

# Tiled Flash Linear Attention

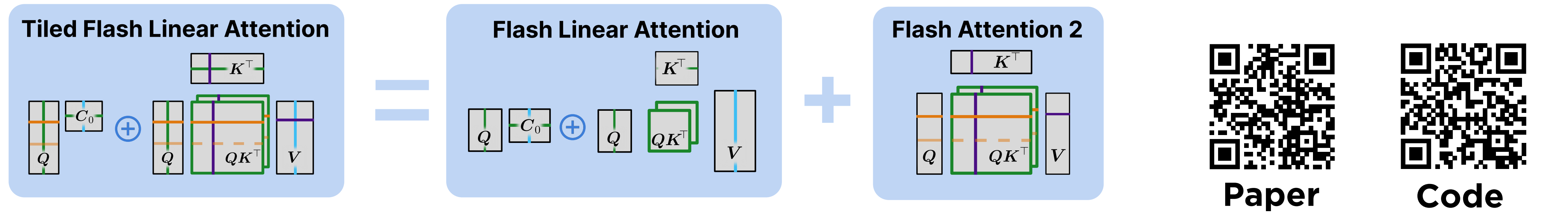
## More Efficient Linear RNN and xLSTM Kernels

Maximilian Beck<sup>1,2</sup>, Korbinian Pöppel<sup>1,2</sup>, Phillip Lippe<sup>2,3</sup>, Sepp Hochreiter<sup>1,2</sup>

1) ELLIS Unit JKU Linz, Institute for Machine Learning 2) NXAI GmbH 3) Now at Google Deepmind

