# xLSTM
## Extended Long Short-Term Memory

Maximilian Beck*[1,2,3], Korbinian Pöppel*[1,2,3], Markus Spanring[1], Andreas Auer[1,2], Oleksandra Prudnikova
Michael Kopp, Günter Klambauer[1,2,3], Johannes Brandstetter[1,2,3], Sepp Hochreiter[1,2,3]
1) Johannes Kepler University Linz, Institute for Machine Learning 2) NXAI Lab Linz 3) NXAI GmbH

## Can LSTMs be scaled to billions of parameters while matching Transformer's capabilities?

Limitations of the LSTM:

### LSTM's inability to revise storage decisions

- Sigmoid input gate is limited → cannot overwrite
- Replace by *exponential input gate*
- Introduce normalizer $\boldsymbol{n}_t$ to re-stabilize

$$c_t = \sigma(\tilde{\mathbf{f}}_t) \odot c_{t-1} + \sigma(\tilde{\mathbf{i}}_t) \odot \tanh(\tilde{z}_t)$$
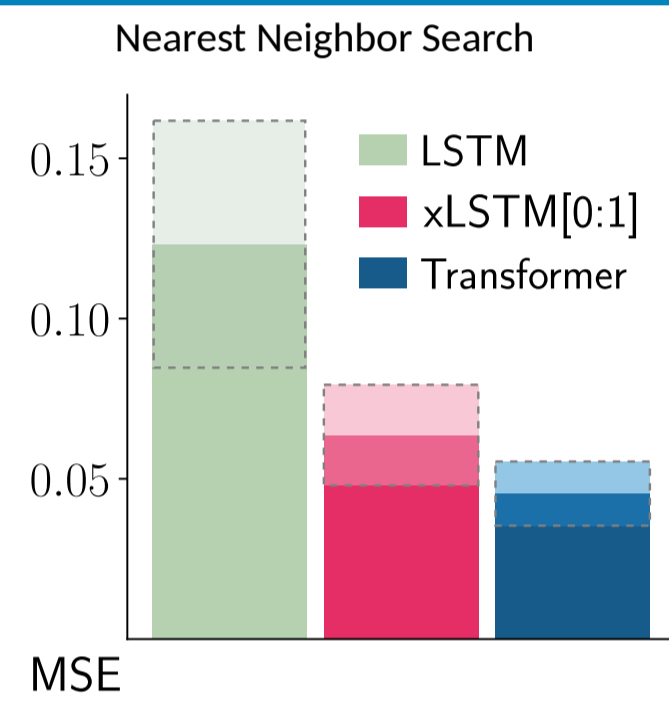$$h_t = \sigma(\tilde{\mathbf{o}}_t) \odot \tanh(c_t)$$

**LSTM**

$$c_t = \sigma(\tilde{\mathbf{f}}_t) \odot c_{t-1} + \exp(\tilde{\mathbf{i}}_t) \odot \tanh(\tilde{z}_t)$$
$$n_t = \sigma(\tilde{\mathbf{f}}_t) \odot n_{t-1} + \exp(\tilde{\mathbf{i}}_t)$$
$$h_t = \sigma(\tilde{\mathbf{o}}_t) \odot c_t / n_t$$

**sLSTM**

Nearest Neighbor Search



Legend: LSTM, xLSTM[0:1], Transformer; axis: MSE

### LSTM's limited storage capacity

- Scalar memory cells, each gated → limited capacity
- Now *matrix memory cell* with *outer product update*
- Use down-projection to hidden state by query vector

$$c_t = \sigma(\tilde{\mathbf{f}}_t) \odot c_{t-1} + \sigma(\tilde{\mathbf{i}}_t) \odot \tanh(\tilde{z}_t)$$
$$h_t = \sigma(\tilde{\mathbf{o}}_t) \odot \tanh(c_t)$$

**LSTM**

$$C_t = \sigma(\tilde{\mathbf{f}}_t) \odot C_{t-1} + \exp(\tilde{\mathbf{i}}_t) \odot \boldsymbol{v}_t \boldsymbol{k}_t^\top$$
$$n_t = \sigma(\tilde{\mathbf{f}}_t) \odot n_{t-1} + \exp(\tilde{\mathbf{i}}_t) \odot \boldsymbol{k}_t$$
$$h_t = \mathbf{o}_t \odot C_t \boldsymbol{q}_t / \max(|\boldsymbol{n}_t^\top \boldsymbol{q}_t|, 1)$$

**mLSTM**

Wiki103 Rare Token PPL



Legend: LSTM, xLSTM[1:0], Transformer; axis PPL: $<10^3$, $10^3$–$10^4$, $>10^4$, all

