

FLASHRNN: OPTIMIZING TRADITIONAL RNNs ON MODERN HARDWARE

Anonymous authors

Paper under double-blind review

ABSTRACT

While Transformers and other sequence-parallelizable neural network architectures seem like the current state of the art in sequence modeling, they specifically lack state-tracking capabilities. These are important for time-series tasks and logical reasoning. Traditional RNNs like LSTMs and GRUs, as well as modern variants like sLSTM do have these capabilities at the cost of strictly sequential processing. While this is often seen as a strong limitation, we show how fast these networks can get with our hardware-optimization FlashRNN in Triton and CUDA, optimizing kernels to the register level on modern GPUs. We extend traditional RNNs with a parallelization variant that processes multiple RNNs of smaller hidden state in parallel, similar to the head-wise processing in Transformers. To enable flexibility on different GPU variants, we introduce a new optimization framework for hardware-internal cache sizes, memory and compute handling. It models the hardware in a setting using polyhedral-like constraints, including the notion of divisibility. This speeds up the solution process in our ConstrINT library for general integer constraint satisfaction problems (integer CSPs). We show that our kernels can achieve 50x speed-ups over a vanilla PyTorch implementation and allow 40x larger hidden sizes compared to our Triton implementation. We will open-source our kernels and the optimization library to boost research in the direction of state-tracking enabled RNNs and sequence modeling.

1 INTRODUCTION

Sequence models are at the core of many applications like time-series modeling, natural language processing, text, audio and video models, and predictions for physical systems based on ODEs or PDEs. Vaswani et al. (2017); Degraive et al. (2022); Nearing et al. (2024) While there are modern sequence-parallelizable architectures like the Transformer (Vaswani et al., 2017), Mamba (Gu & Dao, 2023) or mLSTM (Beck et al., 2024), these lack the state-tracking capabilities (Merrill et al., 2024; Merrill & Sabharwal, 2023; Delétang et al., 2023) of traditional RNNs like LSTM (Hochreiter & Schmidhuber, 1997), GRU (Cho et al., 2014) and sLSTM (Beck et al., 2024).

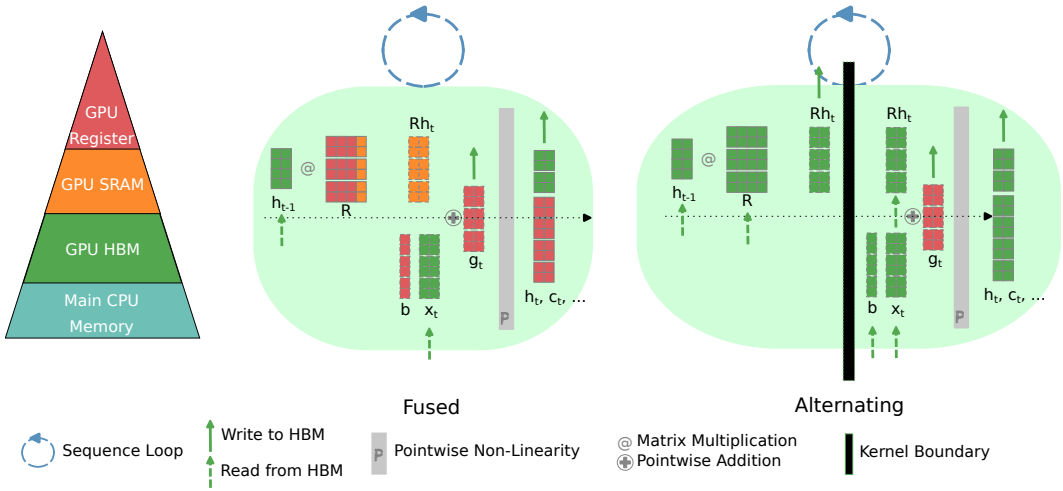
Traditional RNNs include a recurrent connection or memory mixing, that connects the previous hidden state in a non-linear way to the current state update and this way mixes the states of different memory cells. While the sequence has to be processed step by step, computed hidden states and the recurrent matrix weights can stay cached, enabling a large speed optimization. In this work, we introduce FlashRNN as a generic hardware-optimized library for these RNN-style architectures.

Our library facilitates research in the direction of state-tracking enabled RNN architectures, in two ways: Firstly, it enables easier and more efficient use of recent RNN-architectures like sLSTM (Beck et al., 2024). This includes the notion of block-diagonal recurrent matrices that can speed up networks while lowering the number of parameters. Secondly, it can be easily extended to novel RNN-like architecture variants, as it supports generic state and gate numbers per cell. The LSTM (Hochreiter & Schmidhuber, 1997; F.A. Gers, 1999), with its two states and four gates (we consider the cell update as a fourth "gate" for simplicity here), can be implemented as easy as a simple Elman-RNN with one state and one gate (Elman, 1990), or sLSTM with its three states and four gates Beck et al. (2024).

To realize the shown speed-ups, we fuse the recurrent matrix-multiplication part with the point-wise activation part, both wrapped in the sequential loop into one kernel. This can be used on different

054 GPUs and with different state/gate variants, as our library optimizes internal memory sizes and
 055 operations automatically based on the models' hidden sizes and the cache and register sizes of the
 056 hardware.

057 For the auto-optimization we introduce an integer constraint satisfaction library ConstrINT. With
 058 this library, one can model generic integer CSP problems with equality, inequality and divisibil-
 059 ity constraints as these can model size constraints on modern hardware with specific tensor-core,
 060 register and SRAM memory sizes.
 061



072
 073
 074
 075
 076
 077
 078
 079 Figure 1: FlashRNN Kernel overview: Left: Basic Memory Hierarchy in modern GPUs. Center:
 080 Fused Kernel (forward) leveraging all caching options for maximal speed. Right: Alternating Ker-
 081 nels (forward) for maximum hidden sizes, with two kernel calls per time step. The colors show the
 082 caching level of the different tensors, the batch dimension is depicted to the right (except for R), the
 083 hidden / gate dimension vertically.

084
 085
 086 **2 RELATED WORK**

087 Hardware-aware algorithms and their open-source implementations of common sequence modeling
 088 primitives have been focused primarily around the Transformer architecture (Vaswani et al., 2017)
 089 and its attention operation because of its ubiquity in language modeling. FlashAttention (Dao et al.,
 090 2022) introduced an IO-aware attention algorithm and CUDA implementation that uses tiling to re-
 091 duce the number of memory reads/writes between GPU high bandwidth memory (HBM) and GPU
 092 on-chip SRAM, and achieves significant memory savings. FlashAttention2 (Dao, 2024) improves
 093 FlashAttention with better work partitioning and the additional parallelization over the sequence di-
 094 mension. FlashAttention3 (Shah et al., 2024) takes advantage of new capabilities, such as asynchrony
 095 and FP8 low precision support of the recent NVIDIA Hopper GPU generation.

096 Recently, novel sequence models taking inspiration of Linear Attention (Katharopoulos et al., 2020)
 097 have shown promising performance compared to Transformer Attention (Beck et al., 2024; Yang
 098 et al., 2024; Dao & Gu, 2024). Yang et al. (2024) provide an hardware-efficient algorithm and
 099 implementation in Triton for Gated Linear Attention that trades off memory movement against par-
 100 allelizability and show that it is faster than FlashAttention2.

101 Traditional RNNs like LSTMs (Hochreiter & Schmidhuber, 1997) or GRUs (Cho et al., 2014) are
 102 still widely used in many applications, such as for example time series modeling or reinforcement
 103 learning (Nearing et al., 2024; Degraeve et al., 2022). Many of these applications rely on optimized
 104 closed-source implementations of these RNN operations such as in the NVIDIA cuDNN¹ library,
 105 which is integrated in PyTorch. Sharvil (2020) provide an open-source alternative in CUDA for
 106 specific LSTM and GRU variants in their HASTE library, which served as inspiration for this work.
 107

¹<https://developer.nvidia.com/cudnn>

HASTE is limited in speed due to a sequence of alternating calls of matrix multiplication and point-wise kernels.

Our work FlashRNN overcomes this limitation by fusing the recurrent matrix multiplication with the pointwise operations into a single persistent kernel with custom caching of the recurrent weights in registers. FlashRNN also supports the bfloat16 dtype and block-diagonal recurrent matrices. By open-sourcing our CUDA and Triton kernels we aim to enable researchers to quickly reach similar speeds compared to optimized closed source libraries.

3 GENERIC RECURRENT NEURAL NETWORK ARCHITECTURE WITH MEMORY MIXING

A generic RNN architecture that we aim to optimize has N_s states $\mathbf{s}^{(i)} \in \mathbb{R}^d$, and N_g gates (or pre-activations) $\mathbf{g}^{(j)} \in \mathbb{R}^d$, with d being the embedding dimension or hidden size of the RNN. For example the LSTM (Hochreiter & Schmidhuber, 1997) has $N_s = 2$ states and $N_g = 4$ gates.

