "
1098 9     "i32, "
1099 10    "i32, "
1100 11    "*i32, "
1101 12    "%USE_INITIAL_STATE, "
1102 13    "%H, "
1103 14    "%K, "
1104 15    "%V, "
1105 16    "%CHUNK, "
1106 17    "%BLK_K, "
1107 18    "%BLK_V, "
1108 19    "%USE_TMA"
1109 20 )
1110 21
1111 22 def get_fuse_chunk_decay_kernel_info(B: int, T: int, H: int, K: int, V:
1112 23 int):
1113 24     """Static dispatcher for fuse_chunk_decay_kernel. Auto-generated."""
1114 25     D = [K, V]
1115 26     key = (B, T, H, D)
1116 27     if key in (
1117 28         (1, 4096, 4, [64, 64]),
1118 29         (2, 1024, 4, [64, 64]),
1119 30     ):
1120 31         BLK_K = 32
1121 32         BLK_V = 32
1122 33         num_warps = 8
1123 34         num_ctas = 1
1124 35         num_stages = 3
1125 36         maxnreg = None
1126 37     elif key in (
1127 38         (1, 1024, 1, [64, 64]),
1128 39         (1, 1024, 1, [128, 128]),
1129 40         (1, 1024, 4, [64, 64]),
1130 41         (1, 1024, 4, [128, 128]),
1131 42         (1, 1024, 8, [64, 64]),
1132 43         (1, 1024, 8, [128, 128]),
1133 44         (1, 1024, 16, [64, 64]),
1134 45         (1, 1024, 32, [64, 64]),
1135 46         (1, 2048, 1, [64, 64]),
1136 47         (1, 2048, 1, [128, 128]),
1137 48         (1, 2048, 4, [64, 64]),
1138 49         (1, 2048, 4, [128, 128]),
1139 50         (1, 2048, 8, [64, 64]),
1140 51         (1, 2048, 8, [128, 128]),
1141 52         (1, 2048, 16, [64, 64]),

```

1134	53	(1, 2048, 32, [64, 64]),
1135	54	(1, 4096, 1, [64, 64]),
1136	55	(1, 4096, 1, [128, 128]),
1137	56	(1, 4096, 4, [128, 128]),
1138	57	(1, 4096, 8, [64, 64]),
1139	58	(1, 4096, 8, [128, 128]),
1140	59	(1, 4096, 16, [64, 64]),
1141	60	(1, 4096, 32, [64, 64]),
1142	61	(1, 8192, 1, [64, 64]),
1143	62	(1, 8192, 1, [128, 128]),
1144	63	(1, 8192, 4, [64, 64]),
1145	64	(1, 8192, 4, [128, 128]),
1146	65	(1, 8192, 8, [64, 64]),
1147	66	(1, 8192, 8, [128, 128]),
1148	67	(1, 8192, 16, [64, 64]),
1149	68	(1, 8192, 32, [64, 64]),
1150	69	(1, 16384, 1, [64, 64]),
1151	70	(1, 16384, 1, [128, 128]),
1152	71	(1, 16384, 4, [64, 64]),
1153	72	(1, 16384, 4, [128, 128]),
1154	73	(1, 16384, 8, [64, 64]),
1155	74	(1, 16384, 8, [128, 128]),
1156	75	(1, 16384, 16, [64, 64]),
1157	76	(1, 16384, 32, [64, 64]),
1158	77	(1, 32768, 1, [64, 64]),
1159	78	(1, 32768, 1, [128, 128]),
1160	79	(1, 32768, 4, [64, 64]),
1161	80	(1, 32768, 4, [128, 128]),
1162	81	(1, 32768, 8, [64, 64]),
1163	82	(1, 32768, 8, [128, 128]),
1164	83	(1, 32768, 16, [64, 64]),
1165	84	(1, 32768, 32, [64, 64]),
1166	85	(1, 65536, 1, [64, 64]),
1167	86	(1, 65536, 1, [128, 128]),
1168	87	(1, 65536, 4, [64, 64]),
1169	88	(1, 65536, 4, [128, 128]),
1170	89	(1, 65536, 8, [64, 64]),
1171	90	(1, 65536, 8, [128, 128]),
1172	91	(1, 65536, 16, [64, 64]),
1173	92	(1, 65536, 32, [64, 64]),
1174	93	(1, 131072, 1, [64, 64]),
1175	94	(1, 131072, 1, [128, 128]),
1176	95	(1, 131072, 4, [64, 64]),
1177	96	(1, 131072, 4, [128, 128]),
1178	97	(1, 131072, 8, [64, 64]),
1179	98	(1, 131072, 8, [128, 128]),
1180	99	(1, 131072, 16, [64, 64]),
1181	100	(1, 131072, 32, [64, 64]),
1182	101	(1, 262144, 1, [64, 64]),
1183	102	(1, 262144, 1, [128, 128]),
1184	103	(1, 262144, 4, [64, 64]),
1185	104	(1, 262144, 4, [128, 128]),
1186	105	(1, 262144, 8, [64, 64]),
1187	106	(1, 262144, 8, [128, 128]),
	107	(1, 262144, 16, [64, 64]),
	108	(1, 262144, 32, [64, 64]),
	109	(2, 1024, 1, [64, 64]),
	110	(2, 1024, 1, [128, 128]),
	111	(2, 1024, 4, [128, 128]),
	112	(2, 1024, 8, [64, 64]),
	113	(2, 1024, 16, [64, 64]),
	114	(2, 2048, 1, [64, 64]),
	115	(2, 2048, 1, [128, 128]),
	116	(2, 2048, 4, [64, 64]),
	117	(2, 2048, 4, [128, 128]),

```

1188 (2, 2048, 8, [64, 64]),
1189 (2, 2048, 16, [64, 64]),
1190 (2, 4096, 1, [64, 64]),
1191 (2, 4096, 1, [128, 128]),
1192 (2, 4096, 4, [64, 64]),
1193 (2, 4096, 4, [128, 128]),
1194 (2, 4096, 8, [64, 64]),
1195 (2, 4096, 16, [64, 64]),
1196 (2, 8192, 1, [64, 64]),
1197 (2, 8192, 1, [128, 128]),
1198 (2, 8192, 4, [64, 64]),
1199 (2, 8192, 4, [128, 128]),
1200 (2, 8192, 8, [64, 64]),
1201 (2, 8192, 16, [64, 64]),
1202 (2, 16384, 1, [64, 64]),
1203 (2, 16384, 1, [128, 128]),
1204 (2, 16384, 4, [64, 64]),
1205 (2, 16384, 4, [128, 128]),
1206 (2, 16384, 8, [64, 64]),
1207 (2, 16384, 16, [64, 64]),
1208 (2, 32768, 1, [64, 64]),
1209 (2, 32768, 1, [128, 128]),
1210 (2, 32768, 4, [64, 64]),
1211 (2, 32768, 4, [128, 128]),
1212 (2, 32768, 8, [64, 64]),
1213 (2, 32768, 16, [64, 64]),
1214 (2, 65536, 1, [64, 64]),
1215 (2, 65536, 1, [128, 128]),
1216 (2, 65536, 4, [64, 64]),
1217 (2, 65536, 4, [128, 128]),
1218 (2, 65536, 8, [64, 64]),
1219 (2, 65536, 16, [64, 64]),
1220 (2, 131072, 1, [64, 64]),
1221 (2, 131072, 1, [128, 128]),
1222 (2, 131072, 4, [64, 64]),
1223 (2, 131072, 4, [128, 128]),
1224 (2, 131072, 8, [64, 64]),
1225 (2, 131072, 16, [64, 64]),
1226 ):
1227     BLK_K = 32
1228     BLK_V = 32
1229     num_warps = 8
1230     num_ctas = 1
1231     num_stages = 4
1232     maxnreg = None
1233     elif key in (
1234         (1, 1024, 32, [128, 128]),
1235         (1, 2048, 32, [128, 128]),
1236         (1, 4096, 32, [128, 128]),
1237         (1, 8192, 32, [128, 128]),
1238         (1, 16384, 32, [128, 128]),
1239         (1, 32768, 32, [128, 128]),
1240         (1, 65536, 32, [128, 128]),
1241         (1, 131072, 32, [128, 128]),
1242         (1, 262144, 32, [128, 128]),
1243         (2, 1024, 32, [128, 128]),
1244         (2, 2048, 16, [128, 128]),
1245         (2, 2048, 32, [128, 128]),
1246         (2, 4096, 16, [128, 128]),
1247         (2, 4096, 32, [128, 128]),
1248         (2, 8192, 32, [128, 128]),
1249         (2, 16384, 16, [128, 128]),
1250         (2, 16384, 32, [128, 128]),
1251         (2, 32768, 32, [128, 128]),
1252         (2, 131072, 16, [128, 128]),

```