### LSTM's inability of time-parallel training

- Recurrent connection limits parallelizability
- *Remove recurrent connection*
- Or: Block-diagonal / multi-head structure for Recurrent Matrix

$$\{\tilde{\mathbf{i}}, \tilde{\mathbf{f}}, \tilde{\mathbf{z}}, \tilde{\mathbf{o}}\}_t = \boldsymbol{W}_{\{i,f,z,o\}} x_t + \boldsymbol{R}_{\{i,f,z,o\}} h_{t-1}$$ **LSTM**

$$\{\boldsymbol{q}, \boldsymbol{k}, \boldsymbol{v}, \tilde{\mathbf{i}}, \tilde{\mathbf{f}}, \tilde{\mathbf{o}}\}_t = \boldsymbol{W}_{\{q,k,v,i,f,o\}} x_t$$ **mLSTM**

$$\{\tilde{\mathbf{i}}, \tilde{\mathbf{f}}, \tilde{\mathbf{z}}, \tilde{\mathbf{o}}\}_t = \boldsymbol{W}_{\{i,f,z,o\}} x_t + \boldsymbol{R}_{\{i,f,z,o\}} h_{t-1}$$ **sLSTM**

Rearranged and time-parallel:



$\tilde{\mathbf{i}}, \tilde{\mathbf{f}}, \tilde{\mathbf{o}} \in \mathbb{R}^d$ or $\mathbb{R}$ input / forget / output gate preactivation
$\tilde{\boldsymbol{z}}_t, \boldsymbol{q}_t, \boldsymbol{k}_t, \boldsymbol{v}_t \in \mathbb{R}^d$ cell input, query, key, value
$\boldsymbol{x}_t \in \mathbb{R}^d$ input $\quad t \in \{1..T\}$ time
$\boldsymbol{W}_{\{i,f,z,o,q,k,v\}} \in \mathbb{R}^{d \times d}$ weight matrix
$c_t \in \mathbb{R}^d, C_t \in \mathbb{R}^{d \times d}$ cell state
$\boldsymbol{n}_t \in \mathbb{R}^d$ normalizer state
$\boldsymbol{h}_t \in \mathbb{R}^d$ hidden state / output
$\boldsymbol{R}_{\{i,f,z,o\}} \in \mathbb{R}^{d \times d}$ recurrent matrix

## OVERVIEW

**LSTM** → **Memory Cells** → **xLSTM Blocks** → **xLSTM**

**LSTM**

Memory Cells
→ Constant Error Carousel
→ Sigmoid Gating
→ Recurrent Inference
→ Recurrent Training

$$c_t = \mathbf{f}_t\, c_{t-1} + \mathbf{i}_t\, z_t$$
$$h_t = \mathbf{o}_t\, \psi(c_t)$$

**Memory Cells**

sLSTM
+ Exponential Gating
+ New Memory Mixing

mLSTM
+ Exponential Gating
+ Matrix Memory
+ Parallel Training
+ Covariance Update Rule



## SEQUENCE AND LANGUAGE MODELING

- xLSTM with sLSTM (memory mixing) can solve formal language tasks
- keeps the state tracking capabilities of LSTM

- xLSTM Model Structure:
  ResNet-like architecture of
  Pre-Layer-Norm blocks
  mLSTM (pre-up projection) block
  sLSTM (post-up projection) block
  Combination: xLSTM[a:b]
  (mLSTM/sLSTM ratio)

|  | Deterministic Context Free | | | Regular | | |
|---|---|---|---|---|---|---|
|  | Mod Arithmetic (w. Brackets) | Solve Equation | Cycle Nav | Even Pairs | Mod Arithmetic (w/o Brackets) | Parity |
| Llama | 0.02 ± 0.0 | 0.02 ± 0.0 | 0.04 ± 0.01 | 1.0 ± 0.0 | 0.03 ± 0.0 | 0.03 ± 0.0 |
| Mamba | 0.04 ± 0.01 | 0.05 ± 0.02 | 0.86 ± 0.04 | 1.0 ± 0.0 | 0.05 ± 0.02 | 0.13 ± 0.0 |
| RWKV-6 | 0.09 ± 0.01 | 0.09 ± 0.02 | 0.31 ± 0.14 | 1.0 ± 0.0 | 0.16 ± 0.0 | 0.22 ± 0.1 |
| LSTM | 0.72 ± 0.04 | 0.38 ± 0.05 | 0.93 ± 0.07 | 1.0 ± 0.0 | 1.0 ± 0.0 | 1.0 ± 0.0 |
| xLSTM[0:1] | 0.57 ± 0.09 | 0.55 ± 0.09 | 1.0 ± 0.0 | 1.0 ± 0.0 | 1.0 ± 0.0 | 1.0 ± 0.0 |
| xLSTM[1:1] | 0.15 ± 0.06 | 0.24 ± 0.04 | 0.8 ± 0.03 | 1.0 ± 0.0 | 0.6 ± 0.4 | 1.0 ± 0.0 |



