

# Gated Linear Attention Transformers with Hardware-Efficient Training

Songlin Yang<sup>\*1</sup> Bailin Wang<sup>\*1</sup> Yikang Shen<sup>2</sup> Rameswar Panda<sup>2</sup> Yoon Kim<sup>1</sup>

<sup>1</sup>Massachusetts Institute of Technology

<sup>2</sup>MIT-IBM Watson AI Lab

## Summary

Linear attention: removes the softmax in ordinary attention  
 $\Rightarrow$  a linear RNN with matrix-valued hidden states.

|           | Softmax Attention                                                                     | Linear Attention                                                                                             |
|-----------|---------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------|
| Training  | $\mathbf{O} = \text{softmax}((\mathbf{Q}\mathbf{K}^\top) \odot \mathbf{M})\mathbf{V}$ | $\mathbf{O} = ((\mathbf{Q}\mathbf{K}^\top) \odot \mathbf{M})\mathbf{V}$                                      |
| Inference | $\mathbf{o}_t = \sum_{i=1}^t \exp(\mathbf{q}_i \mathbf{k}_i^\top) \mathbf{v}_i$       | $\mathbf{S}_t = \mathbf{S}_{t-1} + \mathbf{k}_t^\top \mathbf{v}_t, \mathbf{o}_t = \mathbf{q}_t \mathbf{S}_t$ |

Issues:

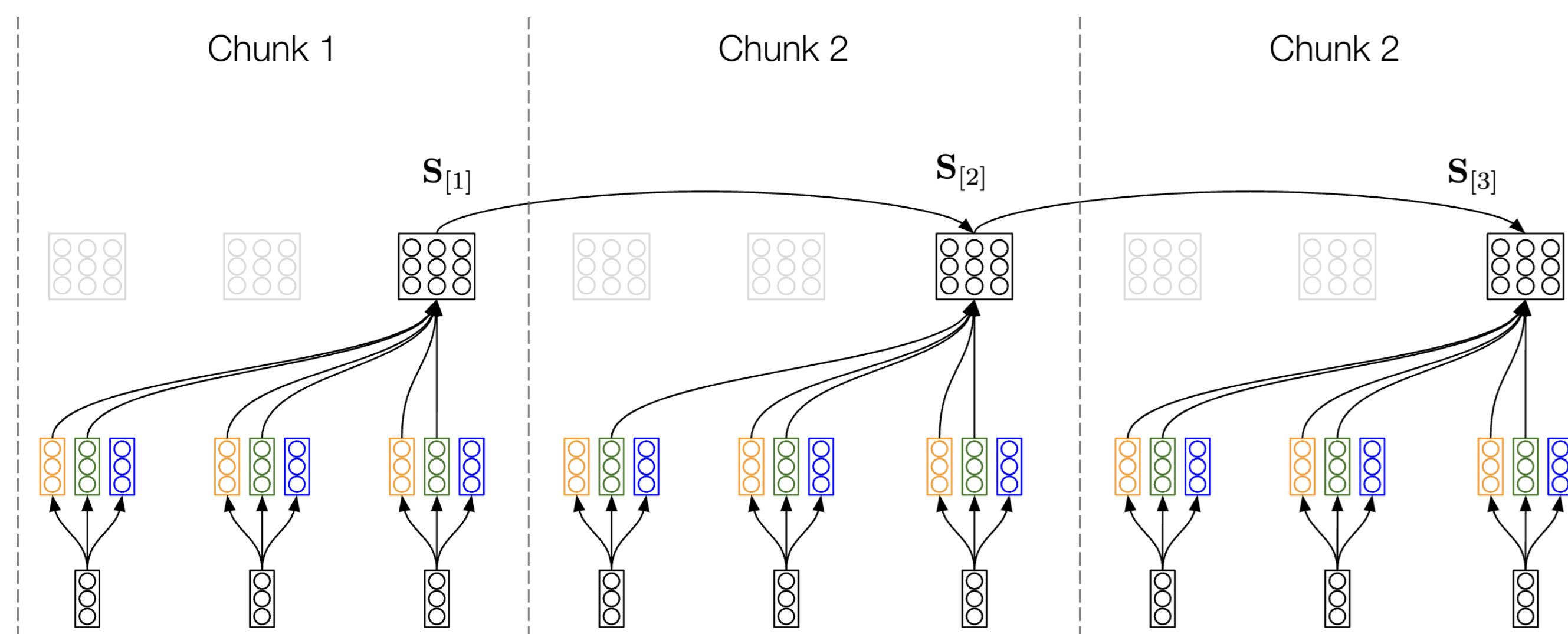
- Slow wall time training speed compared to FlashAttention.
- Poor language modeling performance.