```

1242 183 ) :
1243 184     BLK_K = 32
1244 185     BLK_V = 64
1245 186     num_warps = 4
1246 187     num_ctas = 1
1247 188     num_stages = 3
1248 189     maxnreg = None
1249 190     elif key in (
1250 191         (2, 1024, 16, [128, 128]),
1251 192     ):
1252 193         BLK_K = 32
1253 194         BLK_V = 64
1254 195         num_warps = 8
1255 196         num_ctas = 1
1256 197         num_stages = 3
1257 198         maxnreg = None
1258 199         elif key in (
1259 200             (1, 1024, 16, [128, 128]),
1260 201             (1, 2048, 16, [128, 128]),
1261 202             (1, 4096, 16, [128, 128]),
1262 203             (1, 8192, 16, [128, 128]),
1263 204             (1, 16384, 16, [128, 128]),
1264 205             (1, 32768, 16, [128, 128]),
1265 206             (1, 65536, 16, [128, 128]),
1266 207             (1, 131072, 16, [128, 128]),
1267 208             (1, 262144, 16, [128, 128]),
1268 209             (2, 1024, 8, [128, 128]),
1269 210             (2, 1024, 32, [64, 64]),
1270 211             (2, 2048, 8, [128, 128]),
1271 212             (2, 2048, 32, [64, 64]),
1272 213             (2, 4096, 8, [128, 128]),
1273 214             (2, 4096, 32, [64, 64]),
1274 215             (2, 8192, 8, [128, 128]),
1275 216             (2, 8192, 16, [128, 128]),
1276 217             (2, 8192, 32, [64, 64]),
1277 218             (2, 16384, 8, [128, 128]),
1278 219             (2, 16384, 32, [64, 64]),
1279 220             (2, 32768, 8, [128, 128]),
1280 221             (2, 32768, 16, [128, 128]),
1281 222             (2, 32768, 32, [64, 64]),
1282 223             (2, 65536, 8, [128, 128]),
1283 224             (2, 65536, 16, [128, 128]),
1284 225             (2, 65536, 32, [64, 64]),
1285 226             (2, 131072, 8, [128, 128]),
1286 227             (2, 131072, 32, [64, 64]),
1287 228     ):
1288 229         BLK_K = 32
1289 230         BLK_V = 64
1290 231         num_warps = 8
1291 232         num_ctas = 1
1292 233         num_stages = 4
1293 234         maxnreg = None
1294 235         elif key in (
1295 236             (2, 65536, 32, [128, 128]),
1296 237             (2, 131072, 32, [128, 128]),
1297 238     ):
1298 239         BLK_K = 64
1299 240         BLK_V = 128
1300 241         num_warps = 8
1301 242         num_ctas = 1
1302 243         num_stages = 4
1303 244         maxnreg = None
1304 245     else:
1305 246         raise ValueError(f"Unsupported config for fuse_chunk_decay_kernel
: BTHD={key}")

```

```

1296 247
1297 248     return {
1298 249         "CHUNK": 64,
1299 250         "USE_INITIAL_STATE": True,
1300 251         "H": H,
1301 252         "K": D[0],
1302 253         "V": D[1],
1303 254         "BLK_K": BLK_K,
1304 255         "BLK_V": BLK_V,
1305 256         "num_warps": num_warps,
1306 257         "num_stages": num_stages,
1307 258     }
1308 259
1309 260
1310 261 @aut_compile_spaces({
1311 262     "fuse_chunk_decay_kernel": {
1312 263         "signature": fuse_chunk_decay_kernel_signature,
1313 264         "grid": ["(%K + %BLK_K - 1) / %BLK_K", "(%V + %BLK_V - 1) / %
1314 265         BLK_V", "H_MUL_NS"],
1315 266         "triton_algo_infos": [
1316 267             get_fuse_chunk_decay_kernel_info(B, T, H, K, V)
1317 268             for B,T,H,(K,V) in[
1318 269                 (1, 262144, 32, [64, 64]),
1319 270                 (1, 262144, 8, [128, 128]),
1320 271                 (2, 1024, 16, [128, 128]),
1321 272                 (2, 1024, 4, [64, 64]),
1322 273                 (2, 131072, 1, [128, 128]),
1323 274                 (2, 131072, 1, [64, 64]),
1324 275                 (2, 131072, 16, [128, 128]),
1325 276                 (2, 131072, 16, [64, 64]),
1326 277                 (2, 131072, 32, [128, 128]),
1327 278                 (2, 131072, 32, [64, 64]),
1328 279                 (2, 131072, 4, [128, 128]),
1329 280                 (2, 131072, 4, [64, 64]),
1330 281                 (2, 131072, 8, [128, 128]),
1331 282                 (2, 131072, 8, [64, 64]),
1332 283                 (2, 32768, 32, [128, 128]),
1333 284                 (2, 65536, 16, [128, 128])
1334 285             ]
1335 286         ],
1336 287     },
1337 288 })
1338 289 @triton.autotune(
1339 290     configs=[
1340 291         triton.Config({"BLK_K": BLK_K, "BLK_V": BLK_V, "USE_TMA": USE_TMA
1341 292         }, num_warps=num_warps, num_stages=num_stages)
1342 293         for num_warps in [4, 8]
1343 294         for num_stages in [3, 4]
1344 295         for BLK_K in [128, 64, 32]
1345 296         for BLK_V in [128, 64, 32]
1346 297         for USE_TMA in [True, False]
1347 298     ],
1348 299     key=[],
1349 300 )
1350 301 @triton.jit
1351 302 def fuse_chunk_decay_kernel(
1352 303     k,
1353 304     v,
1354 305     g,
1355 306     prev_s,
1356 307     out_0,
1357 308     out_1,
1358 309     cu_seqlens,

```

```

1350 310 NS,
1351 311 H_MUL_NS,
1352 312 chunk_offsets_with_ini,
1353 313 USE_INITIAL_STATE: tl.constexpr,
1354 314 H: tl.constexpr,
1355 315 K: tl.constexpr,
1356 316 V: tl.constexpr,
1357 317 CHUNK: tl.constexpr,
1358 318 BLK_K: tl.constexpr,
1359 319 BLK_V: tl.constexpr,
1360 320 USE_TMA: tl.constexpr,
1361 321 ):
1362 322 NUM_BLK_K = (K + BLK_K - 1) // BLK_K
1363 323 NUM_BLK_V = (V + BLK_V - 1) // BLK_V
1364 324 NUM_BLK_KV = NUM_BLK_K * NUM_BLK_V
1365 325
1366 326 i_k, i_v, i_sh = tl.program_id(0), tl.program_id(1), tl.program_id(2)
1367 327 i_s, i_h = i_sh // H, i_sh % H
1368 328 bos, eos = tl.load(cu_seqlens + i_s).to(tl.int32), tl.load(cu_seqlens
1369 329 + i_s + 1).to(tl.int32)
1370 330 T = eos - bos
1371 331 NC = tl.cdiv(T, CHUNK)
1372 332 boh = tl.load(chunk_offsets_with_ini + i_s).to(tl.int32)
1373 333
1374 334 out_0 = out_0 + (boh * H + i_h).to(tl.int64) * K * V
1375 335 prev_s = prev_s + (i_s * H + i_h) * K * V
1376 336 initial_decay = 1.0
1377 337 # out_1: [NC + NS, H]. NOTE: **not** in log space
1378 338 ptr_out_1 = out_1 + boh * H + i_h
1379 339 # store initial decay
1380 340 tl.store(ptr_out_1, initial_decay)
1381 341 ptr_out_1 += H
1382 342 # [BK, BV]
1383 343 blk_prev_s = tl.zeros([BLK_K, BLK_V], dtype=tl.float32)
1384 344 if USE_INITIAL_STATE:
1385 345 ptr_prev_s = tl.make_block_ptr(prev_s, (K, V), (V, 1), (i_k *
1386 346 BLK_K, i_v * BLK_V), (BLK_K, BLK_V), (1, 0))
1387 347 blk_prev_s = tl.load(ptr_prev_s, boundary_check=(0, 1)).to(tl.
1388 348 float32)
1389 349 ptr_out_0 = tl.make_block_ptr(out_0, (K, V), (V, 1), (i_k*BLK_K, i_v*
1390 350 BLK_V), (BLK_K, BLK_V), (1, 0)) # fmt: skip
1391 351 tl.store(ptr_out_0, blk_prev_s.to(ptr_out_0.dtype.element_ty))
1392 352 out_0 += H * K * V
1393 353
1394 354 if USE_TMA:
1395 355 k_desc = tl.make_tensor_descriptor(
1396 356 k + bos * H * K + i_h * K,
1397 357 shape=[T, K],
1398 358 strides=[H * K, 1],
1399 359 block_shape=[CHUNK, BLK_K],
1400 360 )
1401 361 v_desc = tl.make_tensor_descriptor(
1402 362 v + bos * H * V + i_h * V,
1403 363 shape=[T, V],
1404 364 strides=[H * V, 1],
1405 365 block_shape=[CHUNK, BLK_V],
1406 366 )
1407 367 out_0_desc = tl.make_tensor_descriptor(
1408 368 out_0,
1409 369 shape=[NC, K, V],
1410 370 strides=[H * K * V, V, 1],
1411 371 block_shape=[1, BLK_K, BLK_V],
1412 372 )
1413 373 for i_c in range(NC):

```