**mLSTM block**
(pre-up projection)

**sLSTM block**
(post-up projection)

| Model | #Params M | SlimPajama (15B) ppl↓ |
|---|---|---|
| GPT-3 | 356 | 14.26 |
| Llama | 407 | 14.25 |
| H3 | 420 | 18.23 |
| Mamba | 423 | 13.70 |
| Hyena | 435 | 17.59 |
| RWKV-4 | 430 | 15.62 |
| RWKV-5 | 456 | 14.25 |
| RWKV-6 | 442 | 15.03 |
| RetNet | 431 | 16.23 |
| HGRN | 411 | 17.59 |
| GLA | 412 | 16.15 |
| HGRN2 | 411 | 14.32 |
| **xLSTM[1:0]** | **409** | **13.43** |
| **xLSTM[7:1]** | **408** | 13.48 |

- Competitive scores on Multi-Query Associative Recall (memory) and Long-Range Arena (long-range) tasks
- State-of-The Art Language Modeling Perplexity on SlimPajama (15B) at 350M parameter scale
- Downstream Performance on LMEval and PALOMA tasks matches PPL performance gap

## SCALING



Legend: Llama, Mamba, RWKV-4, xLSTM[7:1], xLSTM[1:0]
Axes: Number of Parameters ($\times 10^9$) vs Validation Perplexity
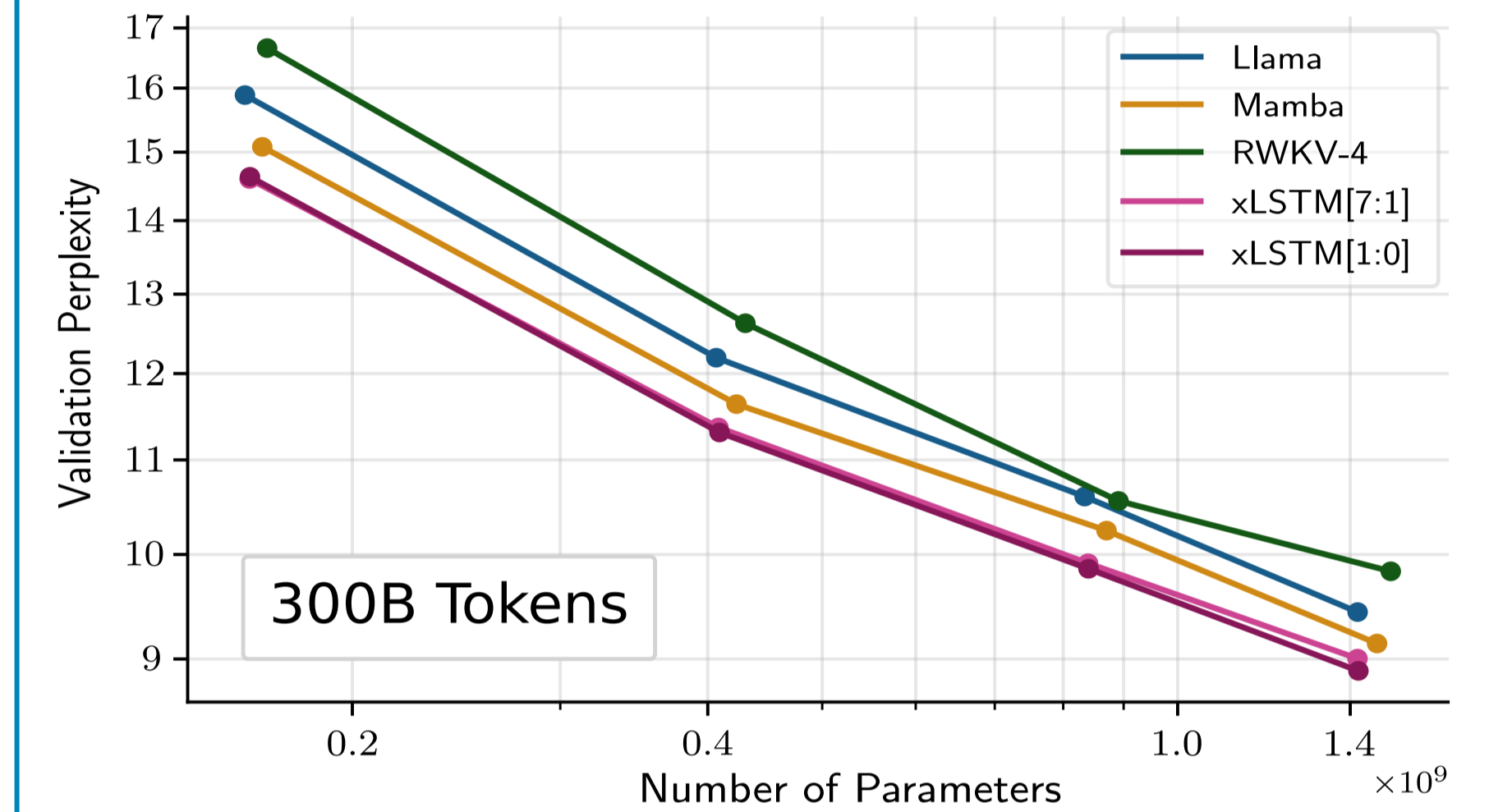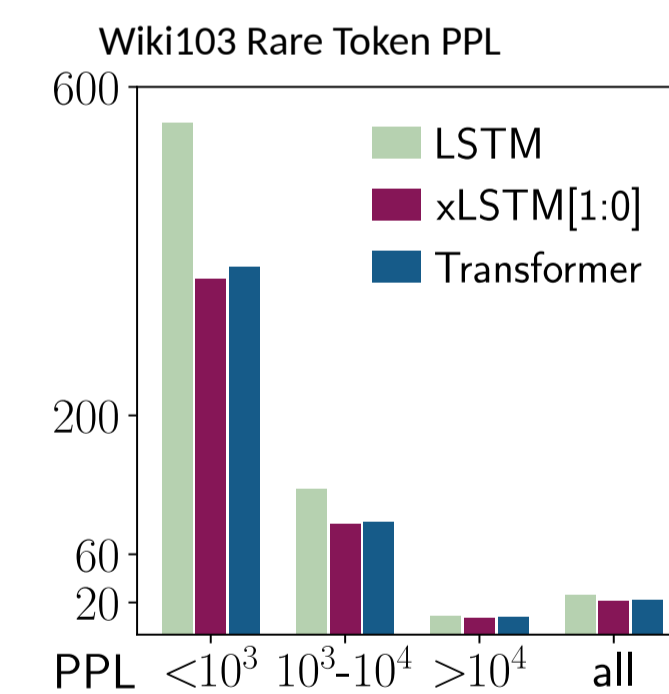