

Parallelizing the **SCYFI** algorithm for GPUs using the Julia Programming Language

Internship of Computer Science B.Sc., University Heidelberg

Kai Rothe

advised by Daniel Durstewitz and Filip Sadlo

Piece-wise Linear Recurrent Neural Networks

Diagram illustrating the structure of a Piece-wise Linear Recurrent Neural Network (PL-RNN) update function:

$$\mathbf{z}_{t+1} = F_{\boldsymbol{\theta}}(\mathbf{z}_t) = \mathbf{A} \mathbf{z}_t + \mathbf{W} \phi(\mathbf{z}_t) + \mathbf{h}$$

Annotations and definitions:

- $\mathbf{z}_t \in \mathbb{R}^M$ (input vector) points to \mathbf{z}_t in the equation.
- $\mathbf{W} \in \mathbb{R}^{M \times M}$ of coupling weights (off-diagonal matrix) points to \mathbf{W} in the equation.
- bias term $\mathbf{h} \in \mathbb{R}^M$ points to \mathbf{h} in the equation.
- diagonal matrix $\mathbf{A} \in \mathbb{R}^{M \times M}$ of auto-regression weights points to \mathbf{A} in the equation.
- element-wise applied rectified linear unit (ReLU) function $\phi(\mathbf{z}_t) = \max(\mathbf{z}_t, 0)$ points to $\phi(\mathbf{z}_t)$ in the equation.

Searching for Fixed Points and Limit Cycles

$$D_{ij}(\mathbf{z}_t) := \begin{cases} 1, & \text{if } i = j \text{ and } z_{t,i} > 0 \\ 0, & \text{else} \end{cases} \quad \{\mathbf{z}_1^*, \dots, \mathbf{z}_k^*\} \text{ of order } k$$

↓

$$\mathbf{z}_{t+1} = F_{\boldsymbol{\theta}}(\mathbf{z}_t) = \mathbf{A} \mathbf{z}_t + \mathbf{W} \mathbf{D}(\mathbf{z}_t) \mathbf{z}_t + \mathbf{h} = \underbrace{(\mathbf{A} + \mathbf{W} \mathbf{D}(\mathbf{z}_t))}_{\mathbf{W}(\mathbf{z}_t)} \mathbf{z}_t + \mathbf{h}$$

$$\mathbf{z}_k^* = F_{\boldsymbol{\theta}}(\mathbf{z}_{k-1}^*) = \mathbf{W}(\mathbf{z}_{k-1}^*) \mathbf{z}_{k-1}^* + \mathbf{h}$$

$$= \mathbf{W}(\mathbf{z}_{k-1}^*) [\mathbf{W}(\mathbf{z}_{k-2}^*) \mathbf{z}_{k-2}^* + \mathbf{h}] + \mathbf{h}$$

$$= \dots$$

$$= \left(\prod_{r=0}^{k-1} \mathbf{W}(\mathbf{z}_{k-r}^*) \right) \mathbf{z}_k^* + \left[\left(\sum_{j=2}^{k-1} \prod_{r=0}^{k-j} \mathbf{W}(\mathbf{z}_{k-r}^*) \right) + \mathbb{1} \right] \mathbf{h},$$

Searching for Fixed Points and Limit Cycles

$$D_{ij}(\mathbf{z}_t) := \begin{cases} 1, & \text{if } i = j \text{ and } z_{t,i} > 0 \\ 0, & \text{else} \end{cases} \quad \{\mathbf{z}_1^*, \dots, \mathbf{z}_k^*\} \text{ of order } k$$