Each gate receives an input $\mathbf{x}^{(j)} \in \mathbb{R}^d$. As learnable parameters the gates have a recurrent matrix $\mathbf{R}^{(j)} \in \mathbb{R}^{d \times d}$ that models the dependency on the previous hidden state $\mathbf{s}_{t-1}^{(0)}$ and a bias $\mathbf{b}^{(j)} \in \mathbb{R}^d$. The state sequence of the RNN is then defined as:

$$\mathbf{g}_t^{(j)} = \mathbf{x}_t^{(j)} + \mathbf{R}^{(j)} \mathbf{s}_{t-1}^{(0)} + \mathbf{b}^{(j)}, \quad (1)$$

$$\mathbf{s}_t^{(i)} = \mathcal{P}^{(i)} \left(\{\mathbf{s}_{t-1}^{(k)}\}_k, \{\mathbf{g}_t^{(j)}\}_j \right), \quad (2)$$

with a point-wise / element-wise function $\mathcal{P}^{(i)}$ that does not mix different cells along the vector dimension (unlike the recurrent weight). In Appendix A we show how this generic formulation translates to the most common RNN variants.

Usually for these networks, the input is modified with another weight matrix \mathbf{W} . We omit this here as it can be moved outside of the basic kernels. In the common training setting, where the whole sequence is given as input, the weight matrix \mathbf{W} can be applied in parallel to all timesteps before processing a sequence in the RNN. Our runtime experiments in Section 6.1 show that this operation has only marginal impact on the overall runtime.

4 GENERIC GRADIENT FOR BACKPROPAGATION THROUGH TIME

In back-propagation through time (Mozer, 1995), the backward pass of this RNN architecture can be recursively defined as well. The backward pass reads:

$$\delta \mathbf{g}_t^{(j)} = \frac{\partial \mathcal{P}^{(l)} \left(\{\mathbf{s}_{t-1}^{(k)}\}_k, \{\mathbf{g}_t^{(j)}\}_j \right)}{\partial \mathbf{g}_t^{(j)}} \delta \mathbf{s}_t^{(l)} \quad (3)$$

$$\delta \mathbf{s}_{t-1}^{(i)} = \frac{\partial \mathcal{P}^{(l)} \left(\{\mathbf{s}_{t-1}^{(k)}\}_k, \{\mathbf{g}_t^{(j)}\}_j \right)}{\partial \mathbf{s}_{t-1}^{(i)}} \delta \mathbf{s}_t^{(l)} + \left(\mathbf{R}^{(j)T} \delta \mathbf{g}_{t-1}^{(j)} \quad \text{if } i = 0 \right) \quad (4)$$

The structure of the gradient shows that, also for the backward pass, we have an alternation of point-wise operations (left) and matrix multiplication (right).

The input gradient is equal to the gate gradients, the bias gradient is the sum of the input gradients and the recurrent weight matrix gradient is the time-wise sum of the outer product of gate gradients with the state values:

$$\delta \mathbf{x}_t^{(j)} = \delta \mathbf{g}_t^{(j)} \quad (5)$$

$$\delta \mathbf{b}^{(j)} = \sum_t \delta \mathbf{g}_t^{(j)} \quad (6)$$

$$\delta \mathbf{R}^{(j)} = \sum_t \delta \mathbf{g}_t^{(j)} \mathbf{s}_t^{(0)T} \quad (7)$$

162 4.1 VANISHING AND EXPLODING GRADIENTS AND GRADIENT MODIFICATIONS

163
164 For a neural network to be stably trainable, there must not be exploding gradients, also vanishing
165 gradients should be prohibited for long context sequence modeling (Hochreiter & Schmidhuber,
166 1997). Still, for the generic structure of Equations 3, there can be exploding components: Firstly,
167 one or more eigenvalues of the point-wise function Jacobian can be greater than one in magnitude.
168 This can be mitigated by a proper choice of the point-wise function. Secondly, the combination
169 of recurrent matrix and gate gradients with the gradient $\frac{\partial \mathcal{P}^{(0)}}{\partial \mathbf{g}_t^{(j)}} \mathbf{R}^{(j)T}$ could have singular values of
170 magnitude > 1 . This case cannot be excluded directly, as the recurrent matrix consists of train-
171 able weights with usually unconstrained magnitude. However, for practical training this is rarely a
172 limitation.

173 In our library, we implement a simple approach for mitigating this at the cost of additional gradient
174 noise, clipping the gradient values on a scalar level after each time step. Specifically, we clip the
175 term containing the recurrent matrix to within a pre-defined magnitude. The gradients can even be
176 cut to zero, leading to typically worse convergence at the benefit of faster training, as the recurrent
177 matrix part in Equation 3 is cut to zero for the backward pass.

179 4.2 HEAD-WISE PARALLELIZATION

180
181 When increasing the size of a neural network, typically the width, i.e. the embedding dimension or
182 hidden size is increased. Vaswani et al. (2017) found that for the attention operation it is beneficial
183 to linearly project the input embedding vectors into multiple smaller input vectors, the so called
184 heads, and then perform attention on each of these small vectors in parallel. This parallelization
185 primitive enables also efficient implementations on GPUs, since each head can be computed on
186 different streaming multiprocessors of the GPU (Dao et al., 2022) in parallel (see also Section 5.1).

187 Many more recent architectures also rely on this head-wise parallelization primitive (Beck et al.,
188 2024; Yang et al., 2024; Dao & Gu, 2024), where the embedding or hidden vector of dimension
189 d is split into N_{heads} heads of smaller dimension $d_{head} = d/N_{heads}$, each of which is processed
190 independently inside the sequential part. In FlashRNN, we apply this primitive to traditional RNNs
191 by dividing the recurrent matrix \mathbf{R} into multiple blocks or heads $\mathbf{R}_{head} \in \mathbb{R}^{d_{head} \times d_{head}}$ rendering
192 the recurrent matrix \mathbf{R} as a block-diagonal matrix.

194 5 HARDWARE-EFFICIENT IMPLEMENTATION

196 5.1 GPU-ACCELERATED COMPUTING

197
198 Modern compute hardware in the form of GPUs enables massive parallelization and accelerated
199 matrix multiplication. This means that both point-wise (scalar) operations can be parallelized and
200 matrix multiplications have good support via BLAS-like libraries (Lawson et al. (1979); Thakkar
201 et al. (2023)), as used for RNN training workloads as defined above.

202 **Execution Model** Specifically, a modern GPU consists of larger computational super-units (i.e.
203 streaming multiprocessors (SMs)) that have some faster memory attached to them. There are three
204 levels of memory, the large HBM which allows global random access from all computational units
205 at the cost of low speed (still fast compared to CPU RAM access), the SRAM which is attached to
206 one computational super-unit and the registers which are tied to a smallest computational unit (i.e.
207 thread). One super-unit usually supports up to 1024 threads in parallel (with varying register sizes)
208 which are typically referred to as a block, multiple blocks executed in parallel on multiple super-units
209 are called the grid.² A NVIDIA H100, for example, consists of 132 streaming multiprocessors, with
210 256 KB SRAM per SM and a SRAM bandwidth of around 33 TB/s (Spector et al.), compared to the
211 up to 3 TB/s for access to the 80 GB of HBM. Starting from the NVIDIA Ampere Architecture and
212 newer, there is hardware acceleration for asynchronous loading and SRAM interconnection, which
213 we did not utilize in this work.³ Beyond the memory levels, a computational super-unit allows for
214 hardware-accelerated matrix multiplication (e.g. via TensorCores, "wmma" operation). Typically it

215 ²https://docs.nvidia.com/cuda/pdf/CUDA_C_Programming_Guide.pdf

³<https://resources.nvidia.com/en-us-tensor-core/gtc22-whitepaper-hopper>

is divided into sub-units (warps) of a certain number of threads (NVIDIA: 32) that act as one for a matrix multiplication. There are certain size limitations for this acceleration which have to be considered in the kernel optimization process. For a NVIDIA H100, this means that only minimal matrices of sizes $32 \times 16 \times 8$, $16 \times 16 \times 16$ or $8 \times 16 \times 32$ can be multiplied for the low-precision bfloat16 or float16 dtypes, larger matrix multiplications have to be composed of those, by parallelization along the outer dimensions and summation along the accumulating dimension.

Performance measures The specific limitation of a computational load falls into two regimes: Being compute-bound or being memory-bound. In the former case, the arithmetic intensity is high, there are many compute operations per loaded byte and therefore the main limitation is the computational part. In the latter case, arithmetic intensity is low and the bottleneck is the memory access to load inputs and store outputs (Williams et al., 2009). Small operations, like applying an activation function in parallel are typically memory bound and should be grouped together into a fused kernel.

Fused Kernels To minimize HBM memory accesses, one combines multiple arithmetic operations in one GPU kernel. A kernel is a set of instructions on the GPU which is executed in parallel on its parts. Only within the execution of one kernel SRAM and registers are kept and can serve as a cache. Therefore, for memory-bound operations it is helpful to fuse multiple arithmetic operations into one kernel to leverage these lower cache levels. While compilers can fuse point-wise operations, an alternation of both point-wise computations and matrix multiplication is non-trivial.