ELLIS  
European Laboratory for Learning and Intelligent Systems

JKU  
NXAI

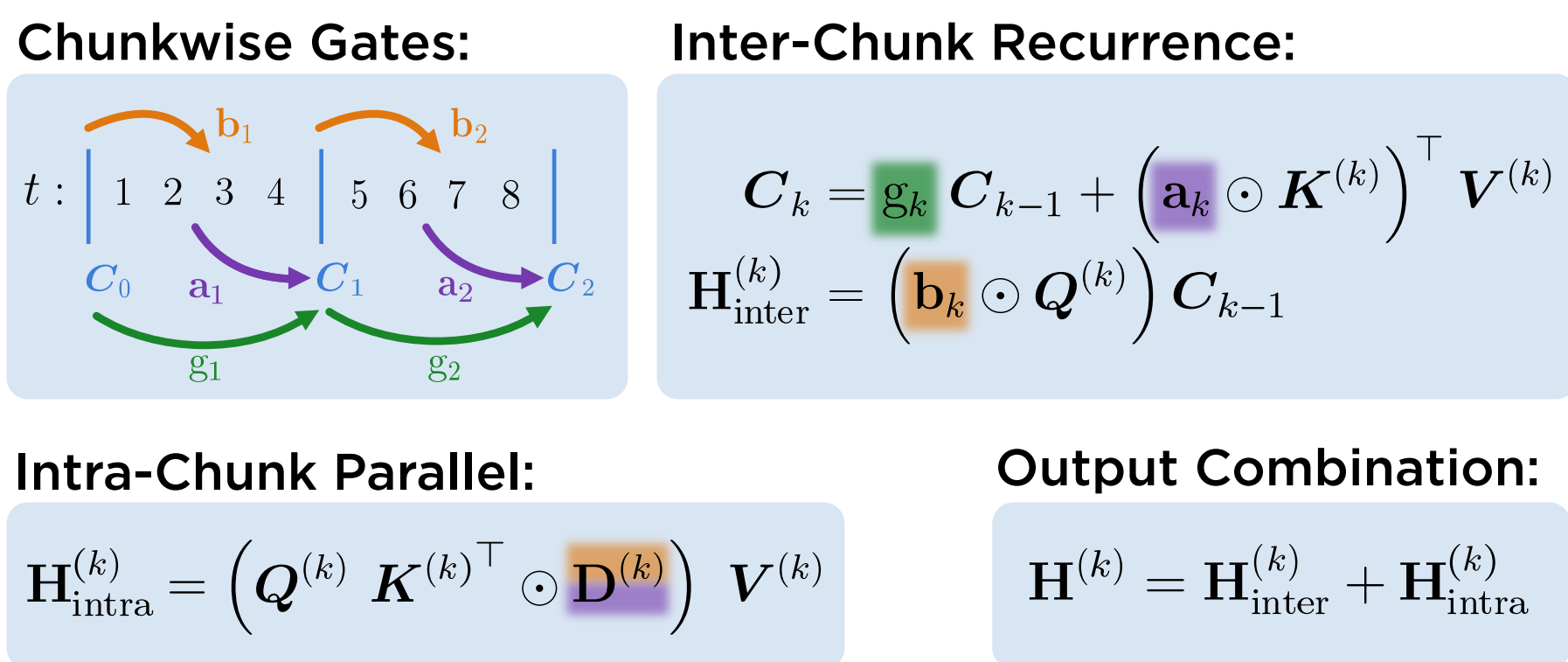


**TL;DR: We combine Flash Linear Attention with Flash Attention 2**

### Motivation

- Gated Linear RNNs become competitive on language modeling: RetNet, Mamba, GLA, xLSTM/mLSTM
- Linear RNNs have a chunkwise formulation which computes intermediate memory states and enables efficient implementations
  - scales linearly with sequence length
- FLA kernels leverage the chunkwise formulation and are faster than Flash Attention kernels
- BUT: FLA is limited in the chunksize by available SRAM on the GPU**
  - Many memory states are materialized in memory, which causes high memory usage & IO and low arithmetic intensity

### Chunkwise-parallel formulation in 4 parts



### TFLA Forward Pass Tiling

Simplified form of the parallel intra-chunk forward pass for chunk  $k$ :

$$H^{(k)} = \underbrace{Q^{(k)}}_{(L_{hq} \times d_{hq})} \underbrace{C_{k-1}}_{(d_{qk} \times d_{hv})} + \left( \underbrace{Q^{(k)}}_{(L_{hq} \times d_{qk})} \underbrace{K^{(k)\top}}_{(d_{qk} \times L_{kv})} \right) \underbrace{V^{(k)}}_{(L_{kv} \times d_{hv})}$$

chunk size dimensions      head dimensions

We parallelize along  $L_{hq}$  hidden-query

We loop over  $L_{kv}$  key-value

$d_{hv}$  hidden-value

$d_{qk}$  query-key

### We apply TFLA to 2 mLSTM variants

#### mLSTMexp

$$m_t = \max \left\{ \log \sigma(\tilde{f}_t) + m_{t-1}, \tilde{i}_t \right\}$$

$$C_t = f_t C_{t-1} + \exp(\tilde{i}_t - m_t) k_t v_t^\top$$

$$n_t = f_t n_{t-1} + \exp(\tilde{i}_t - m_t) k_t$$

$$\tilde{h}_t = \frac{C_t^\top q_t}{\max \left\{ |n_t^\top q_t|, \exp(-m_t) \right\}}$$

$$h_t = o_t \odot \text{NORM}(\tilde{h}_t)$$

#### mLSTMsig

$$C_t = f_t C_{t-1} + \sigma(\tilde{i}_t) k_t v_t^\top$$

$$\tilde{h}_t = C_t^\top q_t$$

$$h_t = o_t \odot \text{NORM}(\tilde{h}_t)$$

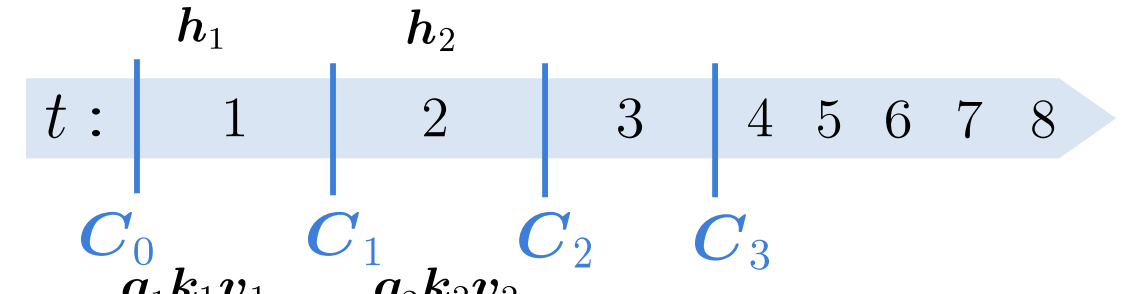
- Exponential input gate
- Additional max and normalizer state for stabilization
- Tested at 7B parameter scales

