

Unlocking State-Tracking in Linear RNNs through Negative Eigenvalues

Riccardo Grazzi*, Julien Siems*, Arbër Zela,
Jörg KH Franke, Frank Hutter, Massimiliano Pontil
(*equal contribution)



State Tracking

Initial State



State Transitions

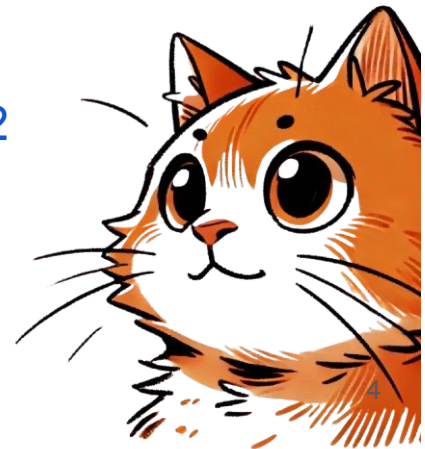
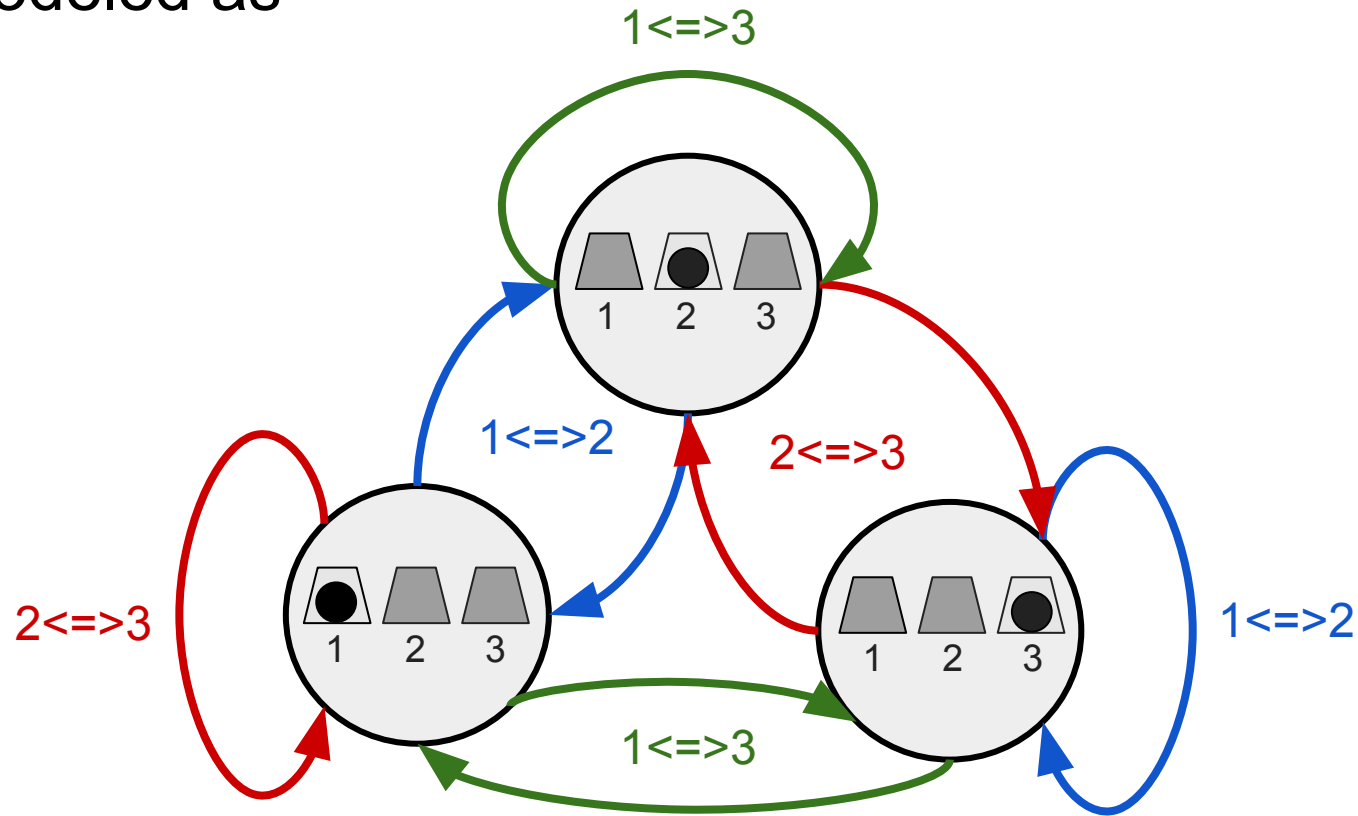


Where is the ball?



