



# 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	$O = \text{softmax}((QK^T) \odot M)V$	$O = ((QK^T) \odot M)V$
Inference	$o_t = \frac{\sum_{i=1}^t \exp(q_t k_i^T) v_i}{\sum_{i=1}^t \exp(q_t k_i^T)}$	$S_t = S_{t-1} + k_t^T v_t, o_t = q_t S_t$

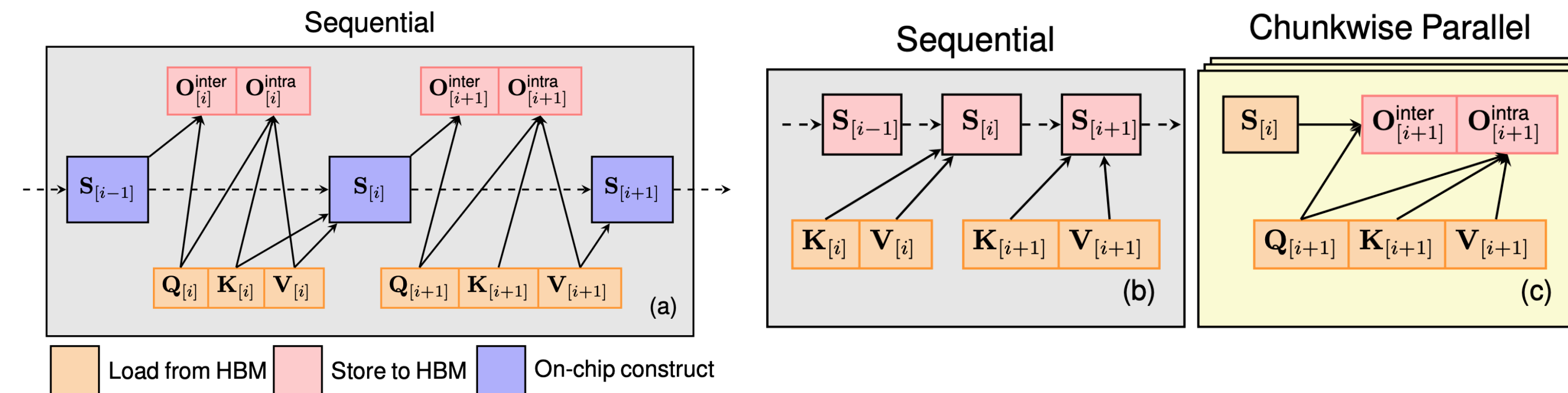
Issues:

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

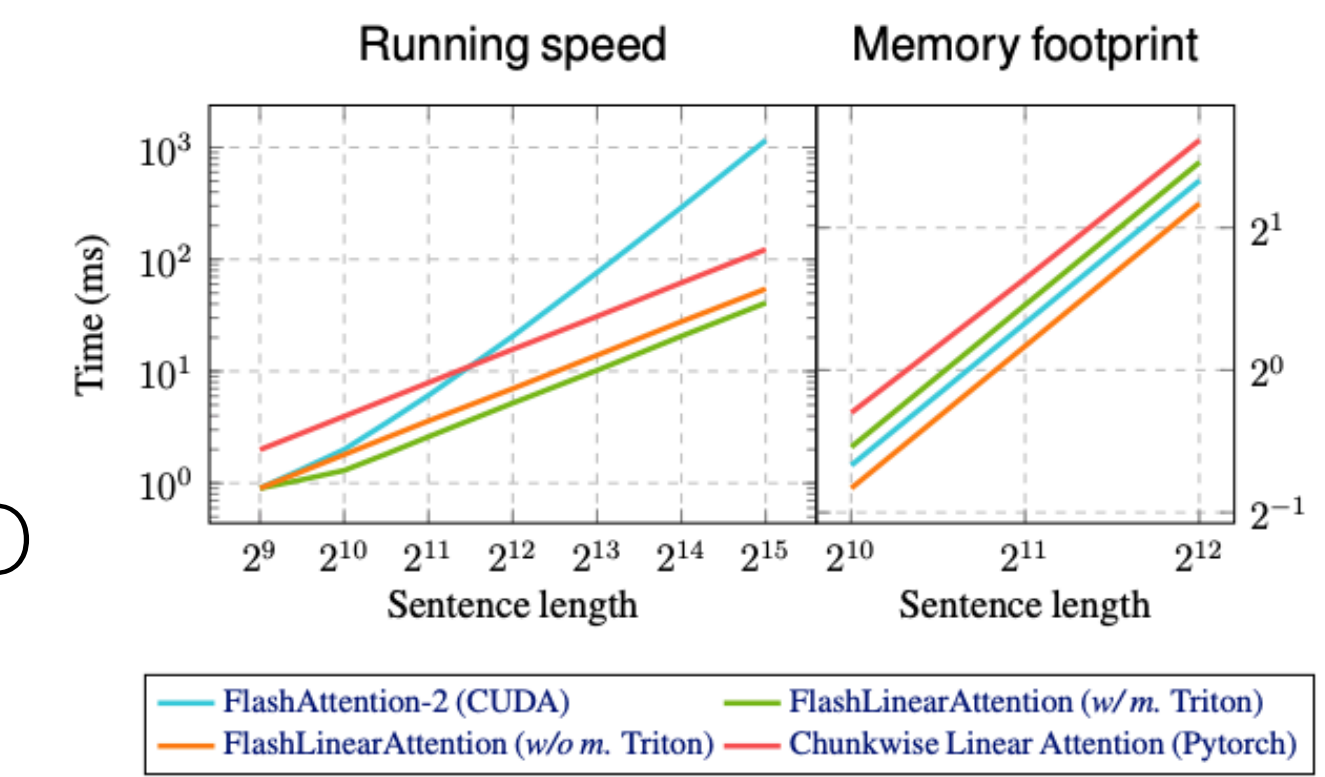
Our contributions

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

## 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  $G_t \in \mathbb{R}^{d \times d}$  to linear attention:

$$S_t = G_t \odot S_{t-1} + k_t^T v_t$$

Different parameterization on  $G_t$  leads to different models:

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

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

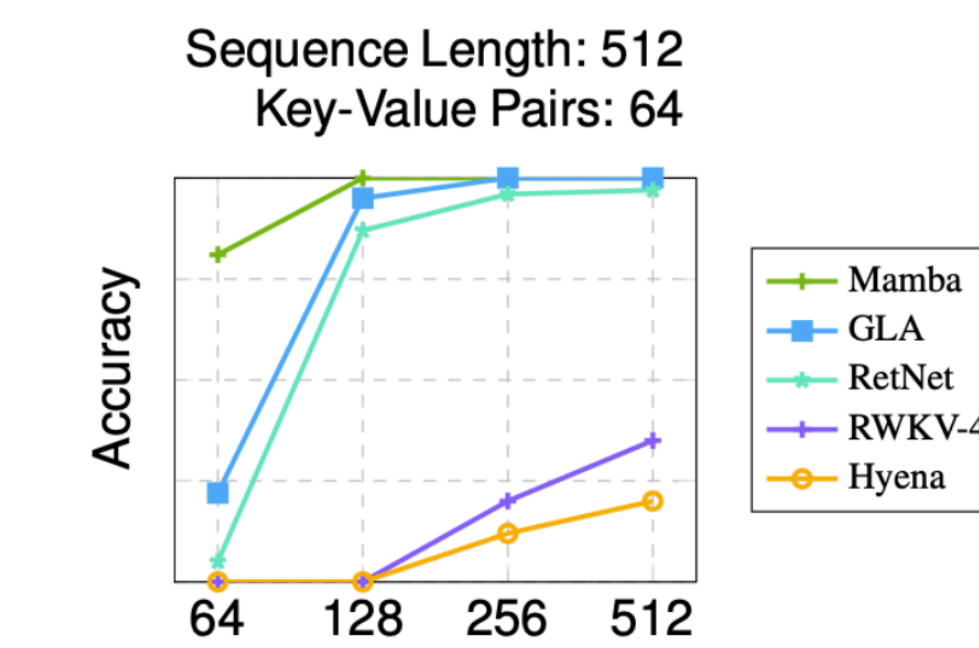
GLA's chunkwise parallel form and fast Triton kernel:

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

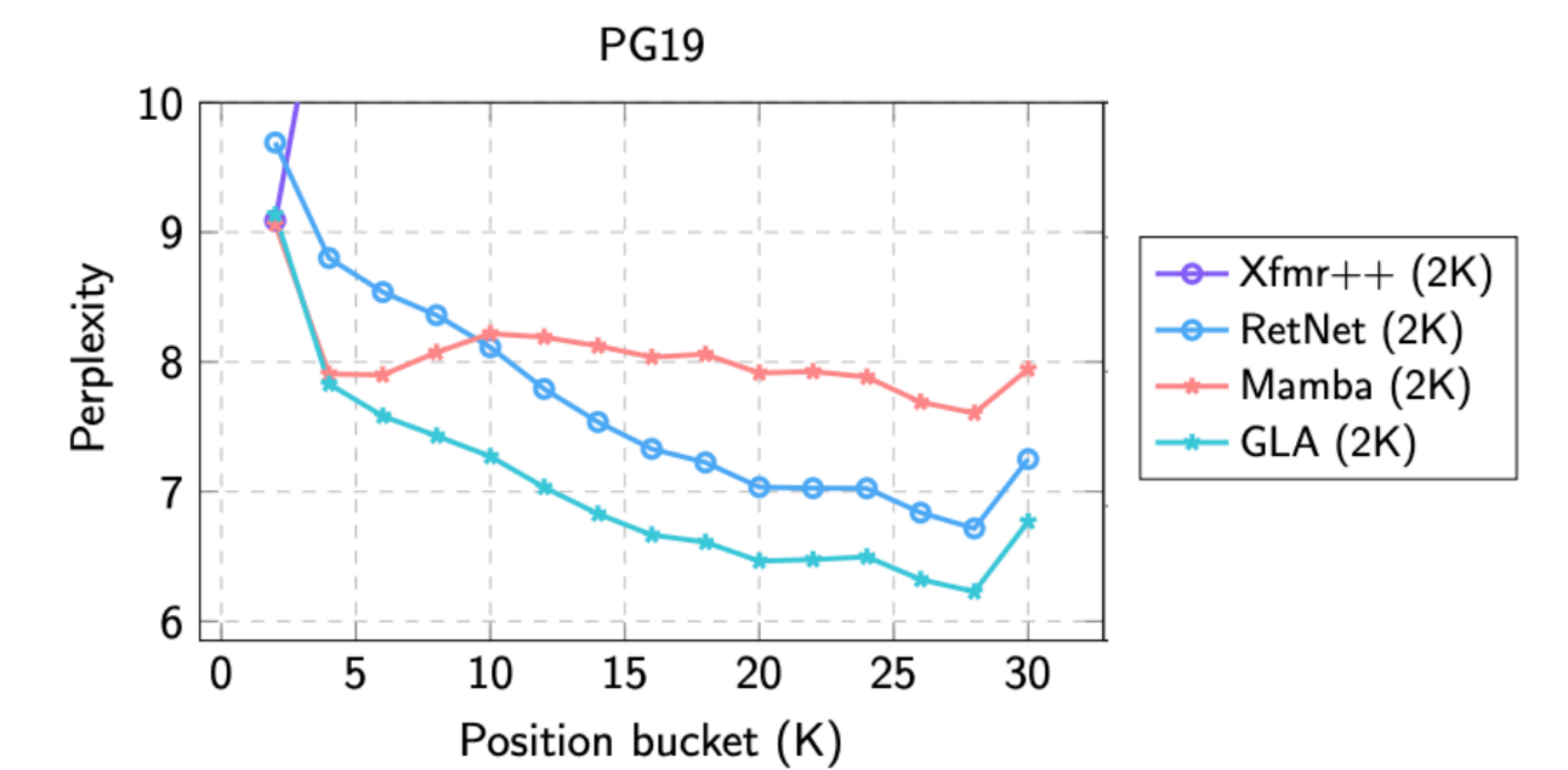
## Performance

Scale	Model	Wiki. ppl ↓	LM Eval. acc. ↑	Recall Tasks		
				FDA	SWD	SQD
340M Params 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 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