```

1404 371     # load (trans) `k`: (T, H, K,) => (BLK_K, CHUNK,)
1405 372     if not USE_TMA:
1406 373         cur_k = k + bos * H * K + i_h * K # fmt: skip
1407 374         ptr_k_0 = tl.make_block_ptr(cur_k, (K, T), (1, H * K), (i_k
1408 * BLK_K, i_c * CHUNK), (BLK_K, CHUNK), (0, 1)) # fmt: skip
1409 375         blk_k_0 = tl.load(ptr_k_0, boundary_check=(0, 1)) # fmt:
1410 skip
1411 376     else:
1412 377         blk_k_0 = k_desc.load([i_c * CHUNK, i_k * BLK_K]).trans()
1413 378
1414 379     # load `v`: (T, H, V,) => (CHUNK, BLK_V,)
1415 380     if not USE_TMA:
1416 381         cur_v = v + bos * H * V + i_h * V # fmt: skip
1417 382         ptr_v_1 = tl.make_block_ptr(cur_v, (T, V), (H * V, 1), (i_c
1418 * CHUNK, i_v * BLK_V), (CHUNK, BLK_V), (1, 0)) # fmt: skip
1419 383         blk_v_1 = tl.load(ptr_v_1, boundary_check=(1, 0)) # fmt:
1420 skip
1421 384     else:
1422 385         blk_v_1 = v_desc.load([i_c * CHUNK, i_v * BLK_V])
1423 386
1424 387
1425 388     # call_external_func: `chunk_local_cumsum` to pre-compute g
1426 389
1427 390     # load `g`: (T, H,) => (CHUNK,)
1428 391     cur_g = g + bos * H + i_h * 1 # fmt: skip
1429 392     ptr_g_2 = tl.make_block_ptr(cur_g, (T), (H), (i_c * CHUNK), (
1430 CHUNK), (0,)) # fmt: skip
1431 393     blk_g_2 = tl.load(ptr_g_2, boundary_check=(0,)) # fmt: skip
1432 394     # load_last_g: => ()
1433 395     # load `g` [(CHUNK - 1)]: (T, H,) => (,)
1434 396     cur_g = g + bos * H + i_h * 1 + (CHUNK - 1) * H # fmt: skip
1435 397     ptr_last_g_3 = cur_g + i_c * CHUNK * H # fmt: skip
1436 398     last_g_3 = tl.load(ptr_last_g_3)
1437 399     # sub: torch.Size([]), ('CHUNK',) => ('CHUNK',)
1438 400     sub_4 = last_g_3 - blk_g_2
1439 401     # unsqueeze: ('CHUNK',) => (1, 'CHUNK')
1440 402     unsqueeze_5 = sub_4[None, :]
1441 403     # exp: (1, 'CHUNK') => (1, 'CHUNK')
1442 404     exp_6 = tl.exp(unsqueeze_5)
1443 405     # to: ('BLK_K', 'CHUNK') => ('BLK_K', 'CHUNK')
1444 406     blk_k_0_float32_7 = blk_k_0.to(tl.float32)
1445 407     # mul: ('BLK_K', 'CHUNK'), (1, 'CHUNK') => ('BLK_K', 'CHUNK')
1446 408     mul_8 = blk_k_0_float32_7 * exp_6
1447 409     # to: ('BLK_K', 'CHUNK') => ('BLK_K', 'CHUNK')
1448 410     mul_8_bfloat16_9 = mul_8.to(tl.bfloat16)
1449 411     # matmul: ('BLK_K', 'CHUNK'), ('CHUNK', 'BLK_V') => ('BLK_K', '
1450 BLK_V')
1451 412     matmul_10 = tl.dot(mul_8_bfloat16_9, blk_v_1).to(mul_8_bfloat16_9
1452 .dtype)
1453 413     # exp: () => ()
1454 414     exp_1_11 = tl.exp(last_g_3)
1455 415     # mul: ('BLK_K', 'BLK_V'), () => ('BLK_K', 'BLK_V')
1456 416     mul_1_12 = blk_prev_s * exp_1_11
1457 417     # add: ('BLK_K', 'BLK_V'), ('BLK_K', 'BLK_V') => ('BLK_K', 'BLK_V
1458 ')
1459 418     add_13 = mul_1_12 + matmul_10
1460 419     # mul: ('s3',), () => ('s3',)
1461 420     mul_tensor_14 = initial_decay * exp_1_11
1462 421     # store => ('BLK_K', 'BLK_V')
1463 422     # assume output layout: [NC_WITH_INI, H, K, V]
1464 423     out_0_ty = out_0.dtype.element_ty
1465 424     if not USE_TMA:
1466 425         ptr_out_0 = tl.make_block_ptr(out_0, (K, V), (V, 1), (i_k *
1467 BLK_K, i_v * BLK_V), (BLK_K, BLK_V), (1, 0)) # fmt: skip

```

```

1458         tl.store(ptr_out_0, add_13.to(out_0_ty), boundary_check=(1,
1459         0))
1460         out_0 += H * K * V
1461     else:
1462         out_0_desc.store([i_c, i_k * BLK_K, i_v * BLK_V], add_13.to(
1463         out_0_ty)[None, :, :])
1464         blk_prev_s = add_13.to(blk_prev_s.dtype)
1465         # store => ('s3',)
1466         # output layout: [NC_WITH_INI, H, s3]
1467         out_1_ty = out_1.dtype.element_ty
1468         tl.store(ptr_out_1 + i_c * H, mul_tensor_14.to(out_1_ty))
1469         initial_decay = mul_tensor_14
1470
1471 def launch_fuse_chunk_decay(
1472     g: torch.Tensor,
1473     prev_s: torch.Tensor,
1474     v: torch.Tensor,
1475     k: torch.Tensor,
1476     cu_seqlens: torch.IntTensor,
1477     cached_results: dict,
1478     dist_scan: DistScanContext = None,
1479     lazy_update: bool = True,
1480 ) -> dict[str, torch.Tensor]:
1481     """
1482     perform "decayed scan" on 'chunked_states'
1483     """
1484     # state_dtype = chunk_state.dtype
1485     chunk_decay = None
1486     prev_rank_state_sum = None
1487
1488     # TODO: now fix chunk_size to 64
1489     CHUNK = 64
1490     chunk_indices = prepare_chunk_indices(cu_seqlens, CHUNK)
1491     chunk_offsets = prepare_chunk_offsets(cu_seqlens, CHUNK)
1492     chunk_offsets_with_ini = prepare_chunk_offsets_with_ini(cu_seqlens,
1493     CHUNK)
1494     NS, NC = len(cu_seqlens) - 1, chunk_indices.shape[0]
1495     NC_WITH_INI = NC + NS
1496     use_aot = os.environ.get('FORGE_USE_AOT', '0')
1497     FORGE_USE_AOT = True if use_aot.lower() in ['1', 'true', 'yes'] else
1498     None
1499     T, H, K, V = *k.shape, v.shape[-1]
1500     updated_states = k.new_empty([NC_WITH_INI, H, K, V])
1501
1502     def alloc_fn(size: int, alignment: int, stream: int):
1503         return torch.empty(size, device='cuda', dtype=torch.int8)
1504
1505     triton.set_allocator(alloc_fn)
1506
1507     def grid(meta):
1508         BLK_K = meta['BLK_K']
1509         BLK_V = meta['BLK_V']
1510         NUM_BLK_K = (K + BLK_K - 1) // BLK_K
1511         NUM_BLK_V = (V + BLK_V - 1) // BLK_V
1512         NUM_BLK_KV = NUM_BLK_K * NUM_BLK_V
1513         H_MUL_NS = H * NS
1514         return (NUM_BLK_K, NUM_BLK_V, H_MUL_NS,)
1515
1516     chunk_decay = torch.empty([NC_WITH_INI, H], dtype=torch.float32)
1517     if FORGE_USE_AOT is None:
1518         fuse_chunk_decay_kernel[grid](
1519             g=g,
1520             prev_s=prev_s,
1521             v=v,
1522             k=k,

```