- **State is not observable:** the ball position is shown only at the start
- The cat needs to **watch the entire sequence of transitions**

Modeled as



Finite State Automata (FSA)

States (Finite set) $\longrightarrow Q = \left\{ \begin{array}{c} \text{trapezoid}_1 \text{ trapezoid}_2 \text{ trapezoid}_3 \\ 1 \quad 2 \quad 3 \end{array}, \begin{array}{c} \text{trapezoid}_1 \text{ trapezoid}_2 \text{ trapezoid}_3 \\ 1 \quad 2 \quad 3 \end{array}, \begin{array}{c} \text{trapezoid}_1 \text{ trapezoid}_2 \text{ trapezoid}_3 \\ 1 \quad 2 \quad 3 \end{array} \right\}$

Alphabet (Finite set) $\longrightarrow \Sigma = \{ 1 \leq 2, 1 \leq 3, 2 \leq 3 \}$

Initial state

$$q_0 \in Q$$

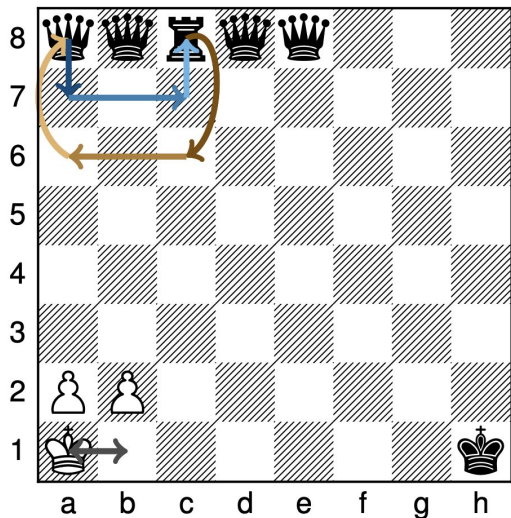
$$\delta : Q \times \Sigma \longrightarrow Q$$

Transition function

State Tracking = mimic an FSA:

map the **sequences of transitions** (input) to **sequences of states** (output).

State Tracking Tasks in Text Data



Code evaluation

Entity Tracking

Tracking a chessboard with non-standard (source, target) notation for moves

(A8, A7), (A1, B1), (C8, C6), (B1, A1), (A7, C7), (A1, B1), (C6, A6), (B1, A1), (C7, C8), (A1, B1), (A6, A8)

Input to the model

```
x = [ 0, 0, 1, 0, 0 ]  
x [ 1 ], x [ 3 ] = x [ 3 ], x [ 1 ] # Swap 1, 3
```

Alice, Bob and Carl each have a coin. Carl is the only one having a penny. Alice and Carl trade coins.

Linear RNNs (One Layer)

State matrix input token Output Channel mix (MLP)

$$\mathbf{H}_i = \mathbf{A}(\mathbf{x}_i) \mathbf{H}_{i-1} + \mathbf{B}(\mathbf{x}_i), \quad \hat{\mathbf{y}}_i = \text{dec}(\mathbf{H}_i, \mathbf{x}_i)$$

State-transition matrix

	$\mathbf{A}(\mathbf{x}_t)$	$\mathbf{B}(\mathbf{x}_t)$
Mamba	$\text{Diag}(\exp(-\Delta_t \odot \exp(\mathbf{w}_{1,i})))$	$k_{t,i} \Delta_t \odot \mathbf{x}_t$
GLA	$\text{Diag}(\alpha_t)$	$\mathbf{k}_t \mathbf{v}_t^\top$
DeltaNet	$\mathbf{I} - \beta_t \mathbf{k}_t \mathbf{k}_t^\top$ ← GH, non-diagonal	$\beta_t \mathbf{k}_t \mathbf{v}_t^\top$

Gu, Albert, and Tri Dao. "Mamba: Linear-time sequence modeling with selective state spaces." *arXiv* (2023).