- Sigmoid input gate
- No max and normalizer state
- Less FLOPS & faster kernels
- Equal performance up to 1.3B parameters

### Background: Linear RNN Formulations

#### Recurrent:

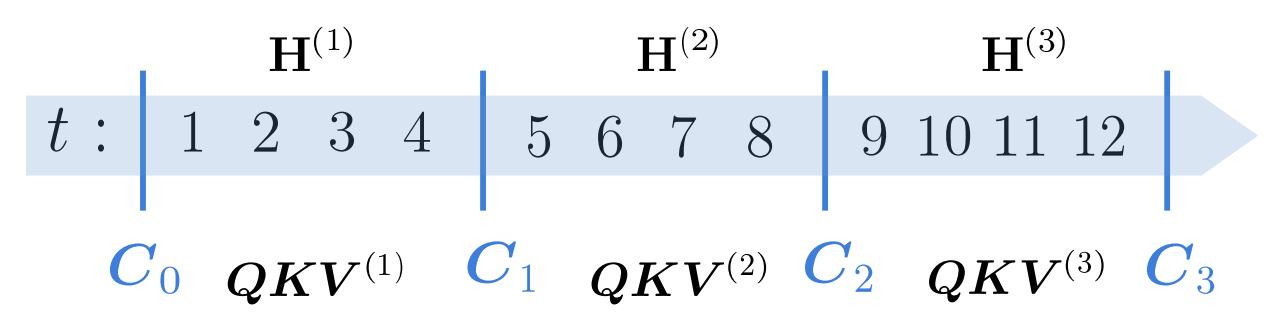
- C state in every time step
- Linear compute



#### Chunkwise-parallel:

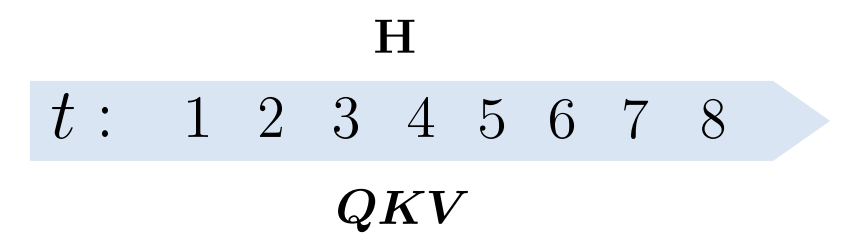
- C state every L-th time step
- Linear compute