Algorithm 1 FlashRNN-fused forward pass

All states are tiled along threads (single ALU) in Warps (for e.g. Matrix Multiplication) in a block (SRAM level, streaming multiprocessor) and blocks in the grid (multiple streaming multiprocessors) - additionally there can be looping levels where the parallelization is resolved to a simple loop.

Require: Recurrent matrix R_{gs} , inputs x_{tg} , biases b_{tg}

Require: Initial state s_{0bs}

Load R_{tbgs} , b_{tg} to registers and SRAM

for i_{batch} in *BatchLoop* **do**

Load s_{0bs} to registers

for $t \in 0..T - 1$ **do**

for Matrix Tiles in Registers **do**

Calculate and Accumulate Matrix product $y_{tbg} = R_{gs} s_{tbs0}$ along s

end for

if more matrix tiles in SRAM **then**

for Matrix Tiles in SRAM **do**

Load Matrix Tile

Calculate and Accumulate Matrix product $y_{tbg} = R_{gs} s_{tbs0}$ along s

end for

end if

Accumulate MatMul results along s in shared memory (Write, Load and Sum)

if state dimension too big for SRAM **then**

Accumulate MatMul results along s in HBM (Write, Grid Sync, Load, Sum)

end if

Sum Gate inputs x_{tbg} with y_{tbg} and biases b_g to gates g_{tbg}

Compute Point-wise Function $s_{t+1bs'} = \mathcal{P}(s_{tbs'}, g_{tbg})$ with aligned states s' and gates g

Write out gates for backward pass and new states to HBM

Grid-Level Sync (for new states to be available across the whole grid)

end for

end for

5.2 FLASHRNN KERNELS

As the RNN operations of Equations 1 and 3 are a sequential alternation between matrix multiplication and pointwise non-linearities, there is a simple speed up variant that optimizes these two primitives separately. Our library implements this variant, in the **alternating backend**. This enables arbitrarily large head dimensions (to the limits of HBM GPU memory). Also, a vanilla PyTorch

270 implementation relying on auto-grads will work in this primitive, but for every time step a separate
 271 state is saved for the backward pass, leading to inefficiencies beyond memory accesses. We show
 272 that moving the time-loop into CUDA can already give large speedups over the vanilla PyTorch
 273 implementation.

274 The downside of the **alternating** implementation is that there are no I/O optimizations beyond a
 275 single time step. For every time step, the current input and last state, as well as the recurrent matrix
 276 and the biases have to be re-loaded. However, both the recurrent matrix \mathbf{R} and the biases remain the
 277 same for the whole time loop and the previous states can stay in memory as they were computed in
 278 the previous time step. Since the structure of the computation remains the same over the time steps,
 279 one can even store most of these values in registers. Registers have the highest memory bandwidth
 280 and, while they can only be used within the lowest computation unit (threads), their total size on a
 281 GPU is comparable to the SRAM (both 256 KB per SM on H100).

282 To reach the maximum speed, we implement FlashRNN **fused** kernels that store the recurrent matrix
 283 \mathbf{R} and the biases \mathbf{b} in registers (and SRAM if register memory is exceeded). The matrix multiplication
 284 results are stored and accumulated in shared memory (or HBM if SRAM sizes are exceeded).
 285 In the forward pass, the computations are mainly tiled along the gate dimension (or the dimension of
 286 the new hidden states). This way, we use the maximum amount of memory along the previous state
 287 dimension. This dimension is the accumulating dimension of the recurrent matrix multiplication.
 288 For the backward pass, the computations are typically tiled along the previous state gradient dimension,
 289 such that the gate dimension, which is accumulated over, is minimally tiled. Algorithm 5.1
 290 shows a simplified representation of the forward pass in pseudo-code.

291 292 5.3 TRITON IMPLEMENTATION

293 With FlashRNN we also implement a Triton⁴ variant of the fused FlashRNN kernel. Triton is a
 294 domain specific language and compiler for parallel programming that provides a Python-based environment
 295 for writing custom GPU kernels.
 296

297 For the Triton kernel we partition the computation along two dimensions the batch dimension and
 298 the head dimension. As described in section 4.2 we partition the embedding dimension into multiple
 299 heads and compute each head in parallel in different programs with no synchronization in between
 300 these programs. Since Triton schedules each program on different streaming multiprocessors (SM),
 301 each SM will hold its recurrent weight matrix \mathbf{R}_{head} and bias \mathbf{b}_{head} in SRAM. In contrast to CUDA,
 302 Triton gives no access to registers on the GPU. Therefore, we cannot apply the custom caching
 303 strategy of the fused CUDA kernels and instead rely on Triton for managing the shared memory and
 304 register cache on each SM. Also, there is no grid synchronization between SMs in Triton, making it
 305 impossible to communicate values between SMs over HBM. In section 6.1 we find that this poses
 306 a limitation on the maximum head dimension of 128 for the forward pass and 64 for the backward
 307 pass on a NVIDIA H100 GPU.

308 The recurrent matrix multiply in equation 1 and 3 is implemented with Triton’s matrix multiply
 309 operation `tl.dot` which gives an interface to the Tensor Core units on GPUs. In Triton minimum
 310 block size of these matrix multiplies is 16x16, which gives a limit on the minimum batch size. In
 311 practice, we enable smaller batch sizes by padding zeros at the cost of efficiency.

312 313 5.4 AUTOMATIC TUNING OF TILING AND LOOPING DIMENSIONS

314 While Algorithm 5.1 describes the algorithmic behaviour, the tile, block and grid sizes and loop
 315 iterations depend on the specific hardware architecture, i.e. the number of computational super-
 316 units (streaming multiprocessors), the SRAM per super-unit, the sizes of matrix-multiplication units,
 317 threads (warps and threads) per super-unit and the number of registers per thread. On NVIDIA
 318 H100s (and most other NVIDIA GPUs), there is a varying amount of registers per thread, depend-
 319 ing on the block size used. The total number of registers on chip per streaming multiprocessor is
 320 physically fixed.

321 These physical constraints can now be reformulated as equalities, inequalities and divisibility con-
 322 straints inside an integer constraint satisfaction problem (integer CSP). Typically this optimization is
 323

⁴<https://triton-lang.org>

done via polyhedral constraint optimization in compilers (Baghdadi et al., 2018). For solving these constraints in FlashRNN, we implement an efficient solver ConstrINT in Python for general integer CSPs going over large number ranges and including the notion of divisibility constraints, which are needed to model the minimal matrix sizes.

For more details on the solution algorithm, see Appendix Section B.

6 EXPERIMENTS

In Section 6.1 we benchmark the runtime of our FlashRNN kernels and compare against the LSTM and Attention implementations provided in PyTorch. In Section 6.2 we measure training time with FlashRNN kernels on language modeling. Finally, in Section 6.3 we confirm that traditional RNNs like LSTM and more recent variants like sLSTM implemented in FlashRNN can solve state tracking problems.

6.1 RUNTIME BENCHMARK

We evaluate the runtime of all backends of our FlashRNN library that implement the LSTM operation:

- **CUDA fused:** CUDA implementation that fuses matrix multiplication and pointwise operations of the LSTM in a single kernel that is persistent over all time iterations.
- **CUDA alternating:** CUDA implementation that implements the time loop in C++ and alternates between a matrix multiply kernel and a LSTM pointwise kernel.
- **Triton fused:** Triton implementation that fuses matrix multiplication and pointwise operations similar to CUDA fused.
- **Vanilla PyTorch:** PyTorch implementation of the LSTM operation with our custom backward pass implementation, which is faster than the PyTorch autograd backward pass. We do not use `torch.compile` due to very long compile times.

We compare our backends to two references from PyTorch:

- **FlashAttention2:** PyTorch Attention⁵ with FlashAttention2 backend. Note that FlashAttention2 is not a recurrent operation and can be parallelized across batch, head, and sequence dimension on the GPU.
- **nn.LSTM:** PyTorch LSTM with NVIDIA cuDNN as backend. In contrast to our FlashRNN LSTM, `nn.LSTM` also integrates the gate pre-activation computation into the function call (not kernel call), which we do not (see Section 3). In Section C.2 in the appendix we provide an comparison where we compare the combination of a linear layer and our FlashRNN LSTM kernel with `nn.LSTM`.

Setup. We assess the impact of the input dimensions batch size (B), sequence length (T) and head dimension (DH) and number of heads (NH). The number of heads together with the head dimension give the embedding dimension $d = NH \times DH$. Except for PyTorch `nn.LSTM` we run all runtime experiments with bfloat16 precision. For `nn.LSTM` we use float16 precision, since this precision yielded the fastest runtimes. For every runtime measurement we do 25 warmup iterations and then report the average across 1000 iterations on NVIDIA H100 GPUs. We use PyTorch 2.4 and with CUDA version 12.4 for our experiments. Further details and additional experiments can be found in Section C in the appendix.

