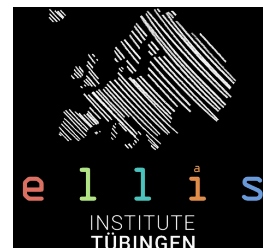


Unlocking State-Tracking in Linear RNNs through Negative Eigenvalues

AutoML Seminar
9th of January 2025

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



Outline

Part 1 (Riccardo)

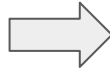
1. State-tracking
2. Linear RNNs and Known Limits of Current Architectures
3. Contribution 1: Limits of Linear RNNs.
4. Contribution 2: How to unlock State-tracking.

Part 2 (Julien)

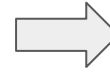
1. Extending the Eigenvalue Range of Mamba and DeltaNet
2. Synthetic Experiments
3. Language Modeling Experiments
4. Conclusion and Future Works

State Tracking

Show Initial State



State Transitions

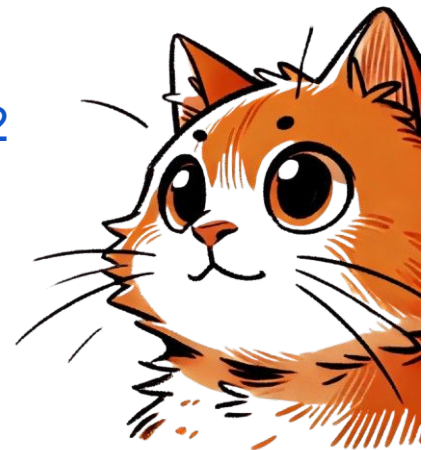
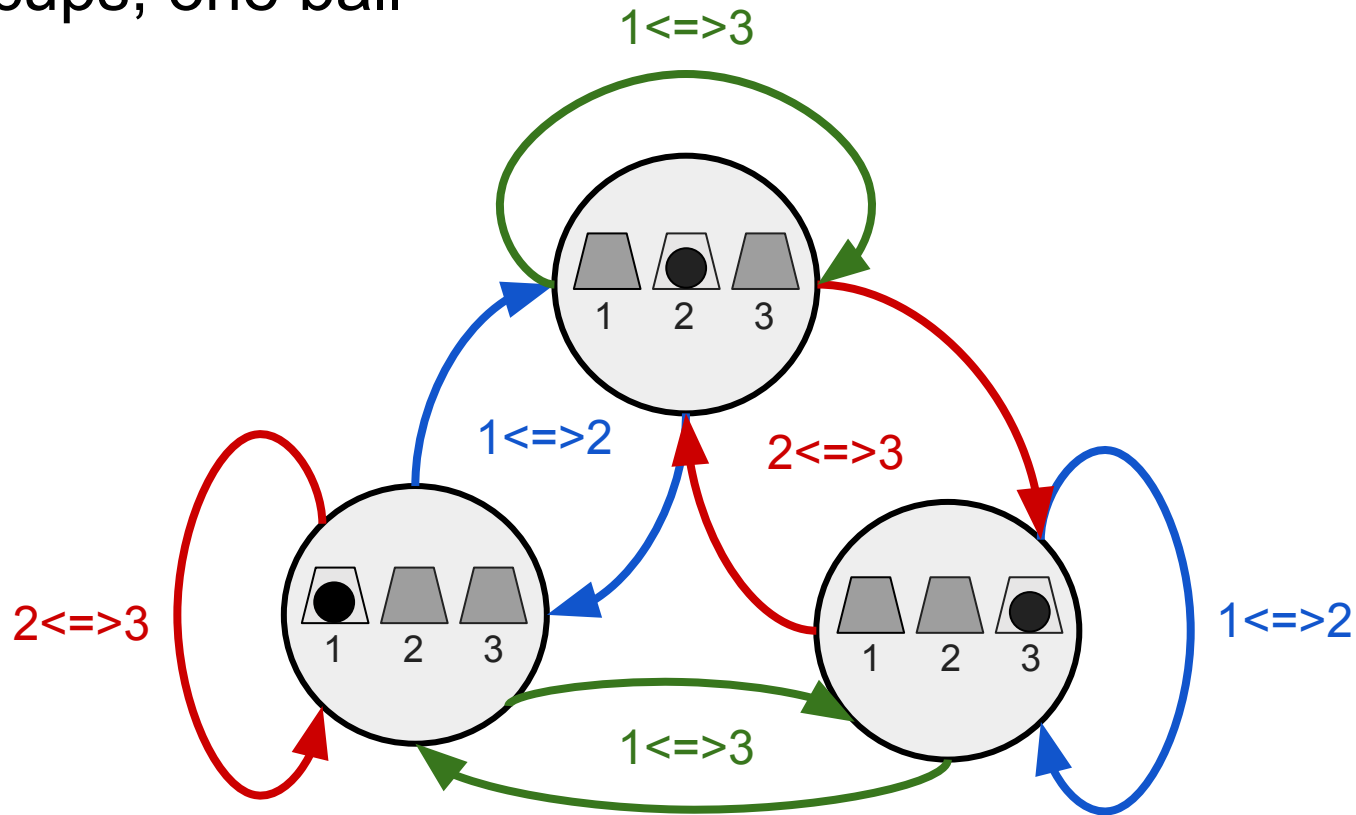