```

1512 out_1=chunk_decay,
1513 cu_seqlens=cu_seqlens,
1514 H=H,
1515 NS=NS,
1516 K=K,
1517 V=V,
1518 CHUNK=CHUNK,
1519 H_MUL_NS=H * NS,
1520 out_0=updated_states,
1521 chunk_offsets_with_ini=chunk_offsets_with_ini,
1522 USE_INITIAL_STATE=True,
1523 )
1524
1525 else:
1526
1527     from fage.aot_utils import forge_aot_ops
1528
1529     algo_info = forge_aot_ops.
1530     fuse_chunk_decay_kernel__triton_algo_info_t()
1531     for _k, _v in get_fuse_chunk_decay_kernel_info(NS, T, H, K, V).
1532     items():
1533         setattr(algo_info, _k, _v)
1534         forge_aot_ops.fuse_chunk_decay_kernel(
1535             0, # torch.cuda.current_stream().cuda_stream,
1536             k.data_ptr(), # k
1537             v.data_ptr(), # v
1538             g.data_ptr(), # g
1539             prev_s.data_ptr(), # prev_s
1540             updated_states.data_ptr(), # out_0
1541             chunk_decay.data_ptr(), # out_1
1542             cu_seqlens.data_ptr(), # cu_seqlens
1543             NS, # NS
1544             H * NS, # H_MUL_NS
1545             chunk_offsets_with_ini.data_ptr(), # chunk_offsets_with_ini
1546             algo_info,
1547         )
1548
1549     final_chunk_indices = chunk_offsets_with_ini[1:] - 1
1550     final_state_local = updated_states[final_chunk_indices, ...]
1551     final_decay_local = chunk_decay[final_chunk_indices, ...]
1552
1553     if dist_scan.pg.size() > 1:
1554         prev_rank_state_sum = dist_scan.forward(
1555             final_state_local=final_state_local,
1556             final_decay_local=final_decay_local,
1557             decay_type=DecayType.SCALAR,
1558             lazy_update=True,
1559         )
1560
1561     cached_results.update({})
1562     if prev_rank_state_sum is not None:
1563         return {
1564             "chunk_state": updated_states,
1565             "chunk_decay": chunk_decay,
1566             "prev_rank_state_sum": prev_rank_state_sum,
1567         }
1568     else:
1569         return {
1570             "chunk_state": updated_states,
1571             "chunk_decay": None,
1572             "prev_rank_state_sum": None,
1573         }
1574
1575 merge_mode_kernel_signature = (
1576     "*bf16:16, "

```

```

1566 551     "*bf16:16, "
1567 552     "*bf16:16, "
1568 553     "*fp32, "
1569 554     "*bf16:16, "
1570 555     "fp32, "
1571 556     "*bf16:16, "
1572 557     "i32:16, "
1573 558     "*bf16:16, "
1574 559     "i32:16, "
1575 560     "*bf16:16, "
1576 561     "*i32, "
1577 562     "i32, "
1578 563     "i32, "
1579 564     "*i32, "
1580 565     "i32, "
1581 566     "%FUSE_SP_STATE_UPDATE, "
1582 567     "%H, "
1583 568     "%K, "
1584 569     "%V, "
1585 570     "%CHUNK, "
1586 571     "%BLK_K, "
1587 572     "%BLK_V, "
1588 573     "%USE_TMA"
1589 574 )
1590 575
1591 576
1592 577 def get_merge_mode_kernel_info(B: int, T: int, H: int, K: int, V: int):
1593 578     """Static dispatcher for merge_mode_kernel. Auto-generated."""
1594 579     D = [K, V]
1595 580     key = (B, T, H, D)
1596 581     if key in (
1597 582         (1, 4096, 4, [128, 128]),
1598 583         (1, 131072, 32, [128, 128]),
1599 584         (1, 262144, 32, [128, 128]),
1600 585         (2, 1024, 8, [128, 128]),
1601 586         (2, 4096, 1, [128, 128]),
1602 587         (2, 8192, 1, [128, 128]),
1603 588     ):
1604 589         BLK_K = 128
1605 590         BLK_V = 128
1606 591         num_warps = 4
1607 592         num_ctas = 1
1608 593         num_stages = 3
1609 594         maxnreg = None
1610 595     elif key in (
1611 596         (1, 4096, 32, [128, 128]),
1612 597         (1, 16384, 8, [128, 128]),
1613 598         (1, 16384, 32, [128, 128]),
1614 599         (1, 32768, 8, [128, 128]),
1615 600         (1, 32768, 16, [128, 128]),
1616 601         (1, 32768, 32, [128, 128]),
1617 602         (1, 65536, 8, [128, 128]),
1618 603         (1, 65536, 32, [128, 128]),
1619 604         (1, 131072, 4, [128, 128]),
1620 605         (1, 131072, 8, [128, 128]),
1621 606         (1, 131072, 16, [128, 128]),
1622 607         (1, 262144, 4, [128, 128]),
1623 608         (1, 262144, 8, [128, 128]),
1624 609         (1, 262144, 16, [128, 128]),
1625 610         (2, 4096, 16, [128, 128]),
1626 611         (2, 4096, 32, [128, 128]),
1627 612         (2, 16384, 8, [128, 128]),
1628 613         (2, 16384, 32, [128, 128]),
1629 614         (2, 32768, 16, [128, 128]),
1630 615         (2, 32768, 32, [128, 128]),

```

```

1620         (2, 65536, 8, [128, 128]),
1621         (2, 65536, 16, [128, 128]),
1622         (2, 131072, 4, [128, 128]),
1623         (2, 131072, 8, [128, 128]),
1624         (2, 131072, 16, [128, 128]),
1625     ):
1626         BLK_K = 128
1627         BLK_V = 128
1628         num_warps = 8
1629         num_ctas = 1
1630         num_stages = 3
1631         maxnreg = None
1632     elif key in (
1633         (1, 1024, 16, [128, 128]),
1634         (1, 1024, 32, [128, 128]),
1635         (1, 2048, 8, [128, 128]),
1636         (1, 2048, 16, [128, 128]),
1637         (1, 16384, 1, [128, 128]),
1638         (1, 32768, 1, [128, 128]),
1639         (2, 2048, 4, [128, 128]),
1640         (2, 65536, 32, [128, 128]),
1641         (2, 131072, 32, [128, 128]),
1642     ):
1643         BLK_K = 128
1644         BLK_V = 128
1645         num_warps = 4
1646         num_ctas = 1
1647         num_stages = 4
1648         maxnreg = None
1649     elif key in (
1650         (1, 2048, 32, [128, 128]),
1651         (1, 4096, 16, [128, 128]),
1652         (1, 8192, 16, [128, 128]),
1653         (1, 8192, 32, [128, 128]),
1654         (1, 16384, 16, [128, 128]),
1655         (1, 65536, 4, [128, 128]),
1656         (1, 65536, 16, [128, 128]),
1657         (2, 1024, 32, [128, 128]),
1658         (2, 2048, 32, [128, 128]),
1659         (2, 8192, 8, [128, 128]),
1660         (2, 8192, 16, [128, 128]),
1661         (2, 8192, 32, [128, 128]),
1662         (2, 16384, 16, [128, 128]),
1663         (2, 32768, 4, [128, 128]),
1664         (2, 32768, 8, [128, 128]),
1665         (2, 65536, 4, [128, 128]),
1666     ):
1667         BLK_K = 128
1668         BLK_V = 128
1669         num_warps = 8
1670         num_ctas = 1
1671         num_stages = 4
1672         maxnreg = None
1673     elif key in (
1674         (1, 1024, 1, [64, 64]),
1675         (1, 4096, 1, [64, 64]),
1676     ):
1677         BLK_K = 128
1678         BLK_V = 32
1679         num_warps = 4
1680         num_ctas = 1
1681         num_stages = 3
1682         maxnreg = None
1683     elif key in (
1684         (1, 1024, 1, [128, 128]),

```