Yang, Songlin, et al. "Gated Linear Attention Transformers with Hardware-Efficient Training." *ICML 2024*

Yang, Songlin, et al. "Parallelizing Linear Transformers with the Delta Rule over Sequence Length.", *NeurIPS 2024*

Linearity + heavily structured matrices make the recurrence efficiently parallelizable

Linear RNNs (One Layer)

State matrix input token Output Channel mix (MLP)

$$\mathbf{H}_i = \mathbf{A}(\mathbf{x}_i) \mathbf{H}_{i-1} + \mathbf{B}(\mathbf{x}_i), \quad \hat{\mathbf{y}}_i = \text{dec}(\mathbf{H}_i, \mathbf{x}_i)$$

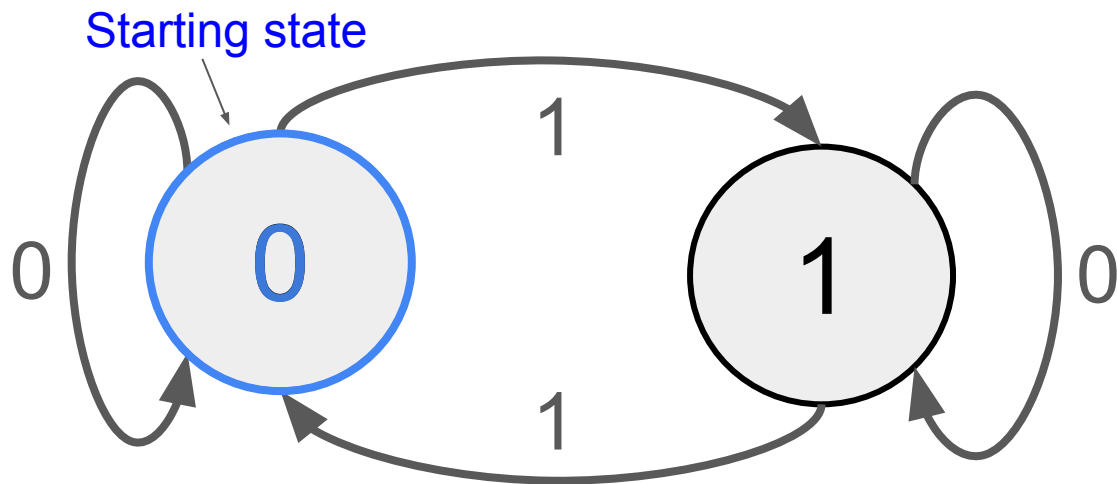
State-transition matrix

	$\mathbf{A}(\mathbf{x}_t)$	$\mathbf{B}(\mathbf{x}_t)$
Mamba	$\text{Diag}(\exp(-\Delta_t \odot \exp(\mathbf{w}_{1,i})))$	$k_{t,i} \Delta_t \odot \mathbf{x}_t$
GLA	$\text{Diag}(\boldsymbol{\alpha}_t)$	$\mathbf{k}_t \mathbf{v}_t^\top$
DeltaNet	$\mathbf{I} - \beta_t \mathbf{k}_t \mathbf{k}_t^\top$	$\beta_t \mathbf{k}_t \mathbf{v}_t^\top$

Transformers are Linear RNNs with infinite dimensional state and $\mathbf{A}(\mathbf{x}_t) = \mathbf{I}$

Katharopoulos, Angelos, et al. "Transformers are rnns: Fast autoregressive transformers with linear attention." ICML 2020.