## Our contributions

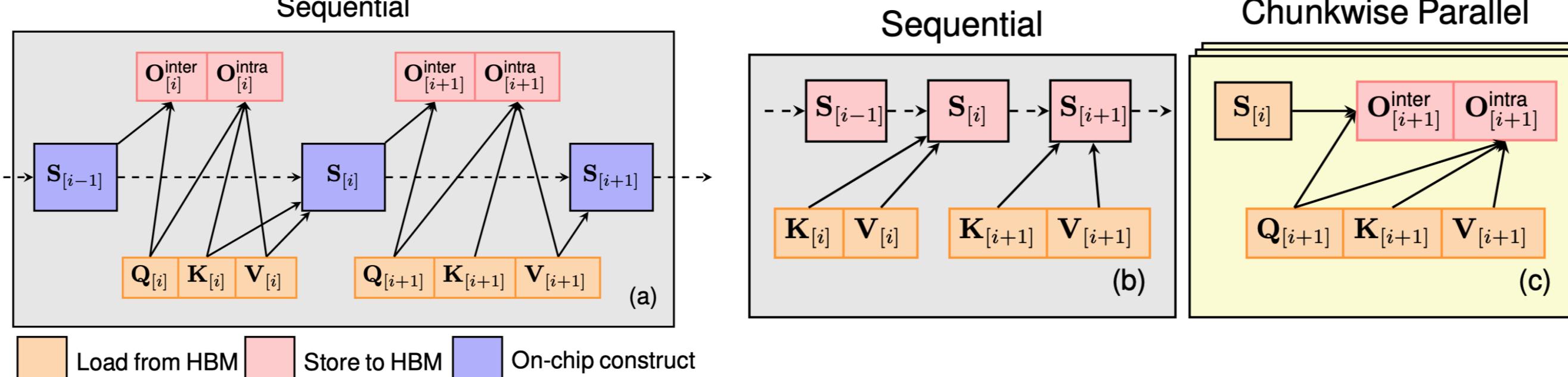
- **FlashLinearAttention**: a hardware-efficient linear attention implementation library.
- **Gated Linear Attention**: improve language modeling performance through a data-dependent gating mechanism.