↓

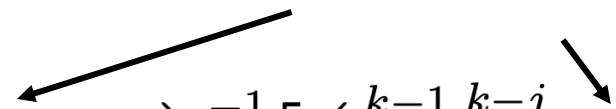
$$\mathbf{z}_{t+1} = F_{\boldsymbol{\theta}}(\mathbf{z}_t) = \mathbf{A} \mathbf{z}_t + \mathbf{W} \mathbf{D}(\mathbf{z}_t) \mathbf{z}_t + \mathbf{h} = \underbrace{(\mathbf{A} + \mathbf{W} \mathbf{D}(\mathbf{z}_t))}_{\mathbf{W}(\mathbf{z}_t)} \mathbf{z}_t + \mathbf{h}$$

$$\mathbf{z}_k^* = \left(\mathbb{1} - \prod_{r=0}^{k-1} \mathbf{W}(\mathbf{z}_{k-r}^*) \right)^{-1} \left[\left(\sum_{j=2}^{k-1} \prod_{r=0}^{k-j} \mathbf{W}(\mathbf{z}_{k-r}^*) \right) + \mathbb{1} \right] \mathbf{h}$$

SCYFI by Eisenman et. al., 2023

1. Randomly guess indicator matrices $\{\mathbf{D}_l\}_{l=1}^k$

2. Solve


$$\mathbf{z}_k^* = \left(\mathbb{1} - \prod_{r=0}^{k-1} \mathbf{W}(\mathbf{z}_{k-r}^*) \right)^{-1} \left[\left(\sum_{j=2}^{k-1} \prod_{r=0}^{k-j} \mathbf{W}(\mathbf{z}_{k-r}^*) \right) + \mathbb{1} \right] \mathbf{h} \quad (1)$$

$$\mathbf{z}_l^* = F^l(\mathbf{z}_k^*) \text{ for } l \in [1, \dots, k-1]$$

3. Check for consistency: $\mathbf{D}(\mathbf{z}_l^*) = \mathbf{D}_l$? $D_{ij}(\mathbf{z}_t) := \begin{cases} 1, & \text{if } i = j \text{ and } \mathbf{z}_{t,i} > 0 \\ 0, & \text{else} \end{cases} \quad (2)$

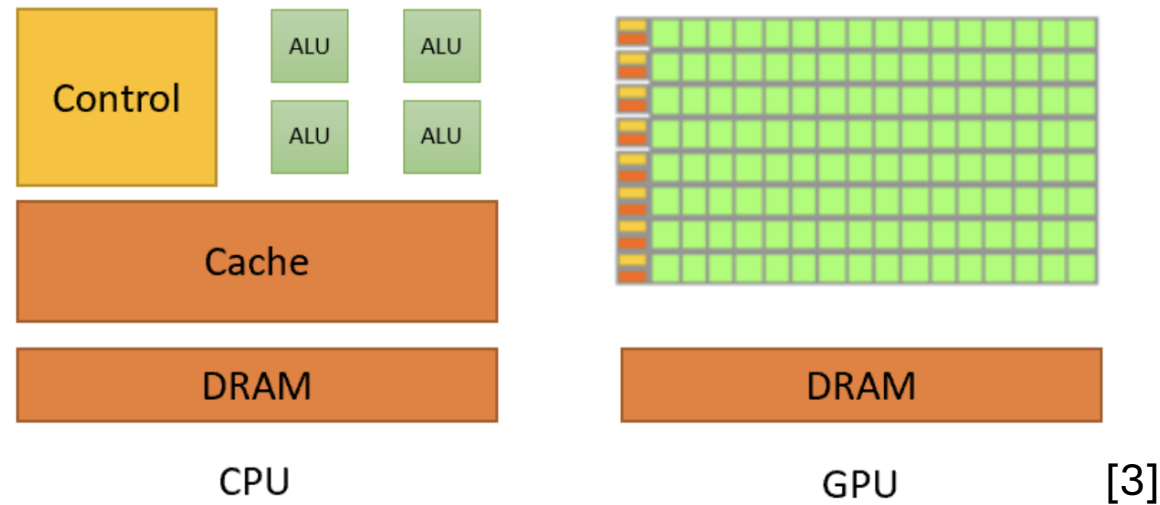
4. Else repeat with $\mathbf{D}_l = \mathbf{D}(\mathbf{z}_l^*)$



CPU Preperation

- profile to find bottlenecks, e.g. flame graphs [4]
- store only diagonals of indicator matrices
- preallocation, in-place & element-wise computations within loops
- simplify control flow, store results of time-consuming computations (avoid repetition), etc.