Parity (2-cups game, addition modulo 2)

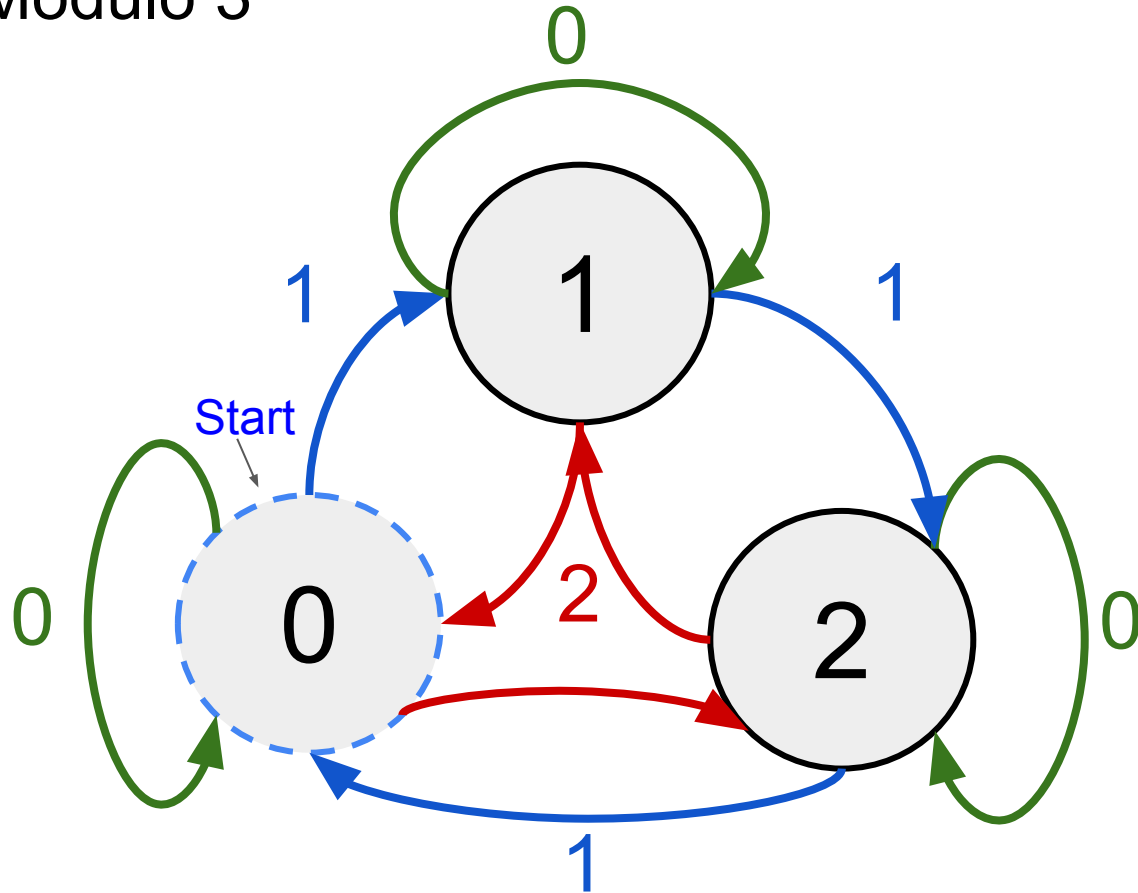


Input bits (transitions)

1	1	0	0	1	0	...
1	0	0	0	1	1	...

Parity (states)

Addition Modulo 3



Solving Parity with a Scalar Linear RNN

$$h_i = a(x_i)h_{i-1} + x_i$$

Solution 1: State = sum of previous values

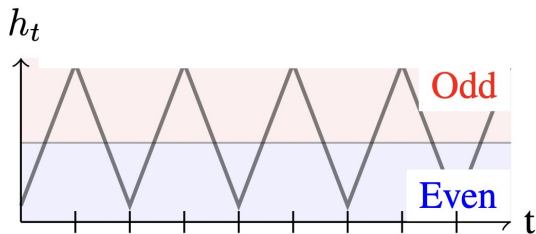
$$a(x_i) = 1$$

$$h_t = \sum_{i=1}^t x_i$$

$$y_t = h_t \bmod 2 \quad (\text{state blows up!})$$

Solution 2: State = parity

$$a(1) = -1, \quad a(0) = 1 \quad y_t = h_t \quad (\text{negative values})$$



Issue with Linear RNNs

	State-transition matrix $\rightarrow \mathbf{A}(\mathbf{x}_t)$	$\mathbf{B}(\mathbf{x}_t)$
Mamba	$\text{Diag}(\exp(-\Delta_t \odot \exp(\mathbf{w}_{1,i})))$	$k_{t,i} \Delta_t \odot \mathbf{x}_t$
GLA	$\text{Diag}(\alpha_t)$	$\mathbf{k}_t \mathbf{v}_t^\top$
DeltaNet	$\mathbf{I} - \beta_t \mathbf{k}_t \mathbf{k}_t^\top$	$\beta_t \mathbf{k}_t \mathbf{v}_t^\top$