```

1674      (1, 1024, 4, [64, 64]),
1675      (1, 1024, 8, [64, 64]),
1676      (1, 2048, 1, [64, 64]),
1677      (1, 2048, 1, [128, 128]),
1678      (1, 2048, 4, [64, 64]),
1679      (2, 1024, 1, [64, 64]),
1680      (2, 2048, 1, [64, 64]),
1681  ):
1682      BLK_K = 128
1683      BLK_V = 32
1684      num_warps = 4
1685      num_ctas = 1
1686      num_stages = 4
1687      maxnreg = None
1688  elif key in (
1689      (1, 1024, 8, [128, 128]),
1690      (1, 1024, 32, [64, 64]),
1691      (1, 2048, 4, [128, 128]),
1692      (1, 4096, 4, [64, 64]),
1693      (1, 4096, 8, [64, 64]),
1694      (1, 4096, 8, [128, 128]),
1695      (1, 8192, 1, [128, 128]),
1696      (1, 8192, 8, [64, 64]),
1697      (1, 8192, 8, [128, 128]),
1698      (1, 8192, 32, [64, 64]),
1699      (1, 16384, 1, [64, 64]),
1700      (1, 16384, 32, [64, 64]),
1701      (1, 32768, 1, [64, 64]),
1702      (1, 32768, 4, [64, 64]),
1703      (1, 32768, 8, [64, 64]),
1704      (1, 32768, 16, [64, 64]),
1705      (1, 65536, 1, [128, 128]),
1706      (1, 131072, 1, [128, 128]),
1707      (1, 131072, 32, [64, 64]),
1708      (1, 262144, 1, [64, 64]),
1709      (1, 262144, 32, [64, 64]),
1710      (2, 1024, 4, [64, 64]),
1711      (2, 1024, 4, [128, 128]),
1712      (2, 1024, 8, [64, 64]),
1713      (2, 1024, 16, [64, 64]),
1714      (2, 1024, 32, [64, 64]),
1715      (2, 2048, 4, [64, 64]),
1716      (2, 2048, 8, [64, 64]),
1717      (2, 2048, 8, [128, 128]),
1718      (2, 2048, 16, [64, 64]),
1719      (2, 2048, 16, [128, 128]),
1720      (2, 2048, 32, [64, 64]),
1721      (2, 4096, 4, [64, 64]),
1722      (2, 4096, 4, [128, 128]),
1723      (2, 4096, 8, [64, 64]),
1724      (2, 4096, 8, [128, 128]),
1725      (2, 4096, 16, [64, 64]),
1726      (2, 4096, 32, [64, 64]),
1727      (2, 8192, 1, [64, 64]),
1728      (2, 8192, 4, [128, 128]),
1729      (2, 16384, 1, [64, 64]),
1730      (2, 16384, 1, [128, 128]),
1731      (2, 16384, 4, [128, 128]),
1732      (2, 16384, 8, [64, 64]),
1733      (2, 16384, 32, [64, 64]),
1734      (2, 32768, 16, [64, 64]),
1735      (2, 32768, 32, [64, 64]),
1736      (2, 65536, 1, [64, 64]),
1737      (2, 65536, 1, [128, 128]),
1738      (2, 65536, 16, [64, 64]),

```



```

1728 746 ) :
1729 747     BLK_K = 128
1730 748     BLK_V = 64
1731 749     num_warps = 4
1732 750     num_ctas = 1
1733 751     num_stages = 3
1734 752     maxnreg = None
1735 753     elif key in (
1736 754         (1, 4096, 1, [128, 128]),
1737 755     ):
1738 756         BLK_K = 128
1739 757         BLK_V = 64
1740 758         num_warps = 8
1741 759         num_ctas = 1
1742 760         num_stages = 3
1743 761         maxnreg = None
1744 762     elif key in (
1745 763         (1, 1024, 4, [128, 128]),
1746 764         (1, 1024, 16, [64, 64]),
1747 765         (1, 2048, 8, [64, 64]),
1748 766         (1, 2048, 16, [64, 64]),
1749 767         (1, 2048, 32, [64, 64]),
1750 768         (1, 4096, 16, [64, 64]),
1751 769         (1, 4096, 32, [64, 64]),
1752 770         (1, 8192, 1, [64, 64]),
1753 771         (1, 8192, 4, [64, 64]),
1754 772         (1, 8192, 4, [128, 128]),
1755 773         (1, 8192, 16, [64, 64]),
1756 774         (1, 16384, 4, [64, 64]),
1757 775         (1, 16384, 4, [128, 128]),
1758 776         (1, 16384, 8, [64, 64]),
1759 777         (1, 16384, 16, [64, 64]),
1760 778         (1, 32768, 4, [128, 128]),
1761 779         (1, 32768, 32, [64, 64]),
1762 780         (1, 65536, 1, [64, 64]),
1763 781         (1, 65536, 4, [64, 64]),
1764 782         (1, 65536, 8, [64, 64]),
1765 783         (1, 65536, 16, [64, 64]),
1766 784         (1, 65536, 32, [64, 64]),
1767 785         (1, 131072, 1, [64, 64]),
1768 786         (1, 131072, 4, [64, 64]),
1769 787         (1, 131072, 8, [64, 64]),
1770 788         (1, 131072, 16, [64, 64]),
1771 789         (1, 262144, 1, [128, 128]),
1772 790         (1, 262144, 4, [64, 64]),
1773 791         (1, 262144, 8, [64, 64]),
1774 792         (1, 262144, 16, [64, 64]),
1775 793         (2, 1024, 16, [128, 128]),
1776 794         (2, 8192, 4, [64, 64]),
1777 795         (2, 8192, 8, [64, 64]),
1778 796         (2, 8192, 16, [64, 64]),
1779 797         (2, 8192, 32, [64, 64]),
1780 798         (2, 16384, 4, [64, 64]),
1781 799         (2, 16384, 16, [64, 64]),
1782 800         (2, 32768, 1, [64, 64]),
1783 801         (2, 32768, 1, [128, 128]),
1784 802         (2, 32768, 4, [64, 64]),
1785 803         (2, 32768, 8, [64, 64]),
1786 804         (2, 65536, 4, [64, 64]),
1787 805         (2, 65536, 8, [64, 64]),
1788 806         (2, 65536, 32, [64, 64]),
1789 807         (2, 131072, 1, [64, 64]),
1790 808         (2, 131072, 1, [128, 128]),
1791 809         (2, 131072, 4, [64, 64]),
1792 810         (2, 131072, 8, [64, 64]),

```