## Three Forms of Linear Attention

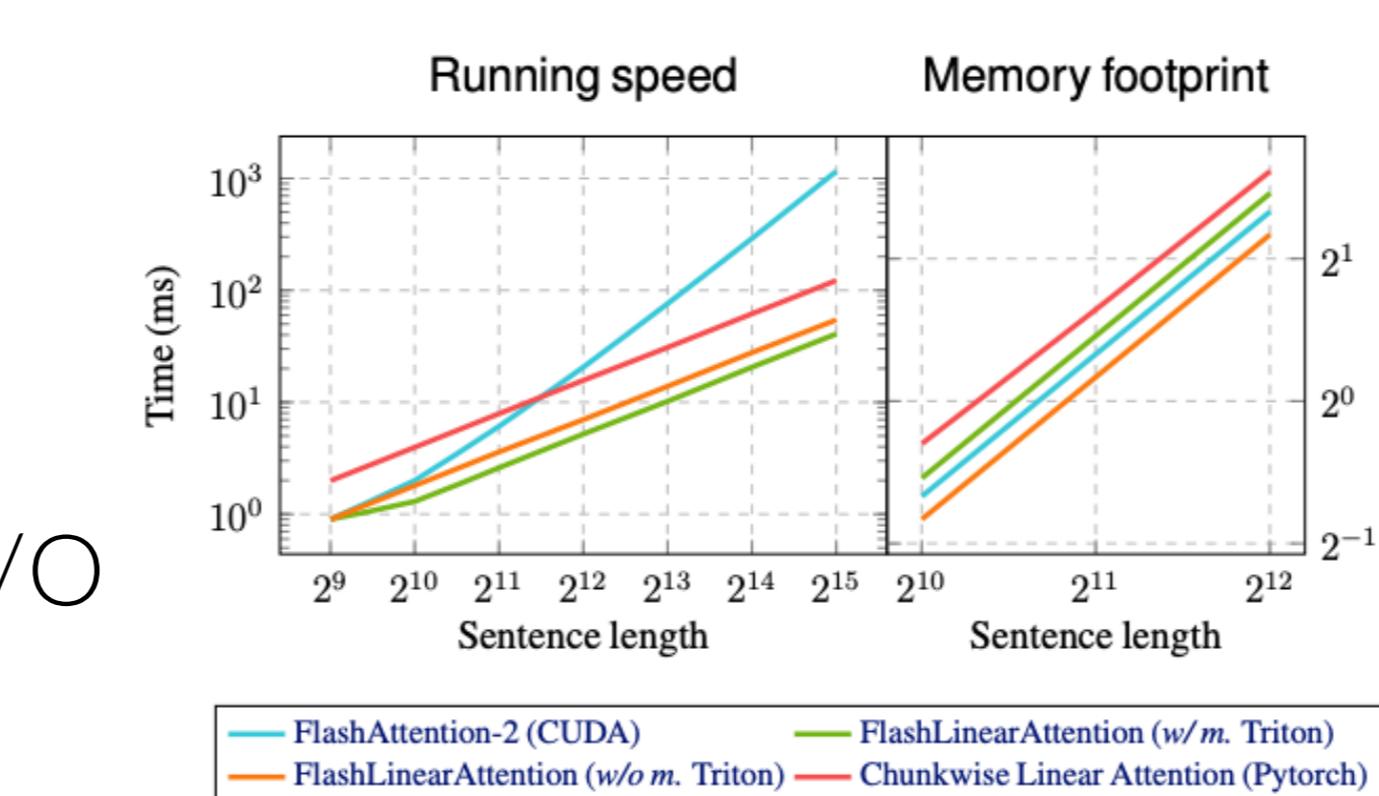
|           | Equation                                                                                                                                                                                                                            | Linear scaling         | Tensor cores | Sequence parallel |
|-----------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------|--------------|-------------------|
| Parallel  | $\mathbf{O} = ((\mathbf{Q}\mathbf{K}^\top) \odot \mathbf{M})\mathbf{V}$                                                                                                                                                             | No, $O(L^2D)$          | Yes          | Yes               |
| Recurrent | $\mathbf{S}_t = \mathbf{S}_{t-1} + \mathbf{k}_t^\top \mathbf{v}_t$<br>$\mathbf{o}_t = \mathbf{q}_t \mathbf{S}_t$                                                                                                                    | Yes, $O(LD^2)$         | No           | No                |
| Chunkwise | $\mathbf{S}_{[i+1]} = \mathbf{S}_{[i]} + \mathbf{K}_{[i]}^\top \mathbf{V}_{[i]}$<br>$\mathbf{O}_{[i+1]} = \mathbf{Q}_{[i+1]} \mathbf{S}_{[i]} + ((\mathbf{Q}_{[i+1]} \mathbf{K}_{[i+1]}^\top) \odot \mathbf{M}) \mathbf{V}_{[i+1]}$ | Yes<br>$O(LCD + LD^2)$ | Yes,<br>Yes  |                   |



## FlashLinearAttention: Efficient Linear Attention



- (a): minimal I/O cost, restricted parallelism
- (b-c): high chunk-level parallelism, slightly higher I/O cost.



## Gated Linear Attention

Introducing 2D forget gate  $\mathbf{G}_t \in \mathbb{R}^{d \times d}$  to linear attention:

$$\mathbf{S}_t = \mathbf{G}_t \odot \mathbf{S}_{t-1} + \mathbf{k}_t^\top \mathbf{v}_t$$

Different parameterization on  $\mathbf{G}_t$  leads to different models:

| Model                          | Parameterization                                                                                                                                                          | Paramet                                        |
|--------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------|
| Mamba [Gu & Dao 2023]          | $\mathbf{G}_t = \exp(-(\mathbf{1}\mathbf{a}_t^\top) \odot \exp(\mathbf{A}))$ , $\mathbf{a}_t = \text{softplus}(\mathbf{x}_t \mathbf{W}_{\alpha_1} \mathbf{W}_{\alpha_2})$ | $\mathbf{A}, \mathbf{W}_{\alpha_1},$           |
| Mamba-2 [Dao & Gu 2024]        | $\mathbf{G}_t = \gamma_t \mathbf{1} \mathbf{1}^\top$ , $\gamma_t = \exp(-\text{softplus}(\mathbf{x}_t \mathbf{W}_\gamma) \exp(a))$                                        | $\mathbf{W}_\gamma, a$                         |
| xLSTM [Beck et al. 2024]       | $\mathbf{G}_t = \gamma_t \mathbf{1} \mathbf{1}^\top$ , $\gamma_t = \sigma(\mathbf{x}_t \mathbf{W}_\gamma)$                                                                | $\mathbf{W}_\gamma$                            |
| GLA [Yang et al. 2023]         | $\mathbf{G}_t = \alpha_t \mathbf{1}^\top$ , $\alpha_t = \sigma(\mathbf{x}_t \mathbf{W}_{\alpha_1} \mathbf{W}_{\alpha_2})^{\frac{1}{2}}$                                   | $\mathbf{W}_{\alpha_1}, \mathbf{W}_{\alpha_2}$ |
| Gated RetNet [Sun et al. 2024] | $\mathbf{G}_t = \gamma_t \mathbf{1} \mathbf{1}^\top$ , $\gamma_t = \sigma(\mathbf{x}_t \mathbf{W}_\gamma)^{\frac{1}{2}}$                                                  | $\mathbf{W}_\gamma$                            |
| HGRN-2 [Qin et al. 2024]       | $\mathbf{G}_t = \alpha_t \mathbf{1}^\top$ , $\alpha_t = \gamma + (1 - \gamma) \sigma(\mathbf{x}_t \mathbf{W}_\alpha)$                                                     | $\mathbf{W}_\alpha, \gamma$                    |
| RWKV-6 [Peng et al. 2024]      | $\mathbf{G}_t = \alpha_t \mathbf{1}^\top$ , $\alpha_t = \exp(-\exp(\mathbf{x}_t \mathbf{W}_\alpha))$                                                                      | $\mathbf{W}_\alpha$                            |
| Gated RFA [Peng et al. 2021]   | $\mathbf{G}_t = \gamma_t \mathbf{1} \mathbf{1}^\top$ , $\gamma_t = \sigma(\mathbf{x}_t \mathbf{W}_\gamma)$                                                                | $\mathbf{W}_\gamma$                            |
| Decaying FW [Mao et al. 2022]  | $\mathbf{G}_t = \alpha_t \beta_t^\top$ , $\alpha_t = \sigma(\mathbf{x}_t \mathbf{W}_\alpha)$ , $\beta_t = \sigma(\mathbf{x}_t \mathbf{W}_\beta)$                          | $\mathbf{W}_\alpha, \mathbf{W}_\beta$          |

## Gated Linear Attention $\subset$ State-Space Models

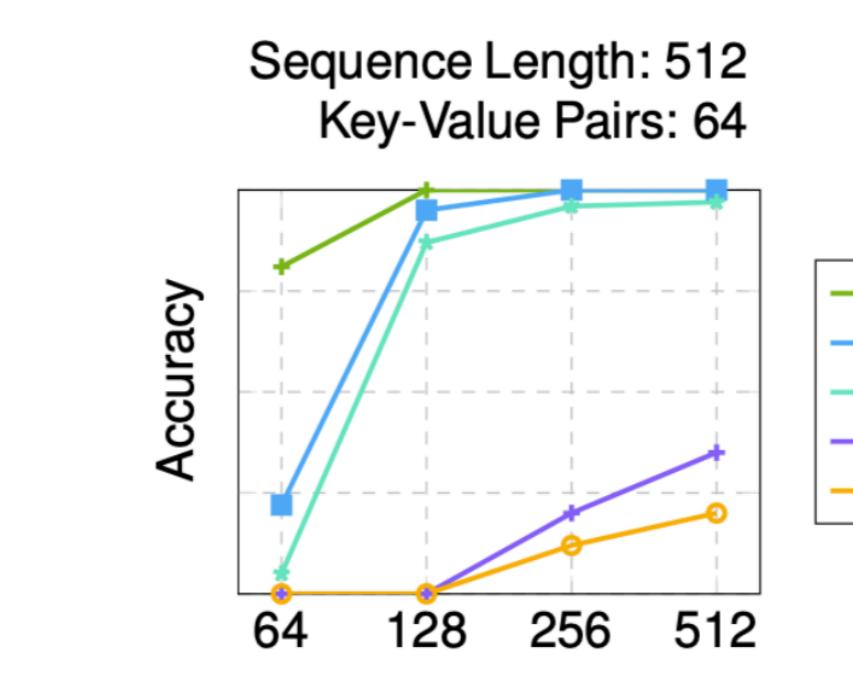
GLA's chunkwise parallel form and fast Triton kernel:

- Support efficient scaling of hidden state size by leveraging tensor cores.
- Facilitate training of recent models like HGRN-2, RWKV-6, Mamba-2.