$$\Delta_{t,i} \geq 0, \quad \alpha_{t,i} \geq 0, \quad \beta_t \in (0, 1), \quad \mathbf{k}_t \in \mathbb{R}^n, \quad \|\mathbf{k}_t\| = 1$$

All state-transition matrices have **positive eigenvalues** in $[0,1]$.

diagonal Linear RNN with positive values *cannot* solve parity in finite precision (Sarrof et al. 2024)

LLMs Struggle to Track States

Transformers and diagonal linear RNNs cannot track states in limited precision and for arbitrary input lengths (Hahn 2020, Merrill et al. 2023, 2024, Sarrof et al. 2024).

In contrast, RNNs and linear RNNs with **full state transition matrices** can track states with only one layer, but cannot be parallelized efficiently.

What about **scalable non-diagonal Linear RNNs** like DeltaNet?

Hahn, Michael. "Theoretical limitations of self-attention in neural sequence models." *Transactions of the Association for Computational Linguistics* 8 (2020): 156-171.

William Merrill and Ashish Sabharwal. The parallelism tradeoff: Limitations of log-precision transformers. *Transactions of the Association for Computational Linguistics*, 11:531–545, 2023.

William Merrill, Jackson Petty, and Ashish Sabharwal. The Illusion of State in State-Space Models. *ICML* 2024.

Yash Sarrof, Yana Veitsman, and Michael Hahn. The Expressive Capacity of State Space Models: A Formal Language Perspective. *NeurIPS* 2024.

Contribution: Limits of Linear RNNs in Finite Precision

Thm. 1 (Parity): Finite precision linear RNNs cannot solve parity at arbitrary input lengths if for all layers

$$\lambda \in \mathbb{R}, \lambda \geq 0 \quad \forall \lambda \in \text{eigs}(\mathbf{A}(\mathbf{x})) \quad \forall \mathbf{x}$$

Thm. 2 (Modular Counting): Finite precision linear RNNs with L layers cannot count modulo m , with m not a power of two, if for every $i \in \{1, \dots, L\}$ the i -th layer satisfies

$$\lambda \in \mathbb{R} \quad \forall \lambda \in \text{eigs}(\mathbf{A}(\mathbf{x}_1) \cdots \mathbf{A}(\mathbf{x}_{2^{i-1}})) \quad \forall \mathbf{x}_1, \dots, \mathbf{x}_{2^{i-1}}$$

⇒ Most linear RNNs cannot solve parity (only positive eigenvalues)

⇒ Diagonal real-valued linear RNNs cannot do modular counting

Trading off Expressivity and Computational Complexity

$$\mathbf{A}(\mathbf{x}_t) = \begin{bmatrix} \text{light blue} & & & \\ & \text{light blue} & & \\ & & \text{dark blue} & \\ & & & \text{dark blue} \end{bmatrix}$$

Diagonal: → Mamba, mLSTM,
GLA

Very fast computation, but
can't go beyond parity

Trading off Expressivity and Computational Complexity

$$\mathbf{A}(\mathbf{x}_t) = \begin{bmatrix} \text{blue} & & & \\ & \text{blue} & & \\ & & \text{blue} & \\ & & & \text{blue} \end{bmatrix} + \begin{bmatrix} \text{red} \\ \text{red} \\ \text{red} \\ \text{red} \end{bmatrix} \times \begin{bmatrix} \text{pink} & \text{pink} & \text{pink} & \text{pink} \end{bmatrix}$$

Rank 1 Update → DeltaNet

$$\begin{bmatrix} \text{blue} & & & \\ & \text{blue} & & \\ & & \text{blue} & \\ & & & \text{blue} \end{bmatrix} + \begin{bmatrix} \text{red} & \text{purple} \\ \text{red} & \text{purple} \\ \text{red} & \text{purple} \\ \text{red} & \text{purple} \end{bmatrix} \times \begin{bmatrix} \text{pink} & \text{pink} & \text{pink} & \text{pink} \\ \text{orange} & \text{orange} & \text{orange} & \text{orange} \end{bmatrix}$$