GPU Computing with CUDA



1. Copy data from host (CPU) memory to device (GPU) memory
2. Load and execute GPU program
3. Copy result from device back to host

GPU Parallelization

- *CuArrays* of CUDA.jl support matrix multiplication and inversion, only rounding operations and eigenvalues on CPU [2]
- avoid CPU fallback
- minimize host-device data transfer: pre-allocation and offloading computations even if standalone slower on GPU
- future improvement (ca. x 1.5): custom SCYFI GPU kernel

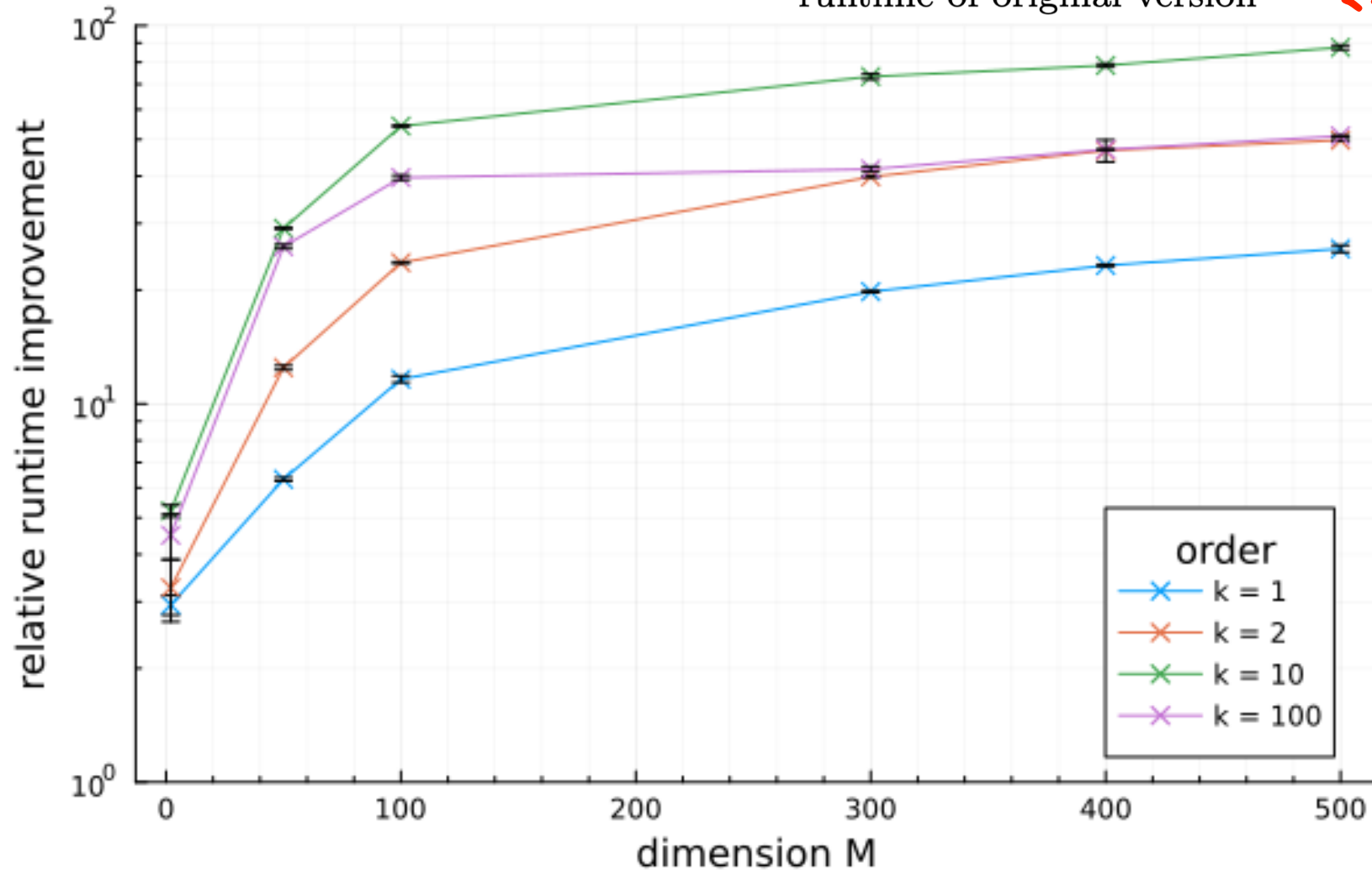
Algorithm 2 GPU implementation of SCYFI

Input: PLRNN parameters \mathbf{A} , \mathbf{W} , \mathbf{h} **Parameters:** $N_{out}(k)$: max. number of random initialisations for searching k -cycles; $N_{in}(k)$: max. number of iterations for searching k -cycles**Output:** \mathcal{L} : tuple of all sets \mathcal{L}_k of all discovered cycles of order $k \in [1, \dots, K_{max}]$


```
1: Preallocate  $\mathcal{L} \rightarrow (\mathcal{L}_1, \dots, \mathcal{L}_{K_{max}})$  as tuple of  $\mathcal{L}_k = \{\}$  for all  $k \in [1, \dots, K_{max}]$ 
2: Transfer  $\mathbf{A}$ ,  $\mathbf{W}$ ,  $\mathbf{h}$  as CuArrays to GPU
3: Preallocate other  $M \times M$  CuArrays (corresponding to factors in (1)) on GPU
4: for  $k \in [1, \dots, K_{max}]$  do