Head dimension. We measure the runtime of all of our FlashRNN kernels and our two references FlashAttention2 and PyTorch `nn.LSTM` for different head dimensions. We fix the embedding dimension $d = NH \times DH$ to 768 and vary the head dimension from 16 to 768. We use batch size 16 and sequence length 1024. In Figure 2 we report the runtime of each the forward pass only on

⁵https://pytorch.org/docs/stable/generated/torch.nn.functional.scaled_dot_product_attention.html

the left and the forward combined with the backward pass. FlashAttention2 does not allow for head dimension larger than 256, due shared memory limitation. The PyTorch `nn.LSTM` does not support multiple heads or blockdiagonal recurrent matrices. Therefore, we only report the runtime for a single head of dimension 768, including the gate pre-activation computation. At this dimension `nn.LSTM` is about 3 times faster than CUDA fused. The Triton kernels are limited to head dimension 128 and 64, but are about two times faster than CUDA fused for small head dimensions 16 and 32. The fused CUDA kernels support all head dimensions up to 768 (actually 2560 on H100) and are about two times faster than the alternating kernels.

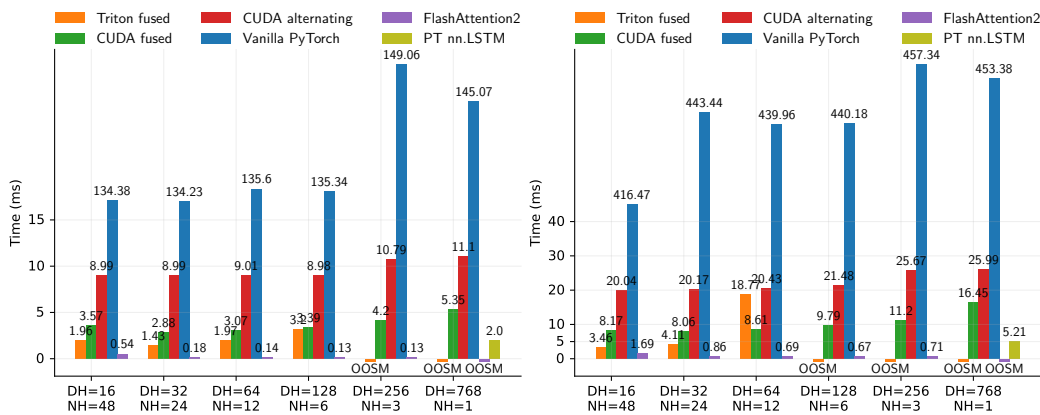


Figure 2: LSTM Runtime for different head dimensions (DH) and number of heads (NH) on a NVIDIA H100. Overall embedding dimension is fixed at 768. We use batch size 16 and sequence length 1024. **Left:** Forward pass. **Right:** Forward + backward pass.

Batch size. We measure the runtime of all LSTM kernels while varying the batch size (B) from 2 to 256 at sequence length 1024. Figure 3 shows the results for NH=12 heads with head dimension DH=64. The CUDA fused backend is optimized for smaller batch sizes and shows a 2x speed up over the alternating backend for batch sizes up to 32. For larger batch sizes than 128 CUDA alternating is faster. Figure 4 shows the results for a single head with head dimension DH=768. At this head dimension CUDA fused is still faster than CUDA alternating up to batch size 32. For larger batch sizes, CUDA alternating is more than two times faster. Comparing to the PyTorch `nn.LSTM`, we find for medium batch sizes from 8 to 64 it is about 3 times faster than and CUDA fused and for larger batch sizes about about 30% faster than CUDA alternating.

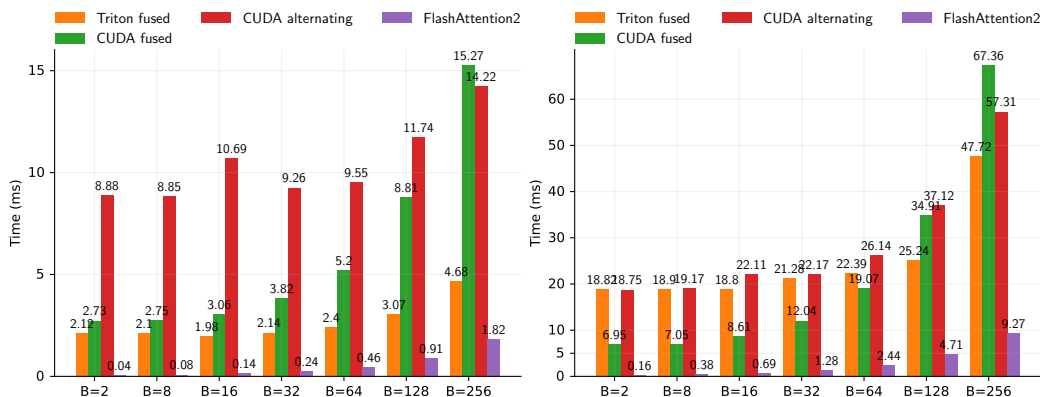


Figure 3: LSTM Runtime for different batch sizes (B) on a NVIDIA H100. We use 12 heads with head dimension 64 and sequence length 1024. **Left:** Forward pass. **Right:** Forward + backward pass.

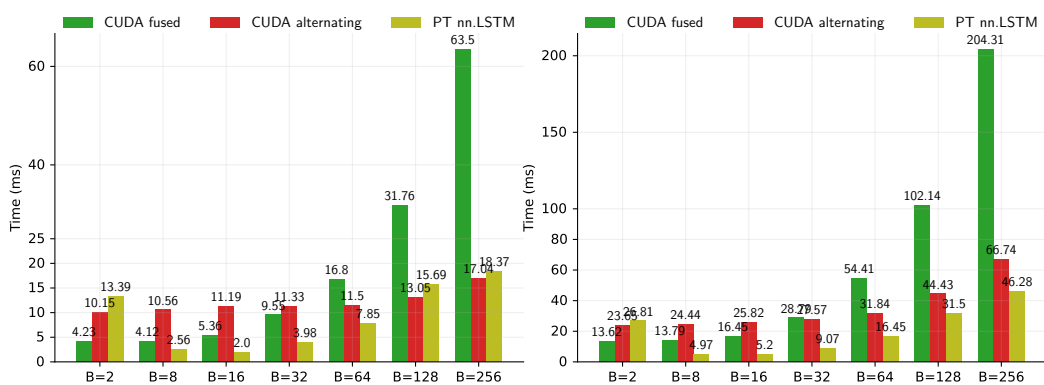


Figure 4: LSTM Runtime for different batch sizes (B) on a NVIDIA H100. We use one head with head dimension 768 and sequence length 1024. **Left:** Forward pass. **Right:** Forward + backward pass.

Additional Runtime Experiments. In section C.1 in the appendix we include experiments on varying sequence lengths. We see the expected linear runtime scaling for our FlashRNN kernels and validate that the above findings transfer to other sequence lengths. In addition, in section C.2 we compare the FlashRNN LSTM kernel in combination with a linear layer that computes the gate pre-activations externally to the PyTorch `nn.LSTM` baseline which integrates the gate pre-activation computation. We find that the gate pre-activation computation has only marginal impact on the overall runtime. Finally, in section C.3 we provide all runtime results also for the sLSTM (Beck et al., 2024).

6.2 LANGUAGE MODELING

Even though we do not expect traditional RNNs to outperform Transformers, the language modeling setting serves as an important benchmark for speed on larger scales. Here, we train models at the 165M parameter scale for a Llama-style Transformer without weight tying, i.e. 12 Transformer blocks with Pre-LayerNorm and a Swish-Gated MLP after the attention layer. We replace attention with FlashRNN LSTM and sLSTM layers for a speed comparison. The results show a slowdown of roughly 25% over attention for equal head dimensions or 140% for one RNN head, see Table 6.2 (H100) and Appendix Table D (A100). In our experiments, we also compare to the cuDNN implementation of LSTM integrated in PyTorch (`torch.nn.LSTM`). While its integration into PyTorch is considerably faster, there are numerical differences to the FlashRNN implementation. With same initialization, FlashRNN LSTMs converge faster in our language experiments (both bfloat16 and float32), even though the differences in a single kernel call are at the expected levels of numerical precision. This deviation should be investigated further and suggests the use of FlashRNN even for the established LSTM architecture.

For larger models, we expect local batch sizes to be smaller and the effective speed difference for fused kernels to be higher compared to the alternating version - as measured in Section 6.1.

6.3 STATE TRACKING TASK

To show state tracking capabilities of traditional RNNs in contrast to Transformers and State Space Models experimentally, we train our implementation on the Parity task and evaluate on longer sequences to measure extrapolation capabilities. (Zhou et al., 2024) This serves as a litmus test for state tracking capabilities (Merrill et al., 2024).

Model	Heads	Param. (M)	Train Time (h)	Median Step (s)	Val PPL (val)
LSTM CUDA fused	1	190	9.9	0.535	22.1
LSTM CUDA altern.	1	190	10.8	0.575	21.9
LSTM PT _{nn} .LSTM	1	190	4.5	0.285	25.8
LSTM CUDA fused	12	164	5.9	0.325	22.2
LSTM CUDA altern.	12	164	9.6	0.511	22.1
sLSTM CUDA fused	1	190	10.1	0.543	21.3
sLSTM CUDA altern.	1	190	10.9	0.577	21.4
sLSTM CUDA fused	12	164	6.8	0.342	21.7
sLSTM CUDA altern.	12	164	9.7	0.509	21.8
Transformer	12	162	2.9	0.190	17.9

Table 1: 165M Model training on 15B tokens of SlimPajama on 8xH100s with two gradient accumulation steps.