```

1782 811         (2, 131072, 16, [64, 64]),
1783 812         (2, 131072, 32, [64, 64]),
1784 813     ):
1785 814         BLK_K = 128
1786 815         BLK_V = 64
1787 816         num_warps = 4
1788 817         num_ctas = 1
1789 818         num_stages = 4
1790 819         maxnreg = None
1791 820     elif key in (
1792 821         (2, 1024, 1, [128, 128]),
1793 822         (2, 2048, 1, [128, 128]),
1794 823         (2, 4096, 1, [64, 64]),
1795 824     ):
1796 825         BLK_K = 128
1797 826         BLK_V = 64
1798 827         num_warps = 8
1799 828         num_ctas = 1
1800 829         num_stages = 4
1801 830         maxnreg = None
1802 831     else:
1803 832         raise ValueError(f"Unsupported config for merge_mode_kernel: BTHD
1804 833         ={key}")
1805 834     return {
1806 835         "CHUNK": 64,
1807 836         "FUSE_SP_STATE_UPDATE": True,
1808 837         "H": H,
1809 838         "K": D[0],
1810 839         "V": D[1],
1811 840         "BLK_K": BLK_K,
1812 841         "BLK_V": BLK_V,
1813 842         "num_warps": num_warps,
1814 843         "num_stages": num_stages,
1815 844     }
1816 845
1817 846
1818 847
1819 848
1820 849 @aot_compile_spaces({
1821 850     "merge_mode_kernel": {
1822 851         "signature": merge_mode_kernel_signature,
1823 852         "grid": [("(%K + %BLK_K - 1) / %BLK_K) * ((%V + %BLK_V - 1) / %
1824 853         BLK_V)", "%H", "NC"],
1825 854         "triton_algos": [
1826 855             get_merge_mode_kernel_info(B, T, H, K, V)
1827 856             for B,T,H,(K,V) in[
1828 857                 (1, 1024, 8, [64, 64]),
1829 858                 (1, 2048, 1, [128, 128]),
1830 859                 (1, 2048, 16, [128, 128]),
1831 860                 (1, 2048, 4, [64, 64]),
1832 861                 (1, 2048, 8, [128, 128]),
1833 862                 (1, 262144, 32, [128, 128]),
1834 863                 (1, 32768, 1, [128, 128]),
1835 864                 (1, 32768, 4, [128, 128]),
1836 865                 (1, 4096, 1, [128, 128]),
1837 866                 (1, 4096, 1, [64, 64]),
1838 867                 (1, 4096, 4, [128, 128]),
1839 868                 (2, 1024, 16, [128, 128]),
1840 869                 (2, 1024, 8, [128, 128]),
1841 870                 (2, 131072, 1, [128, 128]),
1842 871                 (2, 131072, 1, [64, 64]),
1843 872                 (2, 131072, 16, [128, 128]),
1844 873                 (2, 131072, 16, [64, 64]),
1845 874                 (2, 131072, 32, [128, 128]),

```

```

1836         (2, 131072, 32, [64, 64]),
1837         (2, 131072, 4, [128, 128]),
1838         (2, 131072, 4, [64, 64]),
1839         (2, 131072, 8, [128, 128]),
1840         (2, 131072, 8, [64, 64]),
1841         (2, 16384, 16, [128, 128]),
1842         (2, 16384, 4, [128, 128]),
1843         (2, 16384, 8, [64, 64]),
1844         (2, 2048, 1, [128, 128]),
1845         (2, 2048, 1, [64, 64]),
1846         (2, 2048, 16, [128, 128]),
1847         (2, 2048, 4, [128, 128]),
1848         (2, 32768, 32, [128, 128]),
1849         (2, 32768, 32, [64, 64]),
1850         (2, 32768, 8, [128, 128]),
1851         (2, 4096, 1, [64, 64]),
1852         (2, 4096, 4, [64, 64]),
1853         (2, 4096, 8, [128, 128]),
1854         (2, 65536, 1, [128, 128]),
1855         (2, 65536, 1, [64, 64]),
1856         (2, 65536, 16, [64, 64]),
1857         (2, 65536, 4, [128, 128]),
1858         (2, 8192, 1, [128, 128]),
1859         (2, 8192, 32, [128, 128])
1860     ],
1861     ],
1862     })
1863 @triton.autotune(
1864     configs=[
1865         triton.Config({"BLK_K": BLK_K, "BLK_V": BLK_V, "USE_TMA": USE_TMA
1866     }, num_warps=num_warps, num_stages=num_stages)
1867     for num_warps in [4, 8]
1868     for num_stages in [3, 4]
1869     for BLK_K in [128]
1870     for BLK_V in [128, 64, 32]
1871     for USE_TMA in [True, False]
1872 ],
1873     key=[],
1874 )
1875 @triton.jit
1876 def merge_mode_kernel(
1877     q,
1878     k,
1879     v,
1880     g,
1881     chunk_state,
1882     scale,
1883     prev_rank_state_sum,
1884     stride_d0_ns_prev_rank_state_sum,
1885     chunk_decay,
1886     stride_d0_nc_with_ini_chunk_decay,
1887     out_0,
1888     cu_seqlens,
1889     NS,
1890     H_MUL_NS,
1891     chunk_indices,
1892     NC,
1893     FUSE_SP_STATE_UPDATE: tl.constexpr,
1894     H: tl.constexpr,
1895     K: tl.constexpr,
1896     V: tl.constexpr,
1897     CHUNK: tl.constexpr,
1898     BLK_K: tl.constexpr,
1899     BLK_V: tl.constexpr,

```

```

1890 938 USE_TMA: tl.constexpr,
1891 939 ):
1892 940     NUM_BLK_K = (K + BLK_K - 1) // BLK_K
1893 941     NUM_BLK_V = (V + BLK_V - 1) // BLK_V
1894 942     NUM_BLK_KV = NUM_BLK_K * NUM_BLK_V
1895 943
1896 944     i_v, i_h, i_gc = tl.program_id(0), tl.program_id(1), tl.program_id(2)
1897 945     i_s, i_c = tl.load(chunk_indices + i_gc * 2).to(tl.int32), tl.load(
1898 946     chunk_indices + i_gc * 2 + 1).to(tl.int32)
1899 947     bos, eos = tl.load(cu_seqlens + i_s).to(tl.int32), tl.load(cu_seqlens
1900 948     + i_s + 1).to(tl.int32)
1901 949     T = eos - bos
1902 950     # NC = tl.cdiv(T, CHUNK)
1903 951     # FIXME: set 'i_k' temporarily
1904 952     i_k = 0
1905 953     out_0 = out_0 + (bos * H + i_h) * V
1906 954
1907 955     if USE_TMA:
1908 956         q_desc = tl.make_tensor_descriptor(
1909 957             q + bos * H * K + i_h * K,
1910 958             shape=[T, K],
1911 959             strides=[H * K, 1],
1912 960             block_shape=[CHUNK, BLK_K],
1913 961         )
1914 962         k_desc = tl.make_tensor_descriptor(
1915 963             k + bos * H * K + i_h * K,
1916 964             shape=[T, K],
1917 965             strides=[H * K, 1],
1918 966             block_shape=[CHUNK, BLK_K],
1919 967         )
1920 968         v_desc = tl.make_tensor_descriptor(
1921 969             v + bos * H * V + i_h * V,
1922 970             shape=[T, V],
1923 971             strides=[H * V, 1],
1924 972             block_shape=[CHUNK, BLK_V],
1925 973         )
1926 974         chunk_state_desc = tl.make_tensor_descriptor(
1927 975             chunk_state + (i_gc + i_s).to(tl.int64) * H * K * V + i_h * K
1928 976             * V,
1929 977             shape=[K, V],
1930 978             strides=[V, 1],
1931 979             block_shape=[BLK_K, BLK_V],
1932 980         )
1933 981
1934 982     if FUSE_SP_STATE_UPDATE:
1935 983         prev_rank_state_sum_desc = tl.make_tensor_descriptor(
1936 984             prev_rank_state_sum + i_s * stride_d0_ns_prev_rank_state_sum
1937 985             + i_h * K * V,
1938 986             shape=[K, V],
1939 987             strides=[V, 1],
1940 988             block_shape=[BLK_K, BLK_V],
1941 989         )
1942 990
1943 991     # load 'q': (T, H, K,) => (CHUNK, BLK_K,)
1944 992     if not USE_TMA:
1945 993         cur_q = q + bos * H * K + i_h * K # fmt: skip
1946 994         ptr_q_0 = tl.make_block_ptr(cur_q, (T, K), (H * K, 1), (i_c *
1947 995         CHUNK, i_k * BLK_K), (CHUNK, BLK_K), (1, 0)) # fmt: skip
1948 996         blk_q_0 = tl.load(ptr_q_0, boundary_check=(1, 0)) # fmt: skip
1949 997     else:
1950 998         blk_q_0 = q_desc.load([i_c * CHUNK, i_k * BLK_K])
1951 999
1952 1000     # load (trans) 'k': (T, H, K,) => (BLK_K, CHUNK,)
1953 1001     if not USE_TMA:

```