→ between recurrent and parallel

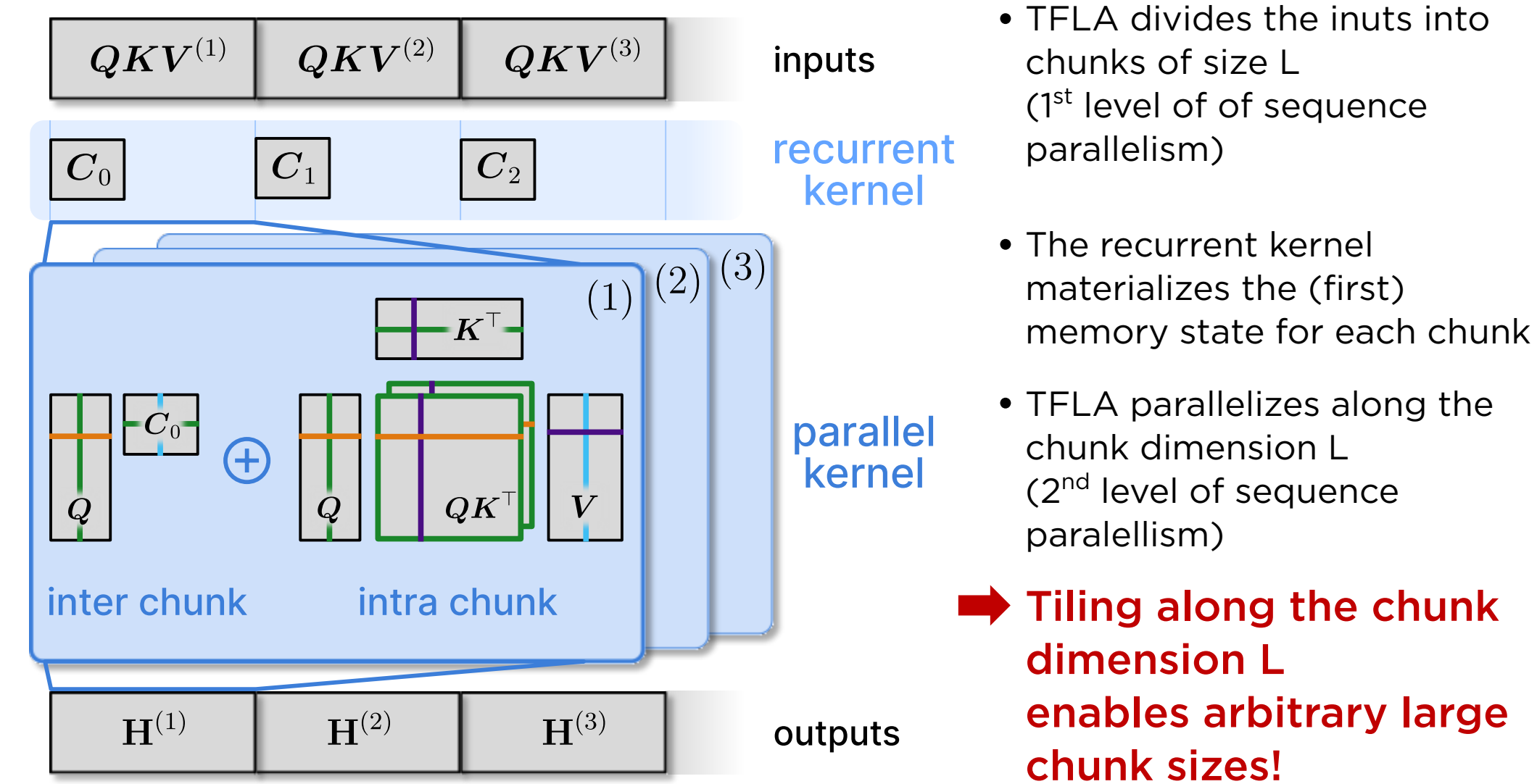


#### Parallel:

- No C states
- Quadratic compute!