Model	Transformer	Mamba	mLSTM	Elman	GRU	LSTM	sLSTM
Acc (Ext.)	0.52	0.56	0.54	1.00	1.00	1.00	1.00

Table 2: Parity Task in Sequence Extrapolation: Transformers, State Space Models and mLSTM fails at this task (close to random chance at 0.5), while traditional recurrent models can learn to extrapolate. Accuracies are averaged over three seeds for the best respective learning rate.

7 CONCLUSION

The FlashRNN library serves as a fast and extendable implementation of traditional RNNs with a recurrent connection or memory mixing. It extends RNNs with the multi-head paradigm introduced by Beck et al. (2024) for sLSTM. FlashRNN provides a speed-up of up to 50x over vanilla PyTorch implementations of RNNs and may serve as a backbone for future RNN architectures that have a recurrent connection.

FlashRNN implements two variants, an alternating version switching between point-wise and matrix-multiplication kernels and a fused implementation, optimizing memory transfers while using hardware-optimized matrix-multiplication. The second leads to a further 3-4x speed-up over the alternating option for small batch sizes. The implementation auto-optimizes its internal sizes for different cache levels via the ConstrINT library - a custom library solving general integer constraint satisfaction problems with equality, inequality and divisibility constraints. This library may be re-used for other optimization problems regarding cache sizes on hardware platforms and beyond.

We show that with FlashRNN, traditional RNNs are not too far in speed from Transformers in practice, even though they are not parallelizable along the sequence dimension. In the future, it may be optimized to leverage asynchronous memory operations and inter-SRAM connections - recent hardware features that promise further speed ups not realized in this work.

ETHICS STATEMENT

We use an open dataset (SlimPajama) that uses publicly crawled internet data for Language Model training. Our implementation speeds up a certain class of Machine Learning models. This may reduce the environmental impact of the research field, in case these architectures remain important in future research. Also it may speed up development of Machine Learning research in the direction of recurrent sequence modeling with state tracking capabilities. The further implications of these impacts may or may not be a benefit for society.

540 REPRODUCIBILITY STATEMENT

541
542 We provide the source code for your implementations along with this paper. The detailed training
543 setup for speed tests is described in Section 6.1. For Language Modeling this setup description is
544 provided in Appendix Section E and uses the open SlimPajama dataset, for the parity task exper-
545 iments in Appendix Section F, the training and test data can be synthetically generated using the
546 mentioned distributions.

547 The observed deviations in language model training compared to the PyTorch LSTM based on
548 cuDNN should be further investigated. The results on A100 and H100, as well as across our differ-
549 ent kernels are within the expected small-scale numerical deviations.

550 The code will also be released on GitHub including training scripts.

551
552 REFERENCES

553 Riyadh Baghdadi, Jessica Ray, Malek Ben Romdhane, Emanuele Del Sozzo, Abdurrahman Akkas,
554 Yunming Zhang, Patricia Suriana, Shoaib Kamil, and Saman Amarasinghe. Tiramisu: A
555 Polyhedral Compiler for Expressing Fast and Portable Code, December 2018. URL <http://arxiv.org/abs/1804.10694>. arXiv:1804.10694 [cs].
556
557

558 M. Beck, K. Pöppel, M. Spanring, A. Auer, O. Prudnikova, M. Kopp, G. Klambauer, J. Brandstetter,
559 and S. Hochreiter. xLSTM: Extended Long Short-Term Memory, May 2024. URL <http://arxiv.org/abs/2405.04517>. arXiv:2405.04517 [cs, stat].
560

561 K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio.
562 Learning phrase representations using RNN encoder–decoder for statistical machine translation.
563 In Alessandro Moschitti, Bo Pang, and Walter Daelemans (eds.), *Proceedings of the 2014 Con-*
564 *ference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1724–1734, Doha,
565 Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1179.
566 URL <https://aclanthology.org/D14-1179>.
567

568 T. Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. In
569 *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=mZn2Xyh9Ec>.
570

571 T. Dao, D. Y. Fu, S. Ermon, A. Rudra, and C. Ré. FlashAttention: Fast and memory-efficient exact
572 attention with IO-awareness. In *Advances in Neural Information Processing Systems (NeurIPS)*,
573 2022.
574

575 Tri Dao and Albert Gu. Transformers are SSMS: Generalized models and efficient algorithms
576 through structured state space duality. In *Forty-first International Conference on Machine Learn-*
577 *ing*, 2024. URL <https://openreview.net/forum?id=ztN8FCR1td>.

578 J. Degraeve, F. Felici, J. Buchli, et al. Magnetic control of tokamak plasmas through deep reinforce-
579 ment learning. *Nature*, 602:414–419, 2022. doi: 10.1038/s41586-021-04301-9.
580

581 G. Delétang, A. Ruoss, J. Grau-Moya, T. Genewein, L. K. Wenliang, E. Catt, C. Cundy, M. Hutter,
582 S. Legg, J. Veness, and P. A. Ortega. Neural networks and the chomsky hierarchy. In *Eleventh*
583 *International Conference on Learning Representations*, 2023.

584 J. L. Elman. Finding Structure in Time. *Cognitive Science*, 14(2):179–211, March 1990. ISSN 0364-
585 0213, 1551-6709. doi: 10.1207/s15516709cog1402_1. URL https://onlinelibrary.wiley.com/doi/10.1207/s15516709cog1402_1.
586

587 F. Cummins F.A. Gers, J. Schmidhuber. Learning to forget: continual prediction with
588 LSTM. In *9th International Conference on Artificial Neural Networks: ICANN '99*, vol-
589 *ume 1999*, pp. 850–855, Edinburgh, UK, 1999. IEE. ISBN 978-0-85296-721-8. doi:
590 10.1049/cp:19991218. URL https://digital-library.theiet.org/content/conferences/10.1049/cp_19991218.
591
592

593 A. Gu and T. Dao. Mamba: Linear-Time Sequence Modeling with Selective State Spaces, December
2023. URL <http://arxiv.org/abs/2312.00752>. arXiv:2312.00752 [cs].

- 594 S. Hochreiter and J. Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9(8):
595 1735–1780, November 1997. ISSN 0899-7667, 1530-888X. doi: 10.1162/neco.1997.9.
596 8.1735. URL [https://www.mitpressjournals.org/doi/abs/10.1162/neco.](https://www.mitpressjournals.org/doi/abs/10.1162/neco.1997.9.8.1735)
597 1997.9.8.1735.
- 598 A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret. Transformers are rnns: Fast autoregressive
599 transformers with linear attention. In *Proceedings of the International Conference on Machine*
600 *Learning (ICML)*, 2020.
- 602 C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh. Basic Linear Algebra Subprograms for
603 Fortran Usage. *ACM Transactions on Mathematical Software*, 5(3):308–323, September 1979.
604 ISSN 0098-3500, 1557-7295. doi: 10.1145/355841.355847. URL [https://dl.acm.org/](https://dl.acm.org/doi/10.1145/355841.355847)
605 [doi/10.1145/355841.355847](https://dl.acm.org/doi/10.1145/355841.355847).
- 606 A. K. Mackworth. Consistency in networks of relations. *Artificial Intelligence*, 8(1):99–118,
607 February 1977. ISSN 00043702. doi: 10.1016/0004-3702(77)90007-8. URL [https://](https://linkinghub.elsevier.com/retrieve/pii/0004370277900078)
608 linkinghub.elsevier.com/retrieve/pii/0004370277900078.
- 609 W. Merrill and A. Sabharwal. The Parallelism Tradeoff: Limitations of Log-Precision Transformers.
610 *Transactions of the Association for Computational Linguistics*, 11:531–545, 06 2023. ISSN 2307-
611 387X. doi: 10.1162/tacl.a.00562. URL [https://doi.org/10.1162/tacl.a.](https://doi.org/10.1162/tacl.a.00562)
612 00562.
- 613 W. Merrill, J. Petty, and A. Sabharwal. The illusion of state in state-space models. In *Forty-first*
614 *International Conference on Machine Learning*, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=QZgo9JZpLq)
615 [forum?id=QZgo9JZpLq](https://openreview.net/forum?id=QZgo9JZpLq).
- 616 M. Mozer. A focused backpropagation algorithm for temporal pattern recognition. *Complex Sys-*
617 *tems*, 3, 01 1995.
- 618 G. Nearing, D. Cohen, V. Dube, M. Gauch, O. Gilon, S. Harrigan, A. Hassidim, D. Klotz, F. Kratzert,
619 A. Metzger, S. Nevo, F. Pappenberger, C. Prudhomme, G. Shalev, S. Shenzi, T. Y. Tekalign,
620 D. Weitzner, and Y. M. B. Kosko. Global prediction of extreme floods in ungauged watersheds.
621 *Nature*, 627:559–563, 2024. doi: 10.1038/s41586-024-07145-1.
- 622 J. Shah, G. Bikshandi, Y. Zhang, V. Thakkar, P. Ramani, and T. Dao. Flashattention-3: Fast and
623 accurate attention with asynchrony and low-precision, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2407.08608)
624 [abs/2407.08608](https://arxiv.org/abs/2407.08608).
- 625 N. Sharvil. Haste: a fast, simple, and open rnn library. [https://github.com/lmnt-com/](https://github.com/lmnt-com/haste/)
626 [haste/](https://github.com/lmnt-com/haste/), Jan 2020.
- 627 B. Spector, A. Singhal, S. Arora, and C. Ré. GPUs Go Brrr. URL [https://hazyresearch.](https://hazyresearch.stanford.edu/blog/2024-05-12-tk)
628 [stanford.edu/blog/2024-05-12-tk](https://hazyresearch.stanford.edu/blog/2024-05-12-tk).
- 629 V. Thakkar, P. Ramani, C. Cecka, A. Shivam, H. Lu, E. Yan, J. Kosaian, M. Hoemmen, H. Wu,
630 A. Kerr, M. Nicely, D. Merrill, D. Blasig, F. Qiao, P. Majcher, P. Springer, M. Hohnerbach,
631 J. Wang, and M. Gupta. Cutlass, 1 2023. URL [https://github.com/NVIDIA/cutlass/](https://github.com/NVIDIA/cutlass/tree/v3.0.0)
632 [tree/v3.0.0](https://github.com/NVIDIA/cutlass/tree/v3.0.0).
- 633 A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N Gomez, Ł. Kaiser,
634 and I. Polosukhin. Attention is All you Need. In I. Guyon, U. Von Luxburg,
635 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Ad-*
636 *vances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,
637 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)
638 [file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 639 S. Williams, A. Waterman, and D. Patterson. Roofline: an insightful visual performance model
640 for multicore architectures. *Communications of the ACM*, 52(4):65–76, April 2009. ISSN 0001-
641 0782, 1557-7317. doi: 10.1145/1498765.1498785. URL [https://dl.acm.org/doi/10.](https://dl.acm.org/doi/10.1145/1498765.1498785)
642 [1145/1498765.1498785](https://dl.acm.org/doi/10.1145/1498765.1498785). Publisher: Association for Computing Machinery (ACM).
- 643 S. Yang, B. Wang, Y. Shen, R. Panda, and Y. Kim. Gated linear attention transformers with
644 hardware-efficient training. In *Forty-first International Conference on Machine Learning*, 2024.
645 URL <https://openreview.net/forum?id=ia5XvxFUJT>.