Where is the ball?



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

3 cups, one ball



Finite State Automata (FSA)

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

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

Initial state

$$q_0 \in Q$$

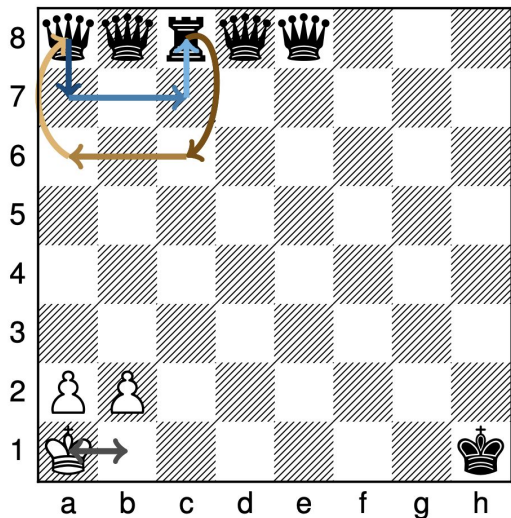
$$\delta : Q \times \Sigma \longrightarrow Q$$

Transition function

State Tracking = mimic an FSA:

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

State Tracking Tasks in Text Data



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

Code evaluation

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

Entity Tracking

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

State-tracking?

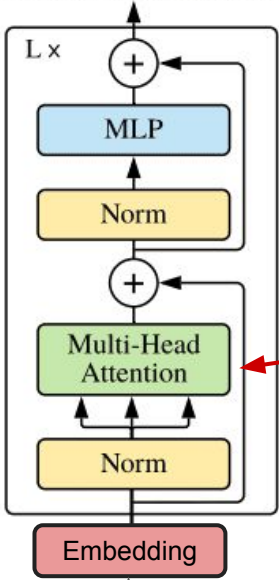


Alice, Bob and Carl each have a coin. Carl is the only one having a penny. Alice and Carl trade coins.

*Alice, Bob and Carl each have a coin. Carl is the only one having a penny. Carl trades **his penny** with Alice.*

Modern Language Modeling Architectures

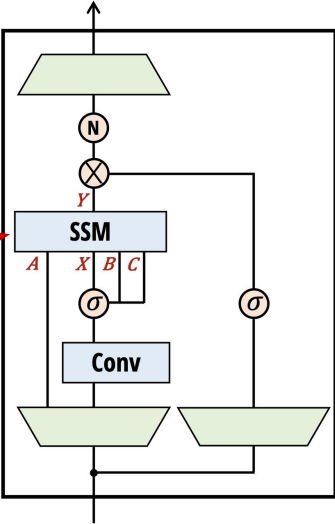
Transformer



Channel Mixing
(Non-Linear)

Token Mixing
(Parallelizable)

Mamba 2



Parallel Mamba Block

Linear RNNs:
More efficient for long sequences

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

GH, non-diagonal: token+channel mix

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 recursion efficiently parallelizable