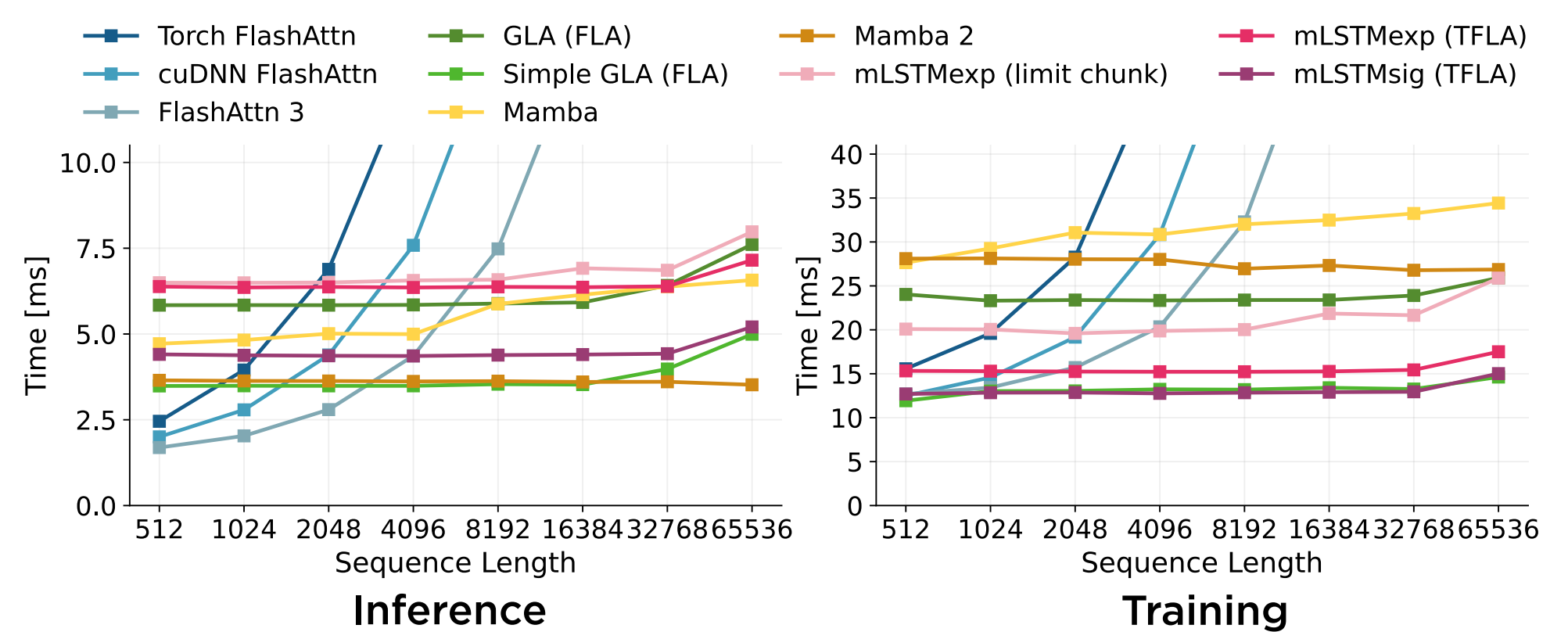


### TFLA Kernels Overview



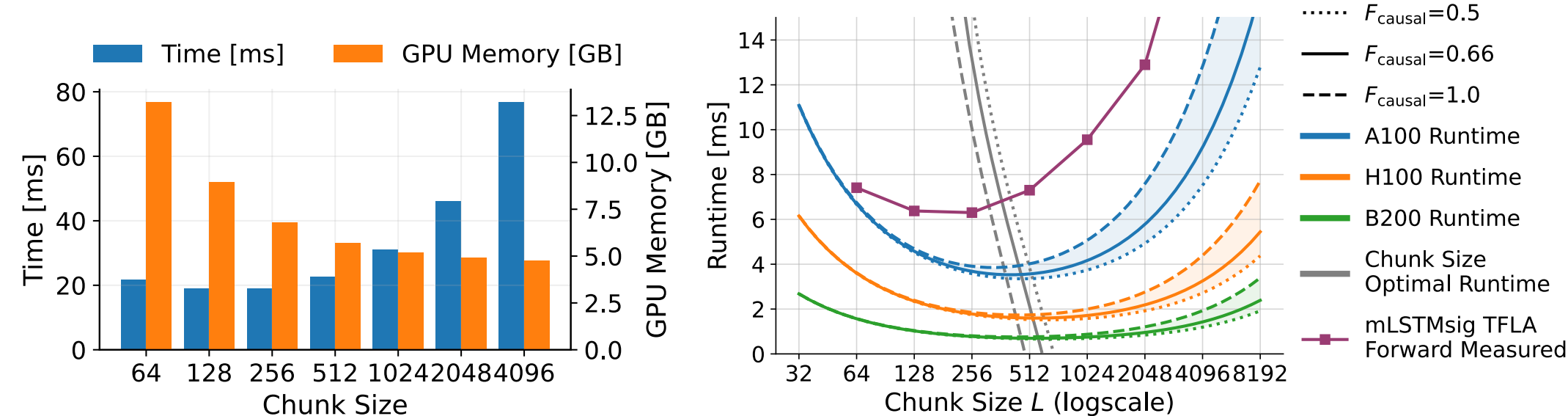
### Benchmark Results

Constant number of 65k tokens with embedding dim 4096



→ TFLA mLSTM kernels are faster than FA3 & 2x faster than Mamba 2

### Chunk size is more than a kernel parameter



#### Runtime vs. Memory Trade-off

The chunk size...

- ...enables trade-off between memory & runtime
- ...interpolates the FLOPs between recurrent and parallel formulation
- ...determines the arithmetic intensity (i.e. compute vs. memory bound)