5:   Preallocate diagonals of indicator matrices  $\mathbf{D}_{init}$ ,  $\mathbf{D}^*$  as  $M \times k$  CuArray on GPU
6:   Preallocate trajectory  $\mathbf{Z}$  as  $M \times k$  CuArray on GPU
7:    $i \rightarrow 0$ 
8:   while  $i < N_{out}(k)$  do on the GPU..
9:     Fill  $\mathbf{D}_{init}$  with random binary values
10:     $c \rightarrow 0$ 
11:    while  $c < N_{in}(k)$  do
12:      Solve (1) inplace for cycle candidate  $\mathbf{Z}$  using  $\mathbf{D}_{init}$  and preallocated arrays
13:      Determine  $\mathbf{D}^*$  from  $\mathbf{Z}$  using (2) inplace
14:      if not self-consistent, i.e.  $\mathbf{D}_{init} \neq \mathbf{D}^*$  then
15:         $\mathbf{D}_{init} \rightarrow \mathbf{D}^*$  inplace
16:      else on the CPU..
17:        if  $\mathbf{Z}$  not found yet, i.e. no column  $\mathbf{Z}_j \in \mathcal{L}_i$  for any  $i, j \in [1, \dots, k]$  then
18:          Transfer the set  $\{\mathbf{Z}_j\}_{j=1}^k$  of all columns of  $\mathbf{Z}$  on CPU and add to  $\mathcal{L}_k$ 
19:        end if
20:        go to line 8, setting  $i \rightarrow c \rightarrow 0$ 
21:      end if
22:       $c \rightarrow c + 1$ 
23:    end while
24:     $i \rightarrow i + 1$ 
25:  end while
26: end for; return  $\mathcal{L}$ 
```