Linear RNNs (One Layer)

State matrix

input token

Output

Channel mix (MLP)

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

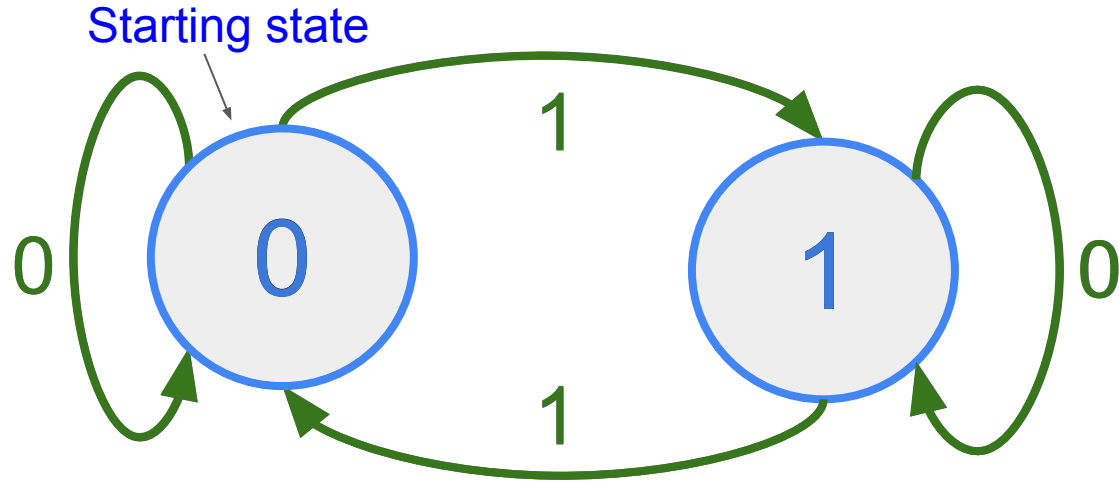
State-transition matrix

	$\mathbf{A}(\mathbf{x}_t)$	$\mathbf{B}(\mathbf{x}_t)$
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)

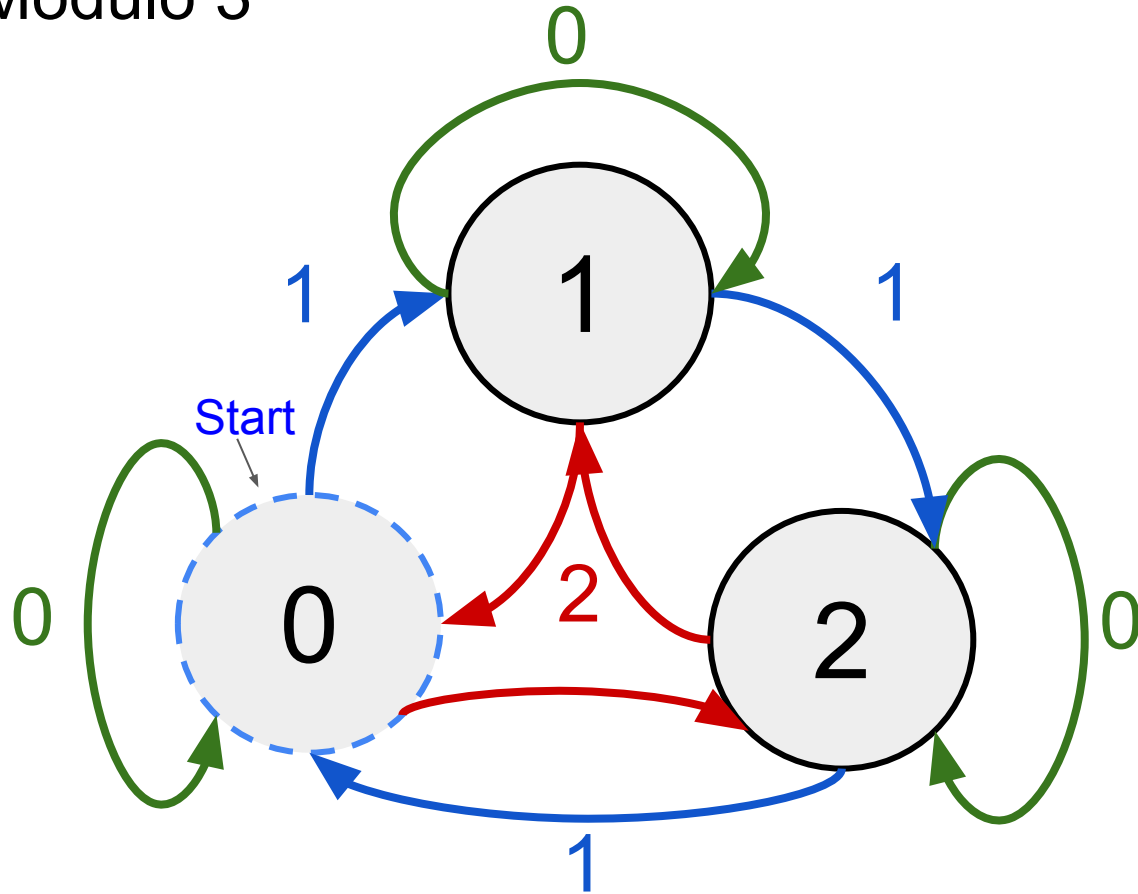
1	1	0	0	1	0	...
---	---	---	---	---	---	-----