Rank 2 Update → DeltaProduct

$$\mathbf{I} - \beta_t \mathbf{k}_t \mathbf{k}_t^\top$$

$$(\mathbf{I} - \beta_{1,t} \mathbf{k}_{1,t} \mathbf{k}_{1,t}^\top) (\mathbf{I} - \beta_{2,t} \mathbf{k}_{2,t} \mathbf{k}_{2,t}^\top)$$

$$\boxed{\begin{aligned} (1 - \beta_{i,t}) \in [-1, 1] &\implies \|\mathbf{A}(\mathbf{x}_t)\| \leq 1 \\ \|\mathbf{k}_{1,t}\| &= 1 \end{aligned}} \quad \text{Stable recurrence!}$$

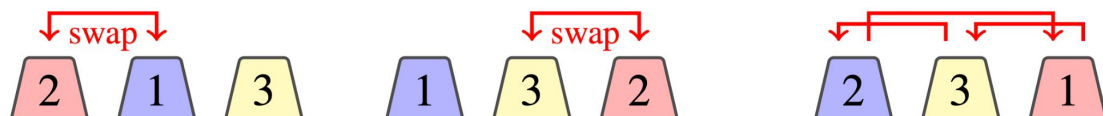
Products of Generalized Householder (GH) Matrices

$$\mathcal{M}_k^n(\Omega) := \{ \mathbf{C}_1 \mathbf{C}_2 \cdots \mathbf{C}_k : \mathbf{C}_i = \mathbf{I} - \beta_i \mathbf{v}_i \mathbf{v}_i^\top, \quad (1 - \beta_i) \in \Omega, \quad \mathbf{v}_i \in \mathbb{R}^n, \|\mathbf{v}_i\| = 1 \}$$

↑
Eigenvalue range

For DeltaNet, $k = 1, \Omega = [0, 1]$

Orthogonal matrices ($\neq \mathbf{I}$) are included only if $-1 \in \Omega$ ($\beta_i = 2$)



$$\begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \times \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}$$

$$\mathbf{I} - 2\mathbf{v}_1 \mathbf{v}_1^\top \quad \mathbf{I} - 2\mathbf{v}_2 \mathbf{v}_2^\top$$


Contribution: Expressivity of Products of GH Matrices

Thm. 3 (Permutations): Finite precision linear RNNs with one layer where state-transition matrices are in $\mathcal{M}_{k-1}^n([-1, 1])$ can model any FSA whose transitions $\delta(\cdot, w) : Q \rightarrow Q$ correspond to permutations of at most k elements, when n is large enough.

Thm. 4 (General FSA): Finite precision linear RNNs with multiple layers where state-transition matrices are in $\mathcal{M}_n^n([-1, 1])$ for a large enough n , can model any finite state automaton.

\Rightarrow We can easily modify DeltaNet to have state transition matrices in $\mathcal{M}_1^n([-1, 1])$ and thus model **swap permutations**

Eigenvalue Extension for Mamba and DeltaNet

	$[0, 1]$	$[-1, 1]$
$\mathbf{A}(\mathbf{x}_t)$ 	Mamba DeltaNet	$\text{Diag}(\mathbf{s}(\mathbf{x}_t))$ $\mathbf{I} - \beta_t \mathbf{k}_t \mathbf{k}_t^\top$

Change for DeltaNet is a *one-liner*!