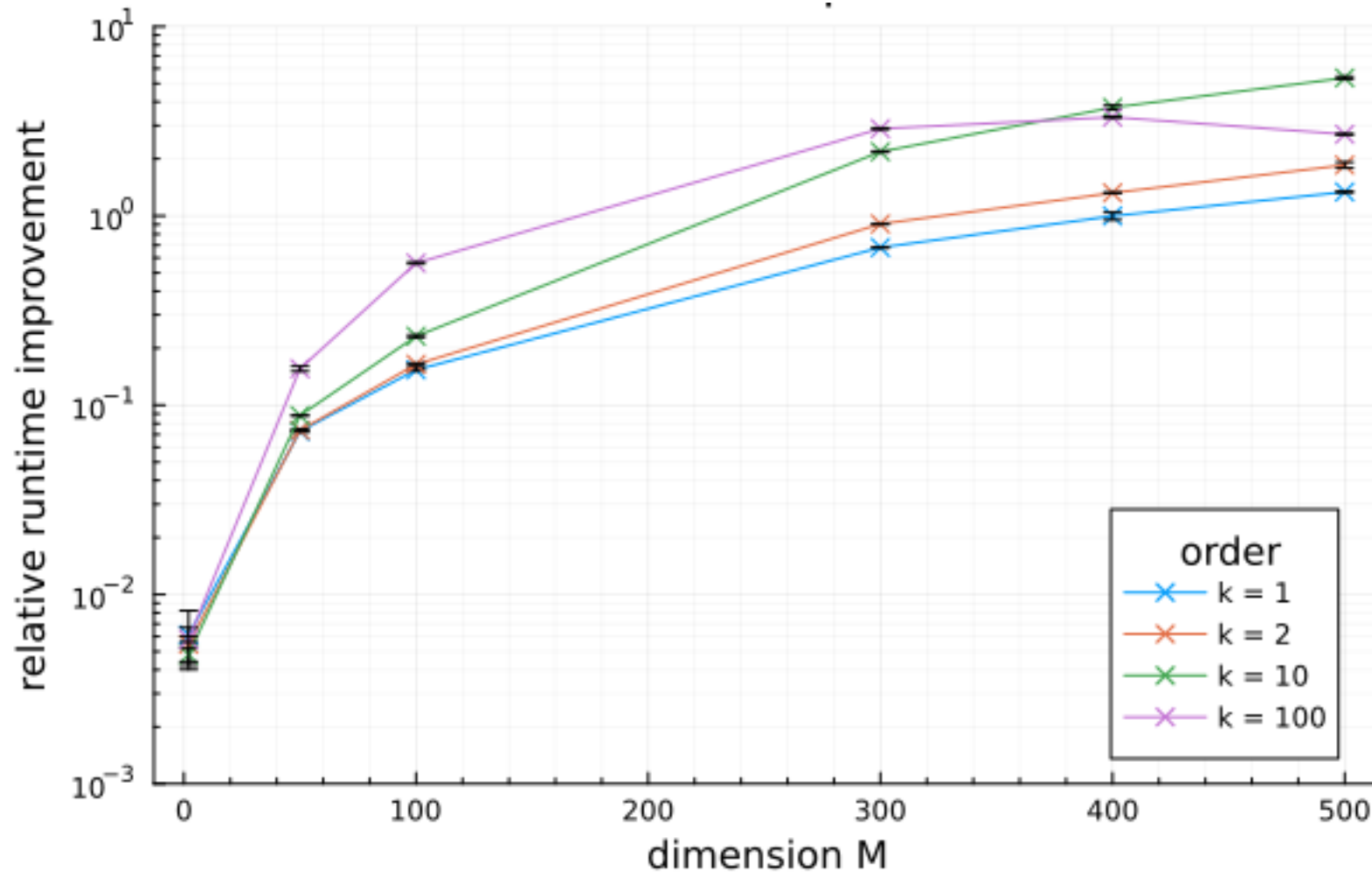
CPU Version

$$\text{CPU improvement} = \frac{\text{runtime of improved CPU version}}{\text{runtime of original version}}$$