```

1944 998      cur_k = k + bos * H * K + i_h * K # fmt: skip
1945 999      ptr_k_1 = tl.make_block_ptr(cur_k, (K, T), (1, H * K), (i_k *
1946      BLK_K, i_c * CHUNK), (BLK_K, CHUNK), (0, 1)) # fmt: skip
1947 1000      blk_k_1 = tl.load(ptr_k_1, boundary_check=(0, 1)) # fmt: skip
1948 1001      else:
1949 1002          blk_k_1 = k_desc.load([i_c * CHUNK, i_k * BLK_K]).trans()
1950 1003
1950 1004      # load 'v': (T, H, V) => (CHUNK, BLK_V,)
1951 1005      if not USE_TMA:
1952 1006          cur_v = v + bos * H * V + i_h * V # fmt: skip
1953 1007          ptr_v_2 = tl.make_block_ptr(cur_v, (T, V), (H * V, 1), (i_c *
1954      CHUNK, i_v * BLK_V), (CHUNK, BLK_V), (1, 0)) # fmt: skip
1955 1008          blk_v_2 = tl.load(ptr_v_2, boundary_check=(1, 0)) # fmt: skip
1956 1009      else:
1957 1010          blk_v_2 = v_desc.load([i_c * CHUNK, i_v * BLK_V])
1958 1011
1958 1012      # call_external_func: 'chunk_local_cumsum' to pre-compute g
1959 1013
1959 1014      # load 'g': (T, H,) => (CHUNK,)
1960 1015      cur_g = g + bos * H + i_h * 1 # fmt: skip
1961 1016      ptr_g_3 = tl.make_block_ptr(cur_g, (T), (H), (i_c * CHUNK), (CHUNK
1962      ,), (0,)) # fmt: skip
1963 1017      blk_g_3 = tl.load(ptr_g_3, boundary_check=(0,)) # fmt: skip
1964 1018      # load 'chunk_state': (NC_WITH_INI, H, K, V) => (BLK_K, BLK_V,)
1965 1019      if not USE_TMA:
1966 1020          cur_chunk_state = chunk_state + (i_gc + i_s).to(tl.int64) * H * K
1967 1021          ptr_chunk_state_4 = tl.make_block_ptr(cur_chunk_state, (K, V), (
1968      V, 1), (i_k * BLK_K, i_v * BLK_V), (BLK_K, BLK_V), (1, 0)) # fmt
1969      : skip
1970 1022          blk_chunk_state_4 = tl.load(ptr_chunk_state_4, boundary_check=(1,
1971      0,)) # fmt: skip
1972 1023      else:
1973 1024          blk_chunk_state_4 = chunk_state_desc.load([i_k * BLK_K, i_v *
1974      BLK_V])
1975 1025
1975 1026      # comments not available for op 'if_beg'
1976 1027      if FUSE_SP_STATE_UPDATE:
1977 1028          # load 'prev_rank_state_sum': (NS, H, K, V) => (BLK_K, BLK_V,)
1978 1029          if not USE_TMA:
1979 1030              cur_prev_rank_state_sum = prev_rank_state_sum + i_s *
1980 1031              stride_d0_ns_prev_rank_state_sum + i_h * K * V # fmt: skip
1981 1032              ptr_prev_rank_state_sum_5 = tl.make_block_ptr(
1982      cur_prev_rank_state_sum, (K, V), (V, 1), (i_k * BLK_K, i_v * BLK_V
1983      ,), (BLK_K, BLK_V), (1, 0)) # fmt: skip
1984 1033              blk_prev_rank_state_sum_5 = tl.load(ptr_prev_rank_state_sum_5
1985      , boundary_check=(1, 0)) # fmt: skip
1986 1034          else:
1987 1035              blk_prev_rank_state_sum_5 = prev_rank_state_sum_desc.load([
1988      i_k * BLK_K, i_v * BLK_V])
1989 1036
1989 1037          # load 'chunk_decay': (NC_WITH_INI, H,) => (,)
1990 1038          cur_chunk_decay = chunk_decay + (i_gc + i_s).to(tl.int64) *
1991 1039          stride_d0_nc_with_ini_chunk_decay + i_h * 1 # fmt: skip
1992 1040          ptr_chunk_decay_6 = cur_chunk_decay # fmt: skip
1993 1041          blk_chunk_decay_6 = tl.load(ptr_chunk_decay_6)
1994 1042          # to: ('BLK_K', 'BLK_V') => ('BLK_K', 'BLK_V')
1995 1043          blk_prev_rank_state_sum_5_float32_7 = blk_prev_rank_state_sum_5.
1996 1044          to(tl.float32)
1997 1045          # mul: (,), ('BLK_K', 'BLK_V') => ('BLK_K', 'BLK_V')
1998          mul_tensor_8 = blk_chunk_decay_6 *
1999          blk_prev_rank_state_sum_5_float32_7
2000          # to: ('BLK_K', 'BLK_V') => ('BLK_K', 'BLK_V')
2001          mul_tensor_8_bfloat16_9 = mul_tensor_8.to(tl.bfloat16)

```

```

1998 1046 # add: ('BLK_K', 'BLK_V'), ('BLK_K', 'BLK_V') => ('BLK_K', 'BLK_V'
1999 ')
2000 1047 add_tensor_10 = mul_tensor_8_bfloat16_9 + blk_chunk_state_4
2001 1048 # comments not available for op 'bind_var'
2002 1049 blk_chunk_state_4 = add_tensor_10
2003 1050 # comments not available for op 'end_if'
2004 1051 # to: ('BLK_K', 'BLK_V') => ('BLK_K', 'BLK_V')
2005 1052 blk_chunk_state_4_bfloat16_11 = blk_chunk_state_4.to(tl.bfloat16)
2006 1053 # matmul: ('CHUNK', 'BLK_K'), ('CHUNK', 'BLK_K') => ('CHUNK', 'CHUNK
2007 ')
2008 1054 matmul_12 = tl.dot(blk_q_0, blk_k_1).to(blk_q_0.dtype)
2009 1055 # tril: ('CHUNK', 'CHUNK') => ('CHUNK', 'CHUNK')
2010 1056 tril_13 = tl.where(tl.arange(0, CHUNK)[:], None] >= tl.arange(0, CHUNK
2011 ) [None, :], matmul_12, 0)
2012 1057 # unsqueeze: ('CHUNK',) => ('CHUNK', 1)
2013 1058 unsqueeze_14 = blk_g_3[:, None]
2014 1059 # unsqueeze: ('CHUNK',) => (1, 'CHUNK')
2015 1060 unsqueeze_1_15 = blk_g_3[None, :]
2016 1061 # sub: ('CHUNK', 1), (1, 'CHUNK') => ('CHUNK', 'CHUNK')
2017 1062 sub_16 = unsqueeze_14 - unsqueeze_1_15
2018 1063 # exp: ('CHUNK', 'CHUNK') => ('CHUNK', 'CHUNK')
2019 1064 exp_17 = tl.exp(sub_16)
2020 1065 # mul: ('CHUNK', 'CHUNK'), ('CHUNK', 'CHUNK') => ('CHUNK', 'CHUNK')
2021 1066 mul_18 = tril_13 * exp_17
2022 1067 # to: ('CHUNK', 'CHUNK') => ('CHUNK', 'CHUNK')
2023 1068 mul_18_bfloat16_19 = mul_18.to(tl.bfloat16)
2024 1069 # matmul: ('CHUNK', 'CHUNK'), ('CHUNK', 'BLK_V') => ('CHUNK', 'BLK_V
2025 ')
2026 1070 matmul_1_20 = tl.dot(mul_18_bfloat16_19, blk_v_2).to(
2027 mul_18_bfloat16_19.dtype)
2028 1071 # matmul: ('CHUNK', 'BLK_K'), ('BLK_K', 'BLK_V') => ('CHUNK', 'BLK_V
2029 ')
2030 1072 matmul_2_21 = tl.dot(blk_q_0, blk_chunk_state_4_bfloat16_11).to(
2031 blk_q_0.dtype)
2032 1073 # exp: ('CHUNK',) => ('CHUNK',)
2033 1074 exp_1_22 = tl.exp(blk_g_3)
2034 1075 # unsqueeze: ('CHUNK',) => ('CHUNK', 1)
2035 1076 unsqueeze_2_23 = exp_1_22[:, None]
2036 1077 # mul: ('CHUNK', 'BLK_V'), ('CHUNK', 1) => ('CHUNK', 'BLK_V')
2037 1078 mul_1_24 = matmul_2_21 * unsqueeze_2_23
2038 1079 # add: ('CHUNK', 'BLK_V'), ('CHUNK', 'BLK_V') => ('CHUNK', 'BLK_V')
2039 1080 add_25 = matmul_1_20 + mul_1_24
2040 1081 # mul: ('CHUNK', 'BLK_V'), () => ('CHUNK', 'BLK_V')
2041 1082 mul_2_26 = add_25 * scale
2042 1083 # store => ('CHUNK', 'BLK_V')
2043 1084 # assume output layout: [T, H, V]
2044 1085 out_0_ty = out_0.dtype.element_ty
2045 1086 ptr_out_0 = tl.make_block_ptr(out_0, (T, V,), (H * V, 1, ), (i_c *
2046 CHUNK, i_v * BLK_V, ), (CHUNK, BLK_V, ), (1, 0)) # fmt: skip
2047 1087 tl.store(ptr_out_0, mul_2_26.to(out_0_ty), boundary_check=(1, 0))
2048 1088
2049 1089 def launch_merge_mode(
2050 1090 chunk_decay: torch.Tensor,
2051 1091 v: torch.Tensor,
2052 1092 k: torch.Tensor,
2053 1093 q: torch.Tensor,
2054 1094 g: torch.Tensor,
2055 1095 prev_rank_state_sum: torch.Tensor,
2056 1096 chunk_state: torch.Tensor,
2057 1097 scale,
2058 1098 cu_seqlens: torch.IntTensor,
2059 1099 cached_results: dict,
2060 1100 fuse_sp_update: bool,
2061 1101 ) -> torch.Tensor:
2062 1102 T, H, K, V = *k.shape, v.shape[-1]

```