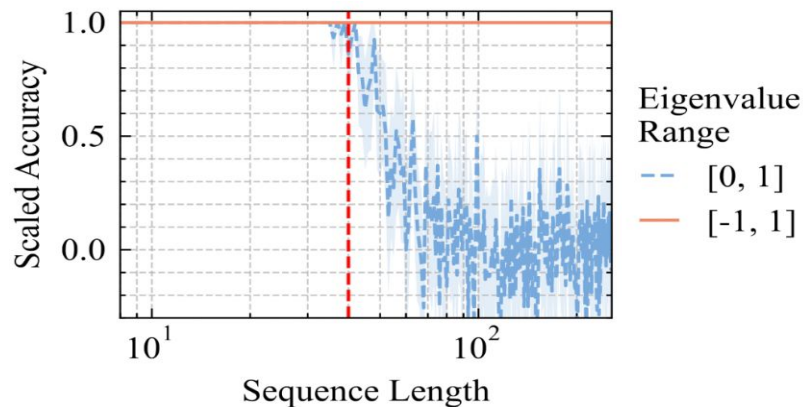
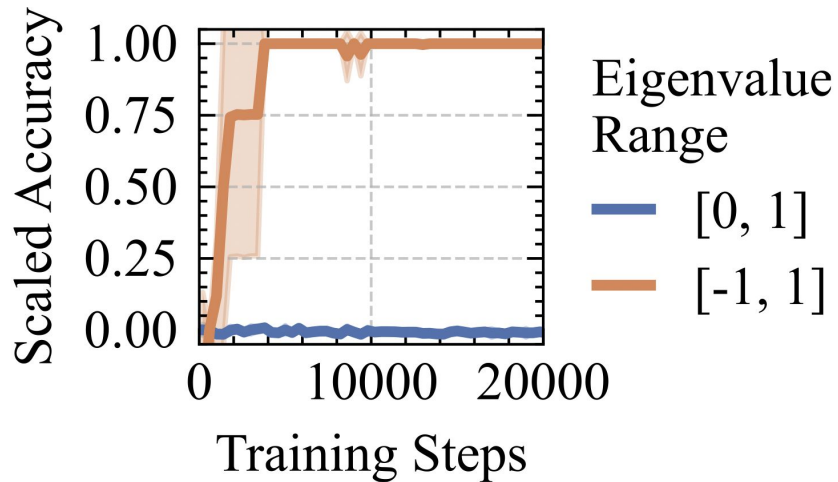
```
if self.use_beta:
- beta = rearrange(self.b_proj(hidden_states), 'b l h -> b h l').sigmoid()
+ beta = 2 * rearrange(self.b_proj(hidden_states), 'b l h -> b h l').sigmoid()
else:
    beta = q.new_ones(q.shape[0], q.shape[1], q.shape[2])
```

Code from [Flash Linear Attention](#) (Yang et al. 2024)

Experiments - Parity

	Parity
Transformer	0.022
mLSTM	0.087 (0.04)
sLSTM	1.000 (1.00)
Mamba [0, 1]	0.000
Mamba [-1, 1]	1.000
DeltaNet [0, 1]	0.017
DeltaNet [-1, 1]	1.000