Parity (states)

						...
--	--	--	--	--	--	-----

We'd like to handle arbitrary long sequences

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

$$a(x_i) = 1$$

$$h_t = \sum_{i=1}^t x_i \quad y_t = h_t \bmod 2 \quad (\text{state blows up!})$$

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	$\mathbf{I} - \beta_t \mathbf{k}_t \mathbf{k}_t^\top$ ← GH, non-diagonal	$\beta_t \mathbf{k}_t \mathbf{v}_t^\top$

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

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 provably 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 **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

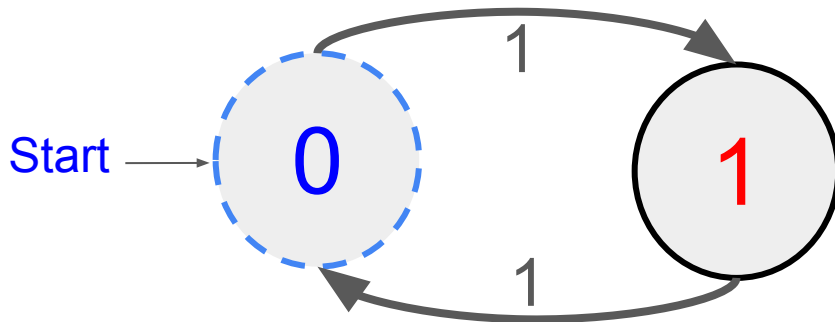
$$\lambda \in \mathbb{R} \quad \forall \lambda \in \text{eigs}(\mathbf{A}(\mathbf{x}_1) \cdots \mathbf{A}(\mathbf{x}_{2^{i-1}})) \quad \forall \mathbf{x}_1, \dots, \mathbf{x}_{2^{i-1}}$$

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

⇒ Diagonal real-valued linear RNNs cannot do modular counting

Theorem 1 - Proof Idea (Same as Sarrof et al. 2024)

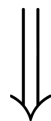
When the input is constant, i.e. $\mathbf{x} = 1, 1, 1, 1, \dots$, then the Linear RNN output in finite precision becomes constant while the state of the parity automaton alternates between 0 and 1.