648 H. Zhou, A. Bradley, E. Littwin, N. Razin, O. Saremi, J. M. Susskind, S. Bengio, and P. Nakkiran.
649 What algorithms can transformers learn? a study in length generalization. In *The Twelfth Interna-*
650 *tional Conference on Learning Representations*, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=AssIuHnmHX)
651 [forum?id=AssIuHnmHX](https://openreview.net/forum?id=AssIuHnmHX).
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A RNN VARIANTS WITH MEMORY MIXING / RECURRENT CONNECTIONS MODELED IN FLASHRNN

Elman RNNs (Elman, 1990) Number of states: 1, Number of gates: 1

$$\mathbf{s}_t^{(0)} = \tanh(\mathbf{g}_t^{(0)}) \quad (8)$$

Here, we omit as well a possible post-processing of the sequence that does not inter-mix states of different time steps. This can as well be parallelized for offline training.

LSTM (Hochreiter & Schmidhuber, 1997; F.A. Gers, 1999) Number of states: 2, Number of gates: 4

States: $\mathbf{h}_t = \mathbf{s}_t^{(0)}$ hidden state, $\mathbf{c}_t = \mathbf{s}_t^{(1)}$ cell state

Gates: $\mathbf{z}_t = \mathbf{g}_t^{(0)}$ cell input, $\mathbf{f}_t = \mathbf{g}_t^{(1)}$ forget gate, $\mathbf{i}_t = \mathbf{g}_t^{(2)}$ input gate, $\mathbf{o}_t = \mathbf{g}_t^{(3)}$ output gate

$$\mathbf{h}_t = \sigma(\mathbf{o}_t) \tanh(\mathbf{c}_t) \quad (9)$$

$$\mathbf{c}_t = \sigma(\mathbf{f}_t) \mathbf{c}_{t-1} + \sigma(\mathbf{i}_t) \tanh(\mathbf{z}_t) \quad (10)$$

GRU (Cho et al., 2014) Number of states: 1, Number of gates: 4 (in the definition of this paper)

States: $\mathbf{s}_t^{(0)}$ hidden state

Gates: $\mathbf{z}_t = \mathbf{g}_t^{(0)}$ cell input, $\mathbf{r}_t = \mathbf{g}_t^{(1)}$ forget gate, $\mathbf{n}_t = \mathbf{g}_t^{(2)}$ input gate, $\mathbf{o}_t = \mathbf{g}_t^{(3)}$ output gate Here, the \mathbf{n}_t gate is not dependent on the previous state, whereas the \mathbf{g}_t gate is not dependent on the input. This behavior can be modeled in FlashRNN as well.

$$\mathbf{h}_t = \sigma(\mathbf{z}_t) \mathbf{h}_{t-1} + (1 - \sigma(\mathbf{z}_t)) \tanh(\mathbf{n}_t + \sigma(\mathbf{r}_t) \tanh(\mathbf{g}_t)) \quad (11)$$

sLSTM (Beck et al., 2024) Number of states: 4, Number of gates: 4 States: $\mathbf{h}_t = \mathbf{s}_t^{(0)}$ hidden state, $\mathbf{c}_t = \mathbf{s}_t^{(1)}$ cell state, $\mathbf{n}_t = \mathbf{s}_t^{(2)}$ normalizer state, $\mathbf{m}_t = \mathbf{s}_t^{(3)}$ stabilizer state

Gates: $\mathbf{z}_t = \mathbf{g}_t^{(0)}$ cell input, $\mathbf{f}_t \mathbf{g}_t^{(1)}$ forget gate, $\mathbf{i}_t = \mathbf{g}_t^{(2)}$ input gate, $\mathbf{o}_t = \mathbf{g}_t^{(3)}$ output gate

$$\mathbf{h}_t = \sigma(\mathbf{o}_t) \frac{\mathbf{c}_t}{\mathbf{n}_t} \quad (12)$$

$$\mathbf{c}_t = \exp((\log \sigma(\mathbf{f}_t) + \mathbf{m}_{t-1} - \mathbf{m}_t) \mathbf{c}_{t-1} + \exp(\mathbf{i}_t - \mathbf{m}_t) \tanh(\mathbf{z}_t)) \quad (13)$$

$$\mathbf{n}_t = \exp((\log \sigma(\mathbf{f}_t) + \mathbf{m}_{t-1} - \mathbf{m}_t) \mathbf{n}_{t-1} + \exp(\mathbf{i}_t - \mathbf{m}_t)) \quad (14)$$

$$\mathbf{m}_t = \max(\log \sigma(\mathbf{f}_t) + \mathbf{m}_{t-1}, \mathbf{i}_t) \quad (15)$$

B CONSTRAINT RESOLUTION ALGORITHMS

To model the hardware constraints, we define IntegerVariables, e.g. a variable describing a tiling size in the FlashRNN algorithm or a constant that defines the total SRAM for one streaming multiprocessor. These can attain a set of numbers (domain), e.g. initially a large range for a so far unconstrained tiling size or a certain value for a constant. These variables can be composed to terms via addition and multiplication, and these terms can be constrained via equalities, inequalities and divisibility constraints.

Specific resolution variables additionally have a heuristic added that defines the behaviour of iteration for choosing among possible values. If the domain of all resolution variables is reduced to a single number, this is a solution. The heuristic gives an order of these variables and for each variable, if smaller or larger values are expected to result in a "better" solution. This helps optimization as there might be many possible solutions, but certain ones promise most speed-ups (e.g. using most TensorCores).