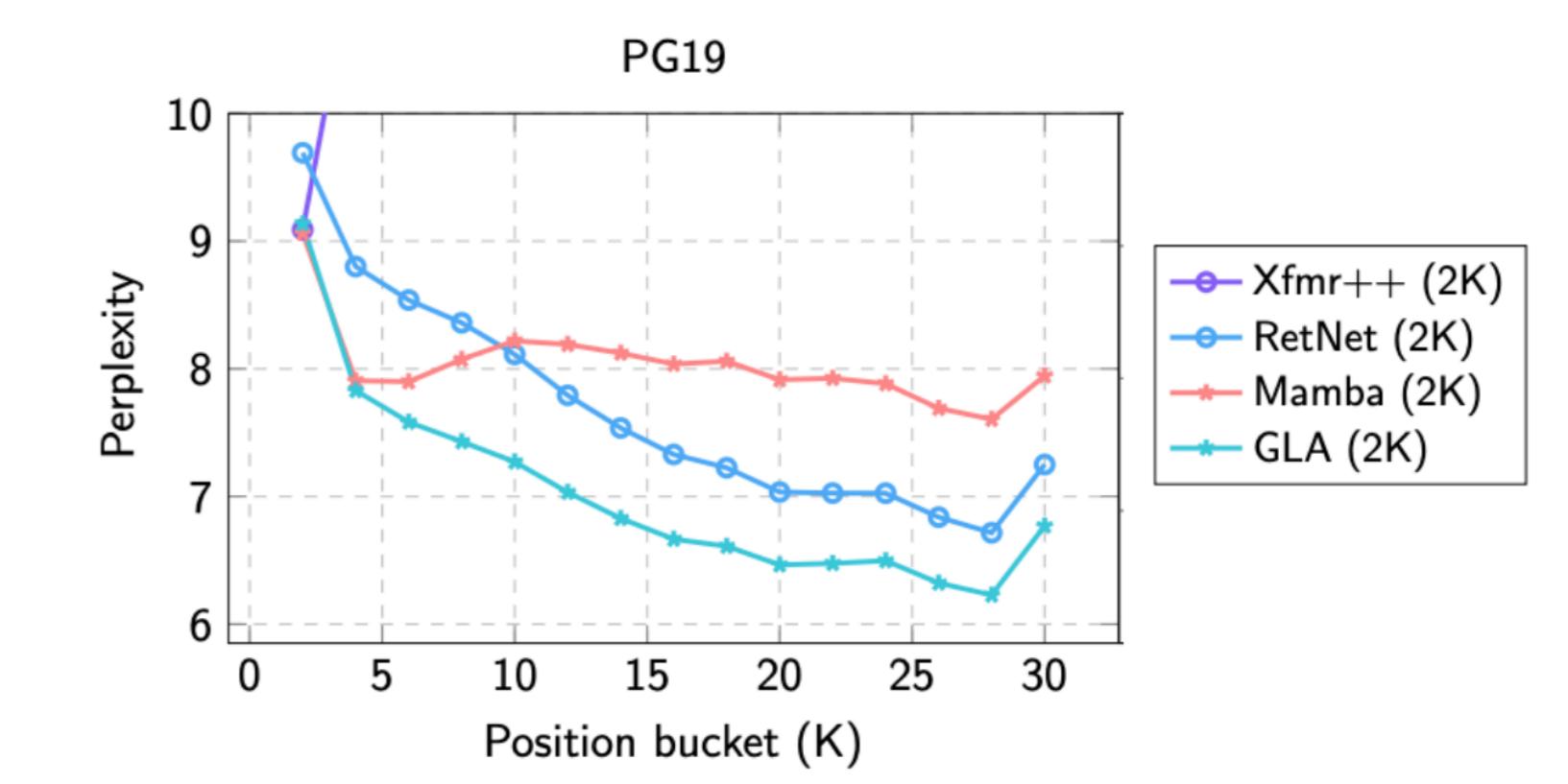
## Performance

| Scale                      | Model         | Wiki. ppl $\downarrow$ | LM Eval. acc. $\uparrow$ | Recall Tasks FDA SWD SQD |
|----------------------------|---------------|------------------------|--------------------------|--------------------------|
| 340M Params<br>15B Tokens  | Transformer++ | 28.39                  | 41.2                     | 21.4 42.2 22.1           |
|                            | RetNet        | 32.33                  | 41.0                     | 2.9 13.3 27.6            |
|                            | Mamba         | 28.39                  | 41.8                     | 2.1 12.4 23.0            |
|                            | GLA           | 28.65                  | 41.5                     | 8.1 18.6 27.2            |
| 1.3B Params<br>100B Tokens | Transformer++ | 16.85                  | 50.9                     | 21.4 42.2 22.1           |
|                            | RetNet        | 18.64                  | 48.9                     | 14.3 42.8 34.7           |
|                            | Mamba         | 17.06                  | 50.0                     | 6.2 41.4 35.2            |
|                            | GLA           | 17.22                  | 51.0                     | 19.9 50.6 42.6           |

## MQAR



## Length extrapolation



## Training Speed / Memory