Input bits (transitions)	1	1	1	1	1	1	1	1	1	1	1
Parity (states)	1	0	1	0	1	0	1	0	1	0	1

Proof Sketch (One Layer): Unrolling the Recurrence

$$\mathbf{x} = (1, 1, 1, 1, \dots, 1)$$



Linear Time Invariant
Recurrence

State \longrightarrow $\mathbf{H}_t = \mathbf{A}(1)\mathbf{H}_{t-1} + \mathbf{B}(1)$

Unroll the recurrence with $\mathbf{H}_0 = 0$ (for simplicity) gives

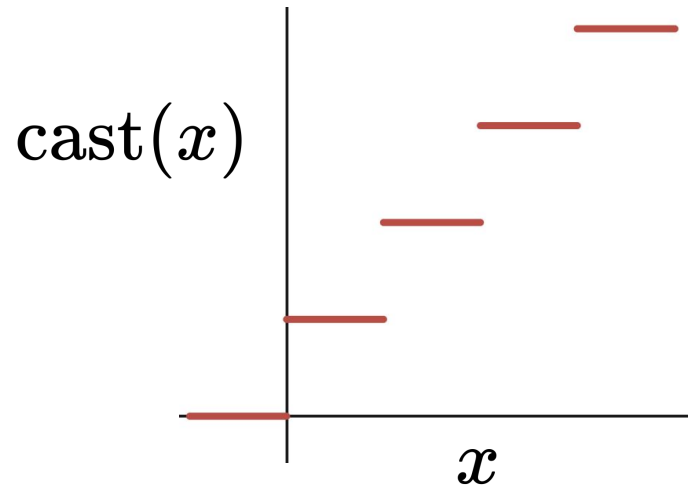
$$\mathbf{H}_t = \sum_{k=0}^{t-1} \mathbf{A}(1)^k \mathbf{B}(1)$$

Finite Precision

$$\text{cast} : \mathbb{R} \rightarrow \mathbb{D}$$

Finite Set of reals

$$\text{cast}(x) \in \arg \min_{z \in \mathbb{D}} |z - x|$$



Extended to matrices by applying it separately to real and imaginary part of each element of the matrix

State in Finite Precision

Infinite precision:
$$\mathbf{H}_t = \sum_{k=0}^{t-1} \mathbf{A}(1)^k \mathbf{B}(1)$$

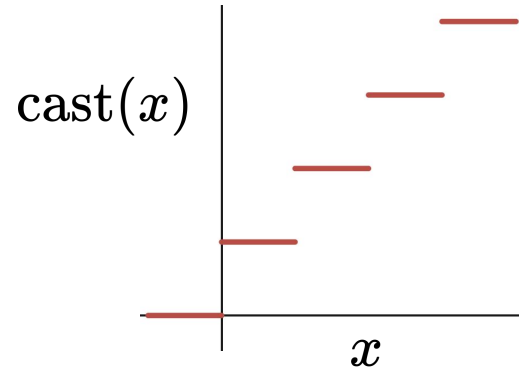
Finite precision:
$$\widehat{\mathbf{H}}_t = \text{cast} \left(\sum_{k=0}^{t-1} \text{cast} (\mathbf{A}(1)^k \mathbf{B}(1)) \right)$$

Simplified model: repeated matrix multiplies and sums are in infinite precision, similar to Merrill et al. (2024)

State in Finite Precision

$$\text{cast}(x) \in \arg \min_{z \in \mathbb{D}} |z - x|$$

$z \in \mathbb{D}$
Finite Set of reals



Cast is applied to matrices by applying it separately to real and imaginary part of each element of matrices and converts real numbers into finite precision:

$$\widehat{\mathbf{H}}_t = \text{cast} \left(\sum_{k=0}^{t-1} \text{cast} \left(\mathbf{A}(1)^k \mathbf{B}(1) \right) \right)$$

Simplified model: repeated matrix multiplies and sums are in infinite precision

Matrix Powers Using the Jordan Canonical Form

$$\mathbf{A}(1)^k = \mathbf{P}\mathbf{J}^k\mathbf{P}^{-1}$$

\mathbf{J}^k Block Diagonal with block i

$$\mathbf{J}_i^k = \begin{bmatrix} \lambda_i^k & \binom{k}{1}\lambda_i^{k-1} & \binom{k}{2}\lambda_i^{k-2} & \cdots & \cdots & \binom{k}{k_i-1}\lambda_i^{k-k_i+1} \\ & \lambda_i^k & \binom{k}{1}\lambda_i^{k-1} & \cdots & \cdots & \binom{k}{k_i-2}\lambda_i^{k-k_i+2} \\ & & \ddots & \ddots & \vdots & \vdots \\ & & & \ddots & \ddots & \vdots \\ & & & & \lambda_i^k & \binom{k}{1}\lambda_i^{k-1} \\ & & & & & \lambda_i^k \end{bmatrix}$$