At the lowest level, a term is composed of two IntegerVariables (or intermediate variables), so constraints on it propagate down to the two summand or factor IntegerVariables. Equality, Inequality and Divisibility constraints propagate to the contained terms as well. For example, since all numbers are strictly positive the upper bound on a sum of two IntegerVariables applies to both the summands - minus one. Applying the constraints iteratively upwards and downwards in the expression parse

tree until convergence (i.e. no change for any variable) leads to an arc-consistent state, which we call "Global ARC-Reduce". The binary "ARC-Reduce" algorithm is part of the "AC-3", a constraint satisfaction problem solver for a more general setting (Mackworth, 1977). An arc-consistent state might still have no solution, it is merely a super-set of all possible solutions. Based on the heuristic the ConstrINT algorithm applies a depth-first tree search with "Global ARC-Reduce" application at each step and backtracking for an empty solution domain. (see also Appendix Section B)

Algorithm 2 ConstrINT Resolution

Require: Input Constants / Variables
Require: Resolution Variables with Heuristic
Require: Equality, Inequality and Divisibility Constraints
 Generate intermediate/background variables for terms that propagate constraints
 Reach arc consistency via "Global ARC-Reduce"
if any variable has empty domain $|D_V| = 0$ **then**
 return "No Solution viable"
end if
 Sort values for each Resolution Variable via Heuristic
while any Resolution Variable has domain $|D_V| > 1$ (not fixed) **do**
 Choose Variable via Heuristic, Increase Index Count
 if Lowest Order Variable has empty domain **then**
 return "No Solution viable"
 end if
 Set Variable Value via Heuristic
 Reach arc consistency via Global ARC-Reduce
 if any variable has empty domain $|D_V| = 0$ **then**
 Backtrack
 end if
end while
return Solution

Algorithm 3 ConstrINT Global ARC-Reduce

Require: Expression Parse Tree of Constraints and Variables
 Status=Not Converged
while Change in Root IntegerVariable or NotConverged **do**
 Propagate Restrictions to SubTerms of Expression
 for Sub-Term in Root-Expression **do** ▷ Top-Down Application of Constraints
 Apply Global ARC-Reduce on Sub-Term - Get changes
 end for
 if any change in values for Sub-Term **then** ▷ Bottom-Up Application of Constraints
 Propagate Restriction from Sub-Terms to Root Variable
 Status=Not Converged
 else
 Status=Converged
 end if
end while
return change in Root IntegerVariable

C ADDITIONAL RUNTIME BENCHMARK EXPERIMENTS

C.1 LSTM SEQUENCE LENGTH RUNTIME EXPERIMENTS

We confirm that the findings from Section 6.1 hold true also for varying sequence lengths from 256 to 2048. We fix the batch size to 16 and measure the runtime for 12 heads with head dimension 64 (see Figure 5) and a single head with head dimension 768 (see Figure 6). In these experiments we see the expected linear scaling of the runtime of all LSTM kernels for increasing sequence lengths. The previous findings transfer across sequence lengths.

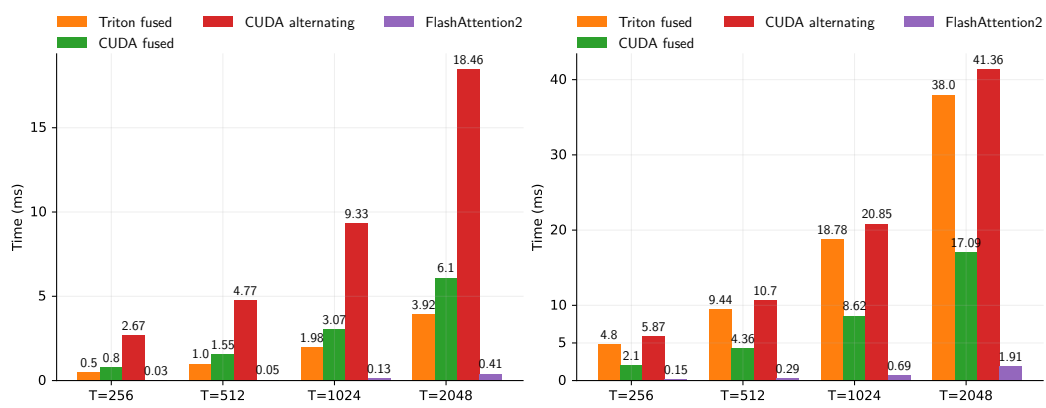


Figure 5: LSTM Runtime for different sequence lengths (T) on a NVIDIA H100. We use 12 heads with head dimension 64 and batch size 16. **Left:** Forward pass. **Right:** Forward + backward pass.

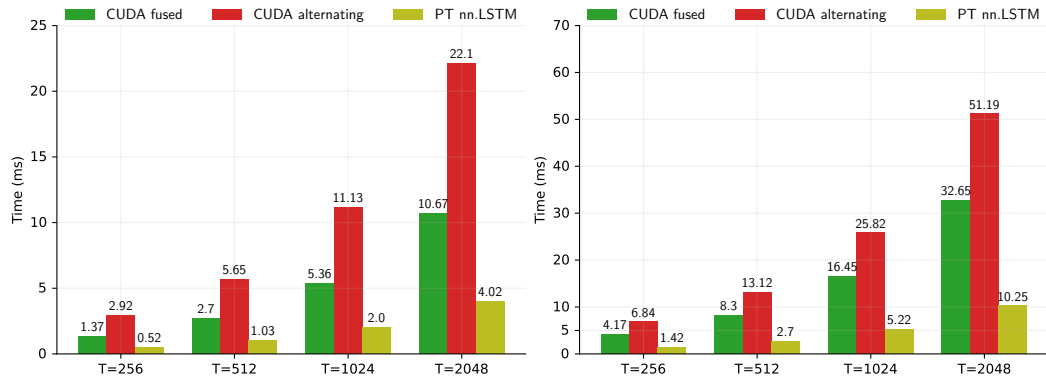


Figure 6: LSTM Runtime for different sequence lengths (T) on a NVIDIA H100. We use one head with head dimension 768 and batch size 16. **Left:** Forward pass. **Right:** Forward + backward pass.

C.2 FLASHRNN WITH EXTERNAL GATE PRE-ACTIVATION COMPUTATION

In Figure 7 we compare the kernel runtimes of the CUDA alternating and CUDA fused kernel with and without the external gate preactivation. *w/ Linear* denotes with external gate preactivation computation via a linear layer. The impact of the gate preactivation computation is marginal compared to the overall kernel runtime.

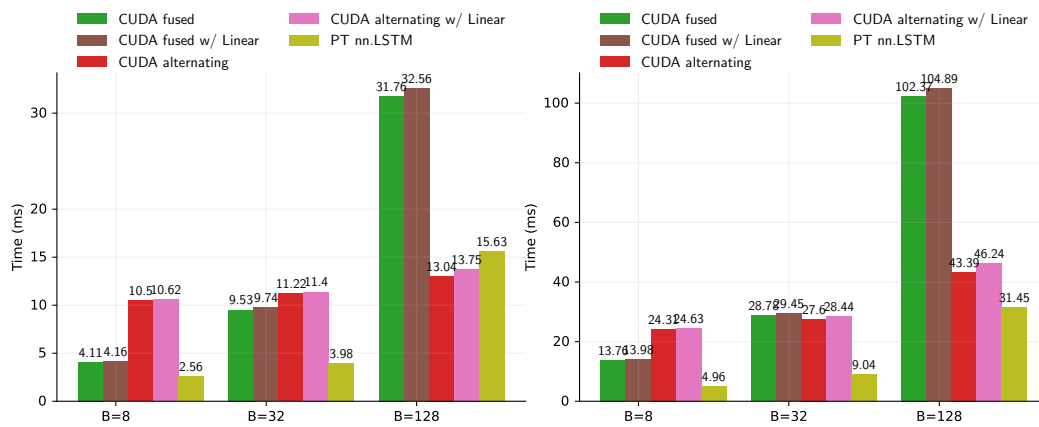


Figure 7: LSTM Runtime for different batch sizes (B) on a NVIDIA H100. We use one head with head dimension 768. We compare the kernel runtime with and without the gate preactivation matrix multiplication. **Left:** Forward pass. **Right:** Forward + backward pass.

C.3 sLSTM RUNTIME EXPERIMENTS

In Figures 8, 9, 10, 11 and 12 we show the results of the experiments from Section 6.1 for the sLSTM (Beck et al., 2024).

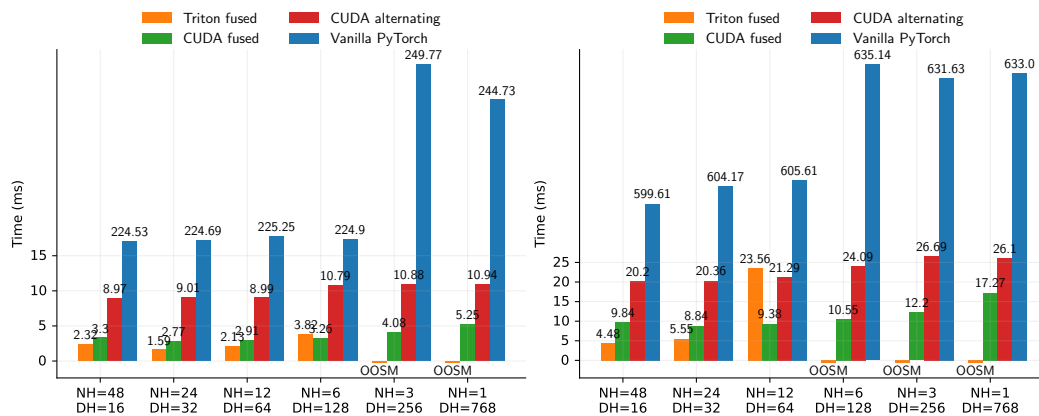
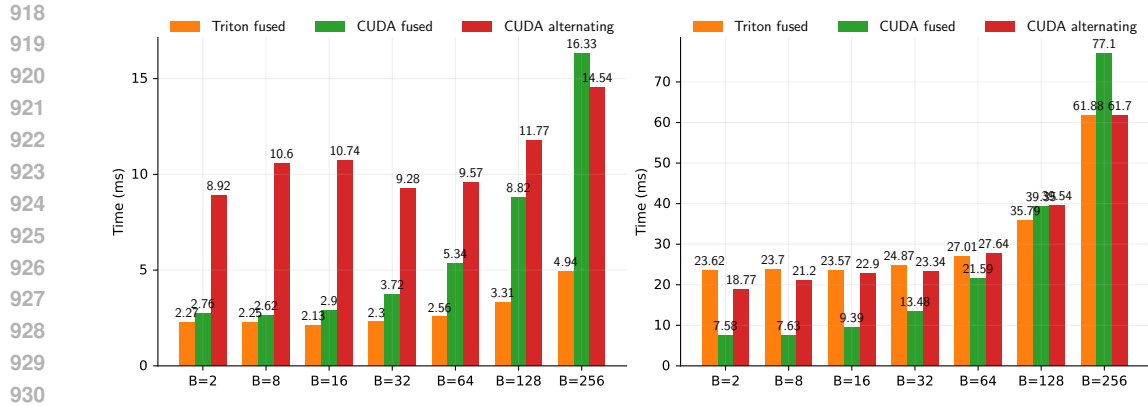