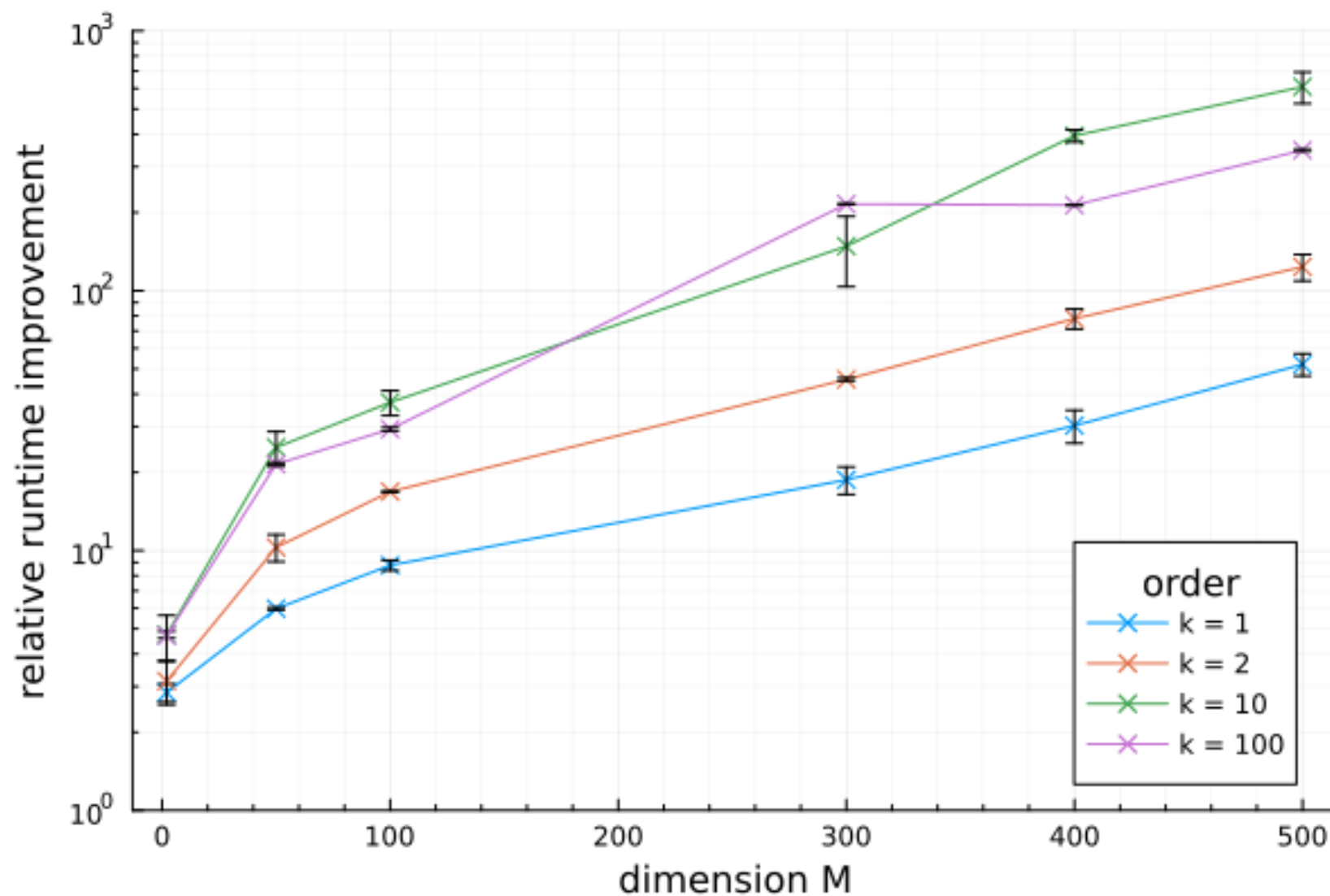


GPU Version

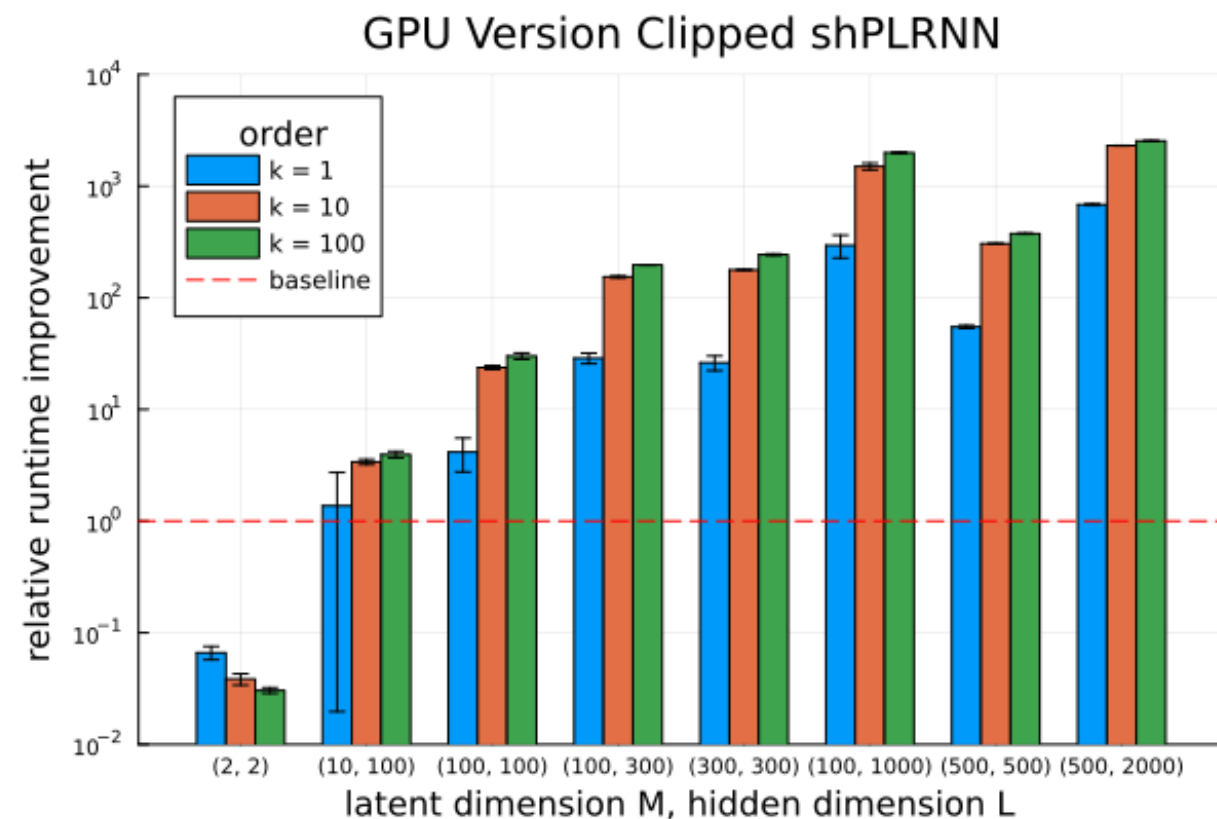
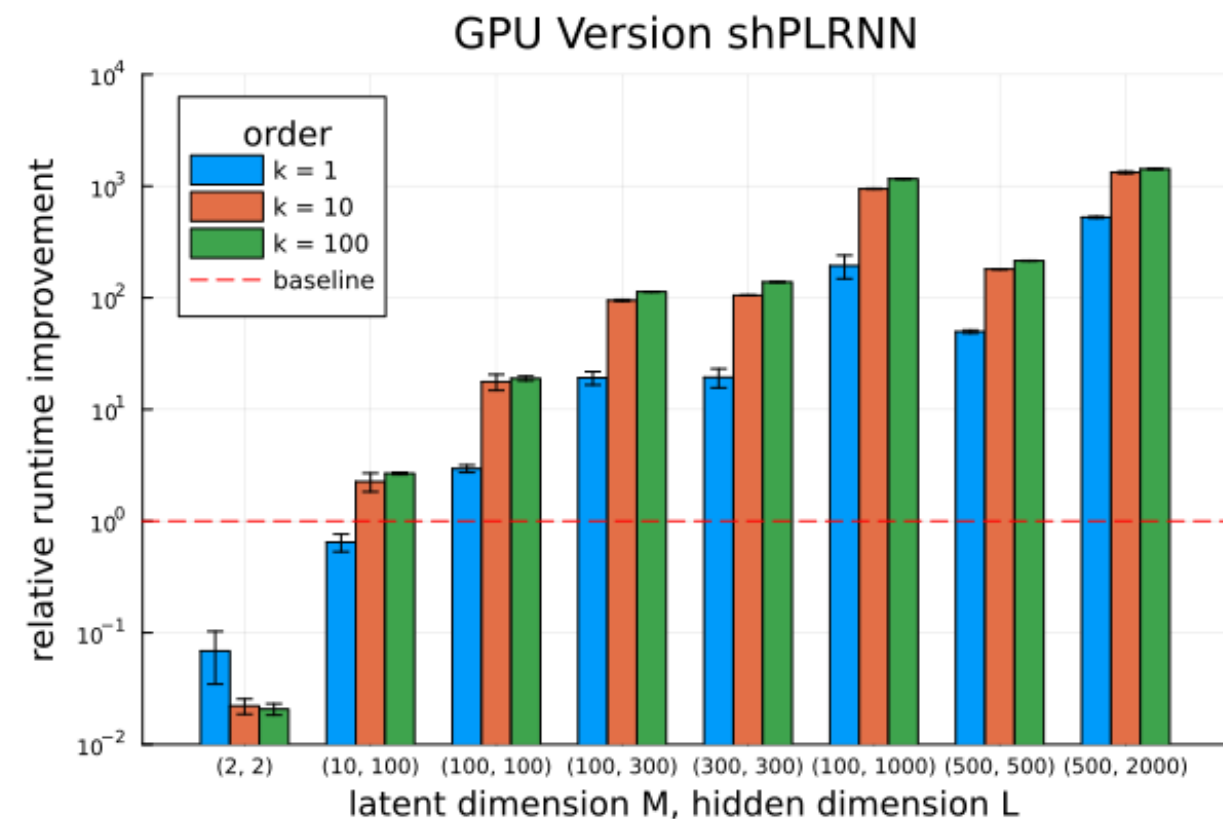
CPU vs. GPU improvement = $\frac{\text{runtime of improved GPU version}}{\text{runtime of improved CPU version}}$ 