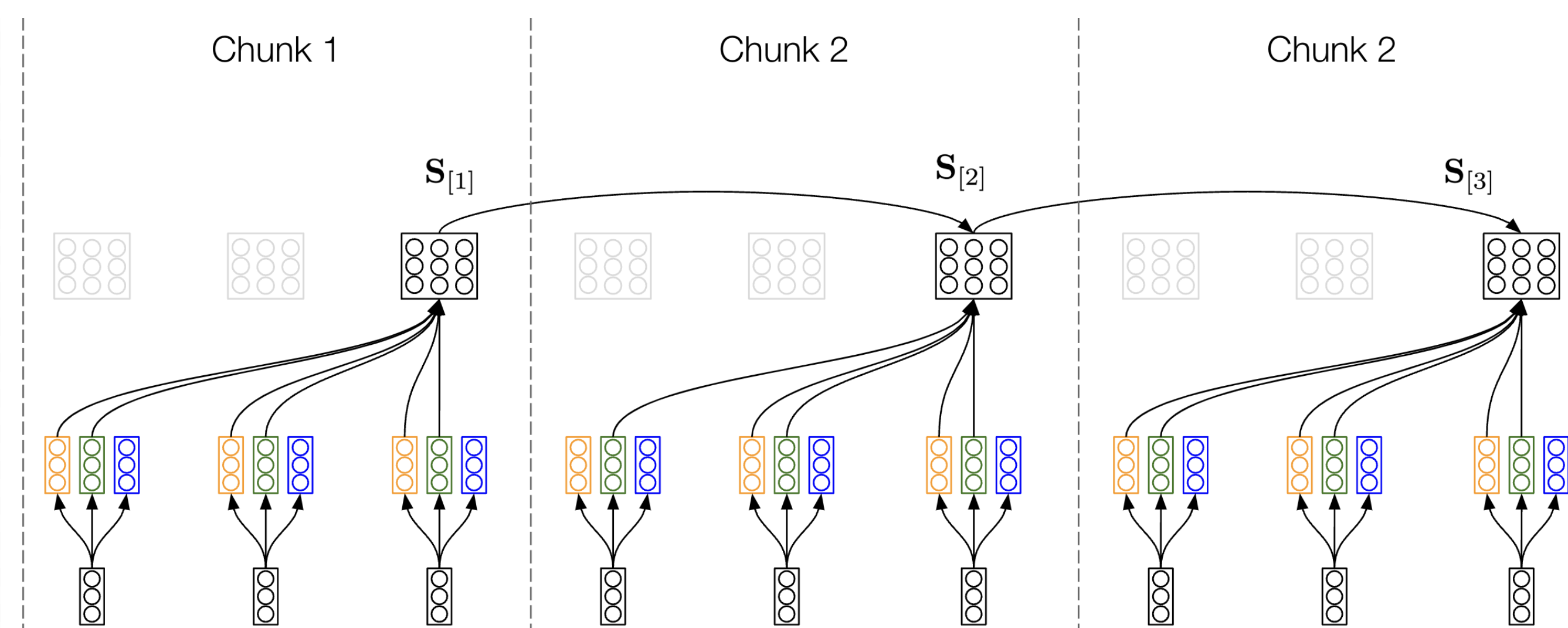


## Length extrapolation



## Three Forms of Linear Attention

	Equation	Linear scaling	Tensor cores	Sequence parallel
Parallel	$O = ((QK^T) \odot M)V$	No, $O(L^2D)$	Yes	Yes
Recurrent	$S_t = S_{t-1} + k_t^T v_t$ $o_t = q_t S_t$	Yes, $O(LD^2)$	No	No
Chunkwise	$S_{[i+1]} = S_{[i]} + K_{[i]}^T V_{[i]}$ $O_{[i+1]} = Q_{[i+1]} S_{[i]} + ((Q_{[i+1]} K_{[i+1]}^T) \odot M) V_{[i+1]}$	Yes $O(LCD + LD^2)$	Yes,	Yes



## Training Speed / Memory