Accuracy

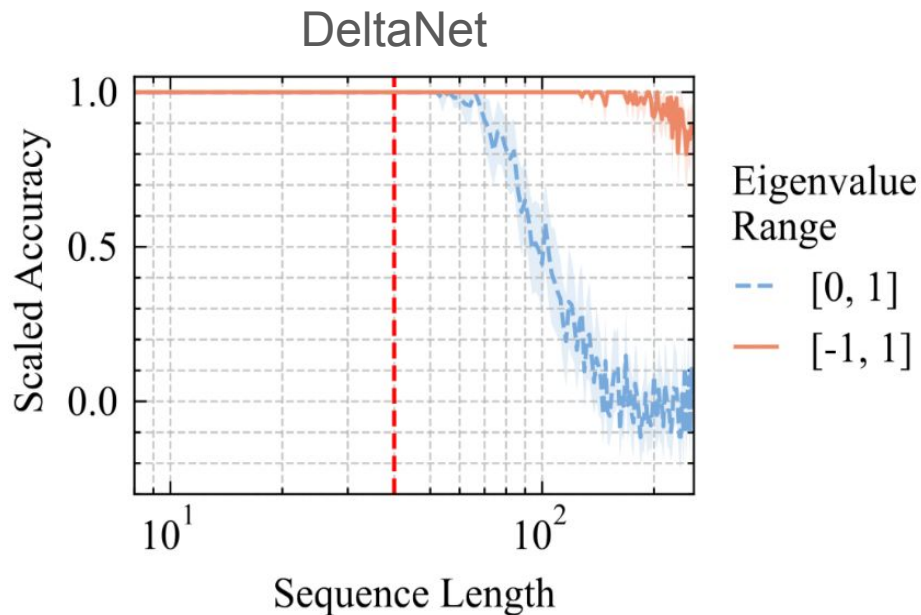


Experiments - Modular Arithmetic

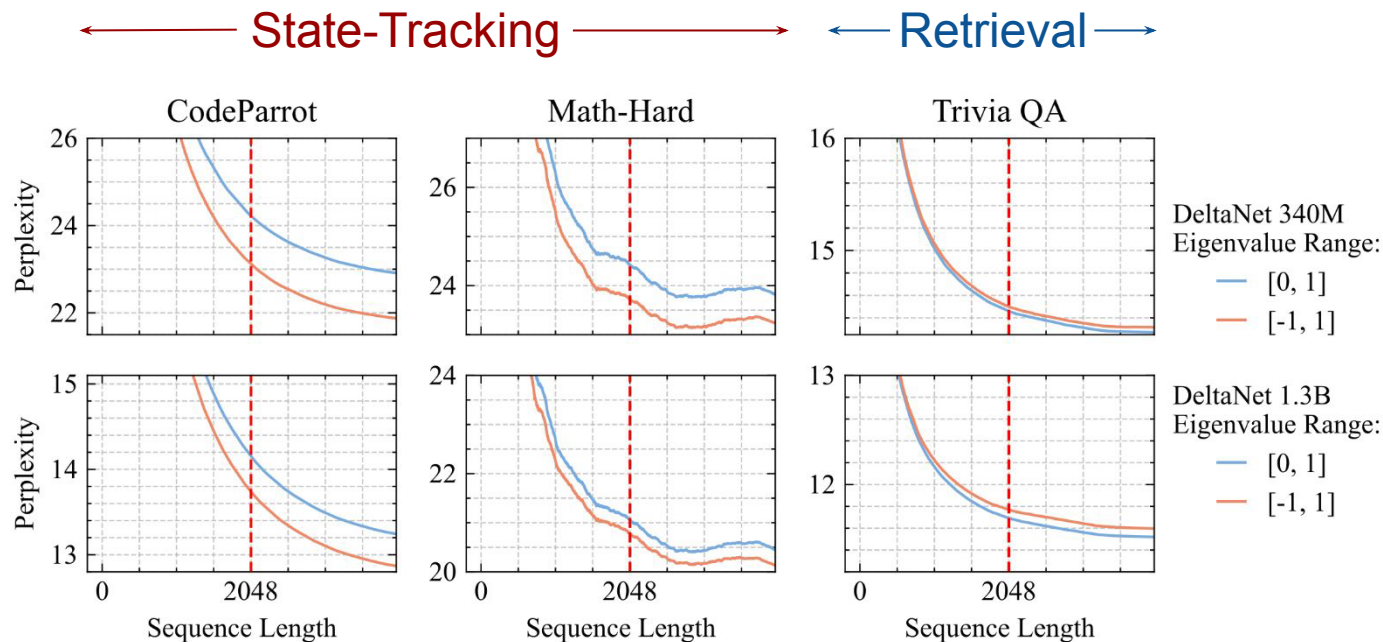
Mod. Arithm. (w/o brackets):

$$2 - 3 - 3 * 2 \bmod 5 = 3$$

	Mod. Arithm. (w/o brackets)
Transformer	0.031
mLSTM	0.040 (0.04)
sLSTM	0.787 (1.00)
Mamba [0, 1]	0.095
Mamba [-1, 1]	0.241
DeltaNet [0, 1]	0.314
DeltaNet [-1, 1]	0.971



Experiments - Language Modelling



→ *Note:* Extended eigenvalue range doesn't cause training instability

Conclusion

- Inclusion of negative eigenvalues expands the expressivity of linear RNNs allowing them to solve state-tracking problems
- Efficient non-diagonal linear RNNs such as DeltaNet and RWKVv7 are promising due to their superior expressivity compared to Mamba.

Future Directions:

- What is the limit of the expressivity of DeltaNet $[-1,1]$?
- Is standard pretraining exploiting the increased expressivity? Are there better ways?
- What is the trade-off between in-context associative recall and state-tracking?