Hybrid Version



(Clipped) Shallow PLRNN Version



Summary & More

- Core performance tricks that worked for me:
 - profile first to find actual bottlenecks and after every change
 - prepare CPU program by using improved data structures, preallocation and in-place computations
 - offload matrix multiplication and inversion to GPU but minimize data transfer between host and device
- More Versions:
 - shallow and clipped shallow PLRNN
 - multiple CPU threads + GPUs
 - initialize subregions visited during data reconstruction

Selected References

- [1] Lukas Eisenmann et al. “Bifurcations and loss jumps in RNN training”. In: Advances in Neural Information Processing Systems. Vol. 36. 2023.
- [2] CUDA.jl Documentation. url: <https://cuda.juliagpu.org/stable/>. (accessed: 2024-02-16).
- [3] Pradeep Gupta. The CUDA Refresher Series. url: <https://developer.nvidia.com/blog/tag/cuda-refresher/>. (accessed: 2024-02-16).
- [4] Brendan Gregg. “The Flame Graph: This visualization of software execution is a new necessity for performance profiling and debugging.” In: Queue 14.2 (Mar. 2016), url: <https://doi.org/10.1145/2927299.2927301>.
- [5] BenchmarkTools.jl Documentation. url: <https://juliaci.github.io/BenchmarkTools.jl> (accessed: 2024-02-16).

Questions?