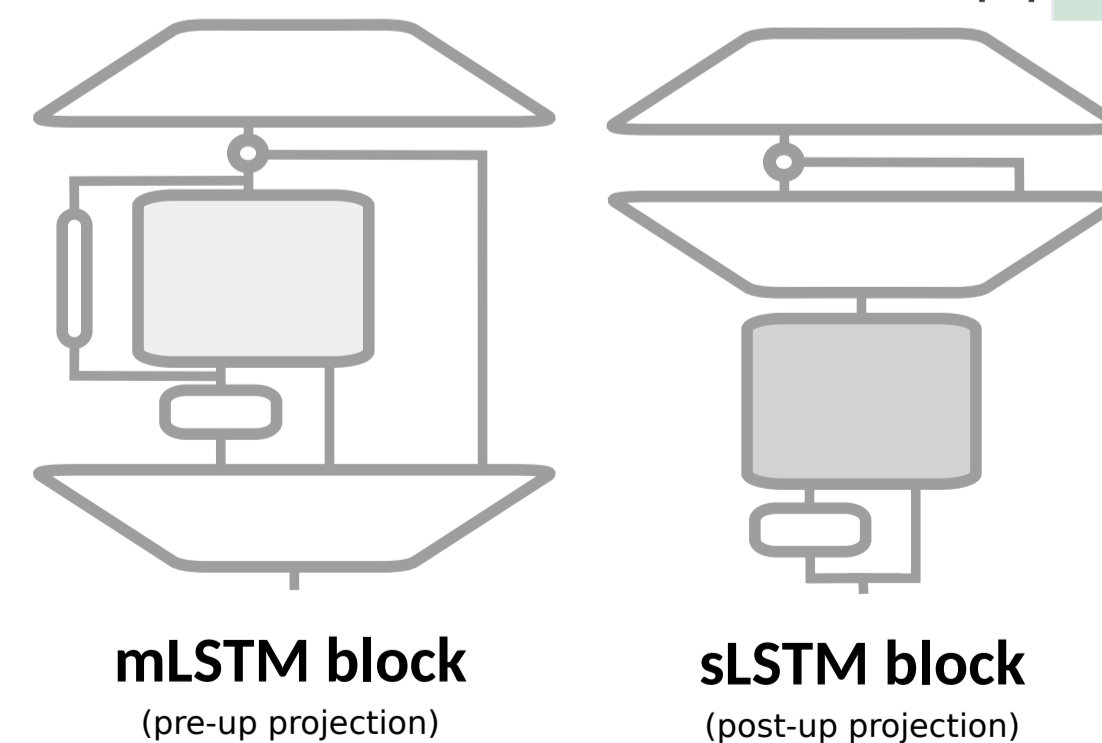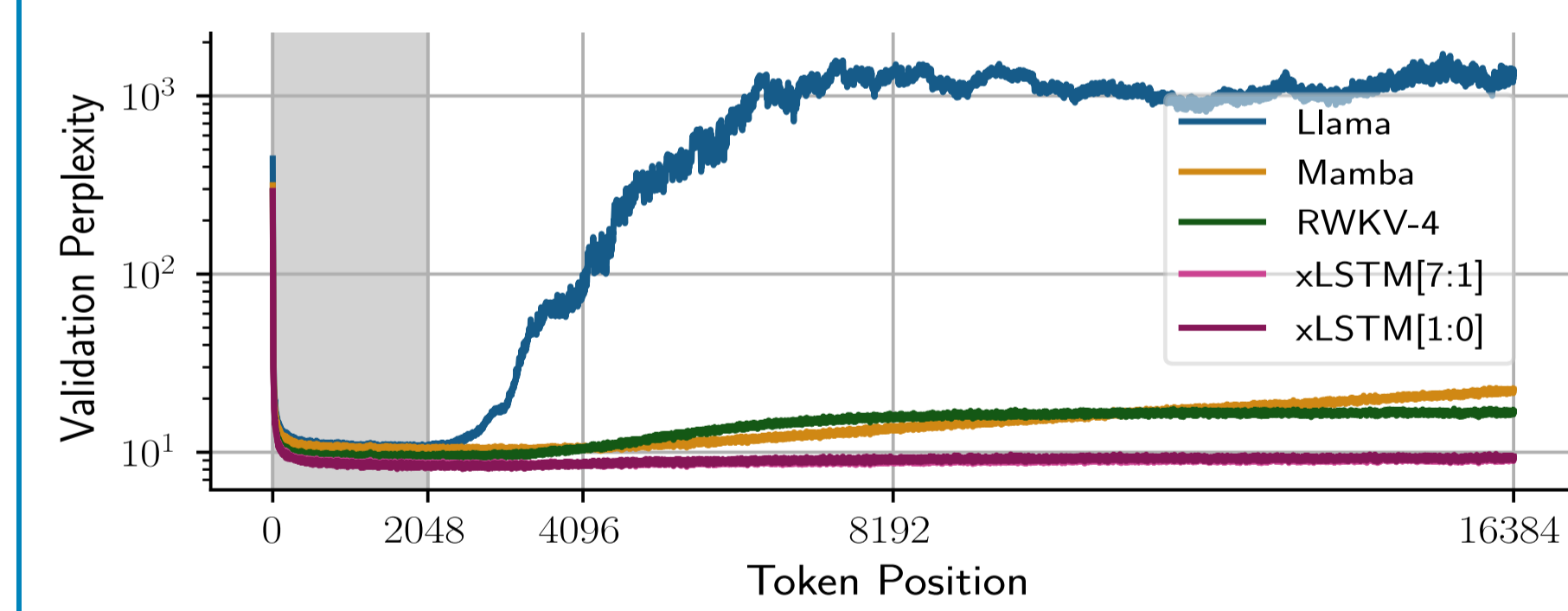300B Tokens

- Perplexity after training on 300B tokens of SlimPajama
- xLSTM scales similar to competitors with more parameters

## LENGTH EXTRAPOLATION



Legend: Llama, Mamba, RWKV-4, xLSTM[7:1], xLSTM[1:0]
Axes: Token Position vs Validation Perplexity

## GENERATION SPEED



Legend: Llama FP16 (prefill 256), Llama FP16 (prefill 2048), xLSTM[1:0] FP16 (t.c), Mamba FP16
Axes: Batch Size vs Tokens / s
Out of Memory