Figure 8: sLSTM Runtime for different head dimensions (DH) and number of heads (NH) on a NVIDIA H100. Overall embedding dimension is fixed at 768. We use batch size 16 and sequence length 1024. **Left:** Forward pass. **Right:** Forward + backward pass.



932 Figure 9: sLSTM Runtime for different batch sizes (B) on a NVIDIA H100. We use 12 heads with
933 head dimension 64 and sequence length 1024. **Left:** Forward pass. **Right:** Forward + backward
934 pass.

935

936

937

938

939

940

941

942

943

944

945

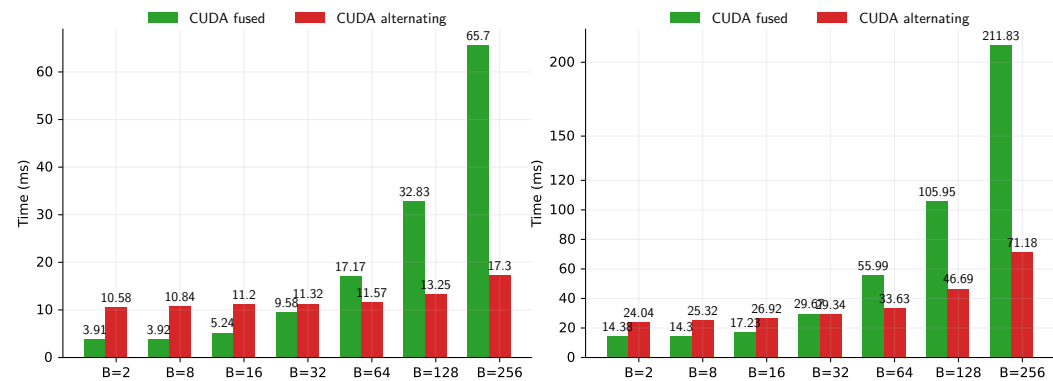
946

947

948

949

950



951 Figure 10: sLSTM Runtime for different batch sizes (B) on a NVIDIA H100. We use one head with
952 head dimension 768 and sequence length 1024. **Left:** Forward pass. **Right:** Forward + backward
953 pass.

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

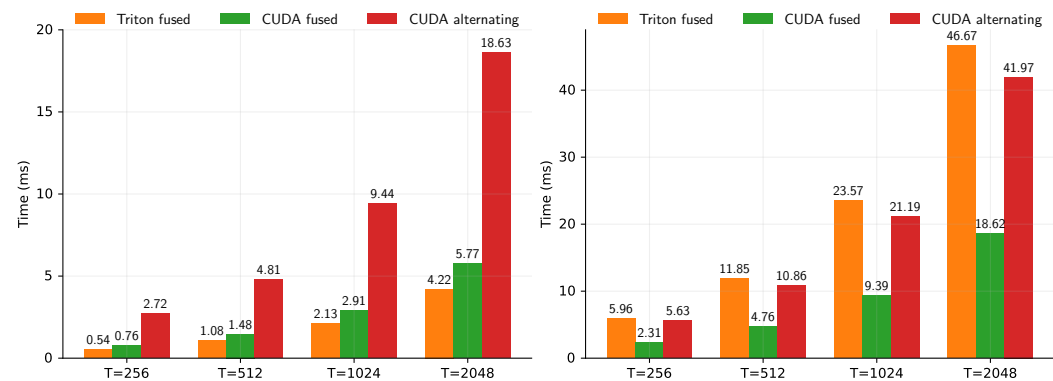


Figure 11: sLSTM Runtime for different sequence lengths (T) on a NVIDIA H100. We use 12 heads
with head dimension 64 and batch size 16. **Left:** Forward pass. **Right:** Forward + backward
pass.

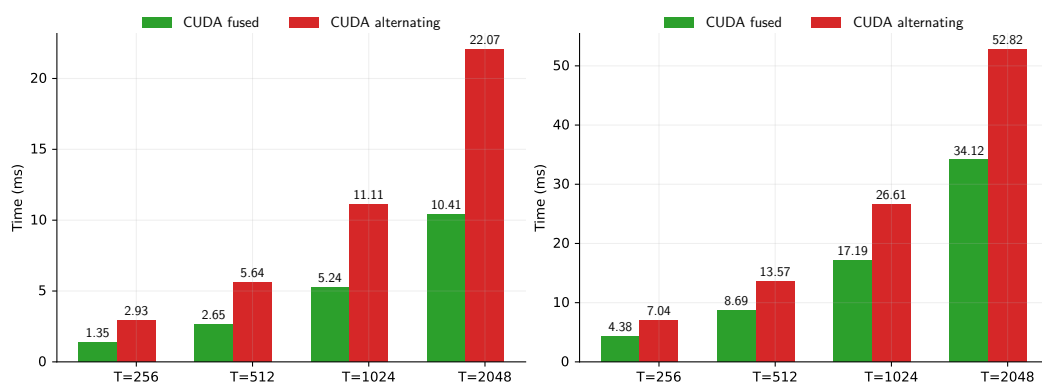


Figure 12: sLSTM Runtime for different sequence lengths (T) on a NVIDIA H100. We use one head with head dimension 768 and batch size 16. **Left:** Forward pass. **Right:** Forward + backward pass.

D LANGUAGE MODEL TRAINING ON A100s

Model	Heads	Param. (M)	Train Time (h)	Median Step (s)	Val PPL (val)
LSTM CUDA fused	1	190	15.4	1.699	22.1
LSTM CUDA altern.	1	190	15.1	1.685	22.1
LSTM PT nn. LSTM	1	190	6.6	0.730	25.9
sLSTM CUDA fused	1	190	15.3	1.707	21.3
sLSTM CUDA altern.	1	190	15.6	1.720	21.3
LSTM CUDA fused	12	164	7.8	0.820	22.3
LSTM CUDA altern.	12	164	7.7	0.809	22.3
sLSTM CUDA fused	12	164	8.0	0.865	21.7
sLSTM CUDA altern.	12	164	8.0	0.852	21.7
Transformer	12	162	6.3	0.688	18.0

Table 3: 165M Model training on 15B tokens of SlimPajama on 8xA100s.

E LANGUAGE TRAINING DETAILS

All models are roughly at 165 M parameter scale, that means 12 Transformer blocks (post-up projection), with a swish-gated MLP, and embedding dimension 768. The Transformer uses RoPE embeddings, whereas the other models do not use any additional positional information. We train with context length 1024 and a global batch size of 512, resulting in roughly 30 k training steps for 15 B tokens of the SlimPajama dataset. We use the GPT-2 tokenizer and learning rate $1e-3$ with linear warmup over 750 steps and cosine decay to $1e-4$ over 30k steps. We use PyTorch in version 2.4.0 and CUDA 12.1 for A100 and 12.4 for H100s. The training uses PyTorch FSDP in the NO.SHARD setting (DDP) with Automated Mixed Precision using bfloat16 and float32 for accumulations.

For the A100 experiments, we use one node of eight A100-SXM (80GB) GPUs and a local batch size of 64. For H100-SXM we reduce the local batch size to 32 and use 2 gradient accumulation steps due to `OutOfMemory` errors, even though they have the same HBM size (80GB).

For the language model trainings, we see more spikes in the training step times for FlashRNN models compared to the PyTorch implementations, which should be investigated further.

F EXPERIMENTAL DETAILS PARITY TASK IN SEQUENCE EXTRAPOLATION

For the parity task we train on the parity task with varying training sequence lengths up to 40. For the reported validation, we evaluate on sequence lengths between 40 and 256. Sequence lengths are uniformly sampled in the respective ranges. We train on three seeds for learning rates $\{1e-2, 1e-3, 1e-4\}$ and choose best learning rates. We train for up to 100k steps with batch size 256 with linear warmup of 10k and cosine decay to 10 % of the peak learning rate. Elman networks and LSTM cannot reach 100 % accuracy on sequence extrapolation for the smallest learning rate. All models reach low losses and high accuracies on the training set.