```

2052 1103
2053 1104     if scale is None:
2054 1105         scale = k.shape[-1] ** -0.5
2055 1106
2056 1107     # TODO: now fix chunk_size to 64
2057 1108     CHUNK = 64
2058 1109     chunk_indices = prepare_chunk_indices(cu_seqlens, CHUNK)
2059 1110     chunk_offsets = prepare_chunk_offsets(cu_seqlens, CHUNK)
2060 1111     chunk_offsets_ini = prepare_chunk_offsets_with_ini(cu_seqlens, CHUNK)
2061 1112     NS, NC = len(cu_seqlens) - 1, chunk_indices.shape[0]
2062 1113     NC_WITH_INI = NC + NS
2063 1114     out = torch.empty_like(v)
2064 1115     use_aot = os.environ.get('FORGE_USE_AOT', '0')
2065 1116     FORGE_USE_AOT = True if use_aot.lower() in ['1', 'true', 'yes'] else
2066 1117     None
2067 1118
2068 1119     def alloc_fn(size: int, alignment: int, stream: int):
2069 1120         return torch.empty(size, device='cuda', dtype=torch.int8)
2070 1121
2071 1122     triton.set_allocator(alloc_fn)
2072 1123
2073 1124     def grid(meta):
2074 1125         BLK_K = meta['BLK_K']
2075 1126         BLK_V = meta['BLK_V']
2076 1127         NUM_BLK_K = (K + BLK_K - 1) // BLK_K
2077 1128         NUM_BLK_V = (V + BLK_V - 1) // BLK_V
2078 1129         NUM_BLK_KV = NUM_BLK_K * NUM_BLK_V
2079 1130         H_MUL_NS = H * NS
2080 1131         return (NUM_BLK_KV, H, NC,)
2081 1132
2082 1133     if FORGE_USE_AOT is None:
2083 1134         merge_mode_kernel[grid](
2084 1135             chunk_decay=chunk_decay,
2085 1136             stride_d0_nc_with_ini_chunk_decay=H,
2086 1137             v=v,
2087 1138             k=k,
2088 1139             q=q,
2089 1140             g=g,
2090 1141             prev_rank_state_sum=prev_rank_state_sum,
2091 1142             stride_d0_ns_prev_rank_state_sum=H*K*V,
2092 1143             chunk_state=chunk_state,
2093 1144             scale=scale,
2094 1145             cu_seqlens=cu_seqlens,
2095 1146             H=H,
2096 1147             NS=NS,
2097 1148             K=K,
2098 1149             V=V,
2099 1150             CHUNK=CHUNK,
2100 1151             H_MUL_NS=H * NS,
2101 1152             out_0=out,
2102 1153             chunk_indices=chunk_indices,
2103 1154             NC=NC,
2104 1155             FUSE_SP_STATE_UPDATE=fuse_sp_update,
2105 1156         )
2106 1157
2107 1158     else:
2108 1159         from forge.aot_utils import forge_aot_ops
2109 1160
2110 1161         algo_info = forge_aot_ops.merge_mode_kernel__triton_algo_info_t()
2111 1162         for _k, _v in get_merge_mode_kernel_info(NS, T, H, K, V).items():
2112 1163             setattr(algo_info, _k, _v)
2113 1164         forge_aot_ops.merge_mode_kernel(
2114 1165             0, # torch.cuda.current_stream().cuda_stream,
2115 1166             q.data_ptr(), # q

```

```

2106      k.data_ptr(), # k
2107      v.data_ptr(), # v
2108      g.data_ptr(), # g
2109      chunk_state.data_ptr(), # chunk_state
2110      scale, # scale
2111      prev_rank_state_sum.data_ptr() if prev_rank_state_sum else 0,
2112      # prev_rank_state_sum
2113      H*K*V, # stride_d0_ns_prev_rank_state_sum
2114      chunk_decay.data_ptr() if chunk_decay else 0, # chunk_decay
2115      H, # stride_d0_nc_with_ini_chunk_decay
2116      out.data_ptr(), # out_0
2117      cu_seqlens.data_ptr(), # cu_seqlens
2118      NS, # NS
2119      H * NS, # H_MUL_NS
2120      chunk_indices.data_ptr(), # chunk_indices
2121      NC, # NC
2122      algo_info,
2123  )
2124  cached_results.update({})
2125  return out
2126
2127  def fused_op():
2128      cached_results = {}
2129      # AUTO: precompute decay outside
2130      g_cumsum = chunk_local_cumsum(g[None, ...], 64, cu_seqlens=cu_seqlens)
2131      ).squeeze(0)
2132      updated_states = launch_fuse_chunk_decay(
2133          g=g_cumsum,
2134          prev_s=prev_s,
2135          v=v,
2136          k=k,
2137          cu_seqlens=cu_seqlens,
2138          cached_results=cached_results,
2139          dist_scan=dist_scan,
2140          lazy_update=lazy_update,
2141      )
2142      chunk_state = updated_states['chunk_state']
2143      prev_rank_state_sum = updated_states.get('prev_rank_state_sum', None)
2144      chunk_decay = updated_states.get('chunk_decay', None)
2145      o = launch_merge_mode(
2146          chunk_decay=chunk_decay,
2147          v=v,
2148          k=k,
2149          q=q,
2150          g=g_cumsum,
2151          prev_rank_state_sum=prev_rank_state_sum,
2152          chunk_state=chunk_state,
2153          scale=scale,
2154          cu_seqlens=cu_seqlens,
2155          cached_results=cached_results,
2156          fuse_sp_update=fuse_sp_update,
2157      )
2158      return o, chunk_state

```

Listing 2: Generated Scalar GLA kernel by Forge

## G.2 IMPLEMENTATION FOR DELTANET

```

2156 1 def chunk_mode_deltanet(k: Tensor, v: Tensor, b: Tensor) -> Tensor:
2157 2     I = forge.identity(k, k.size(0))
2158 3     # Note: forge handles the specific fp32 accumulation
2159 4     T = (I + forge.matmul_out_fp32(k * b[None, None], k.T).tril(-1)).
        inverse().to(k.dtype)

```



```

2160 5     U = T @ (v * b[..., None])
2161 6     W = T @ (k * b[..., None])
2162 7     # Explicitly cache results for reuse in other phases
2163 8     forge.cache_result(U, "u")
2164 9     forge.cache_result(W, "w")
2165 10    S = k.T @ U
2166 11    return S
2167 12
2167 13 def decay_mode_deltanet(prev_s: Tensor, k: Tensor, w: Tensor, chunk_state
2168 : Tensor) -> Tensor:
2169 14     """
2170 15     Inter-Chunk State Propagation
2171 16     """
2172 17     # Calculate decay matrix based on cached W
2173 18     decay = forge.identity(k, k.size(1)) - k.T @ w
2174 19     return decay @ prev_s.to(decay.dtype) + chunk_state
2175 20
2175 21 def merge_mode_deltanet(q: Tensor, k: Tensor, u: Tensor, w: Tensor,
2176 chunk_state: Tensor, scale: Tensor) -> Tensor:
2177 22     """
2178 23     Output Merging
2179 24     Note: 'u' and 'w' are automatically injected from the cached results
2180 25     """
2181 26     new_v = u - w @ chunk_state
2182 27     return ((q @ k.T).tril(0) @ new_v + q @ chunk_state) * scale

```

Listing 3: Forge Implementation of DeltaNet with Intermediate Result Caching