$k_i \in \mathbb{N}, \lambda_i \in \mathbb{C}$
 ↑
 Eigenvalue

Real and Positive Eigenvalues

The imaginary and real part of each element of $\mathbf{A}(1)^k \mathbf{B}(1)$ take the form

$$a_k = \sum_{i=1}^{n^2} c_i \binom{k}{m_i} \lambda_i^k \quad \lambda_i \in \mathbb{C}, c_i \in \mathbb{R}, m_i \in \mathbb{N}$$

↑
Eigenvalue

$$\binom{k}{m_i} \lambda_i^k = O(k^{m_i} \lambda_i^k)$$

↑
Polynomial

↙
Exponential (Complex base)

Real and Positive Eigenvalues

The imaginary and real part of each element of $\mathbf{A}(1)^k \mathbf{B}(1)$ take the form

$$a_k = \sum_{i=1}^{n^2} c_i \binom{k}{m_i} \lambda_i^k \quad \lambda_i \in \mathbb{C}, c_i \in \mathbb{R}, m_i \in \mathbb{N}$$

↑
Eigenvalue

Lemma 1 (values become constant in finite precision)

$$\text{if } \lambda_i \in \mathbb{R}, \lambda_i > 0 \implies \exists \bar{k} \geq 0, \bar{a} \in \mathbb{R} : \text{cast}(a_k) = \bar{a} \quad \forall k \geq \bar{k}$$



state and layer output
become constant in
finite precision

$$\exists \bar{t} \geq 0 : \begin{aligned} \widehat{\mathbf{H}}_t &= \overline{\mathbf{H}} \\ \overline{\mathbf{y}}_t &= \text{dec}(\widehat{\mathbf{H}}_t, \mathbf{x}_t) = \overline{\mathbf{y}} \end{aligned} \quad \forall t \geq \bar{t}$$

Parity Can't Be Modeled - End of Proof Sketch

If $\mathbf{A}(1)$ has only **real positive** eigenvalues, then

input $\mathbf{x} = \dots 1 1 1 1 1 1 1 1 1 1 1 1 \dots$

LRNN output $\bar{\mathbf{y}} = \dots \bar{y} \bar{y} \bar{y} \bar{y} \bar{y} \bar{y} \bar{y} \bar{y} \bar{y} \bar{y} \bar{y} \bar{y} \dots$

parity $\mathbf{y} = \dots 0 1 0 1 0 1 0 1 0 1 0 1 \dots$

The proof can then proceed by induction over the number of layers

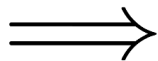
Theorem 2 Proof Sketch - Real Eigenvalues

$$a_k = \sum_{i=1}^{n^2} c_i \binom{k}{m_i} \lambda_i^k \quad \lambda_i \in \mathbb{C}, c_i \in \mathbb{R}, m_i \in \mathbb{N}$$

Lemma 1 (Part 2): if $\lambda_i \in \mathbb{R} \implies \lambda_i^k = \text{sign}(\lambda_i)^k |\lambda_i|^k$



$$\exists \bar{k}, \geq 0, \bar{a}_1, \bar{a}_2 \in \mathbb{R} : \text{cast}(a_{2k}) = \bar{a}_1, \text{cast}(a_{2k+1}) = \bar{a}_2 \quad \forall k \geq \bar{k}$$



$$\exists \bar{t} \geq 0 : \widehat{\mathbf{H}}_{2t} = \overline{\mathbf{H}}_1, \widehat{\mathbf{H}}_{2t+1} = \overline{\mathbf{H}}_2 \quad \forall t \geq \bar{t}$$
$$\overline{\mathbf{y}}_{2t} = \overline{\mathbf{y}}_1, \overline{\mathbf{y}}_{2t+1} = \overline{\mathbf{y}}_2$$

layer output

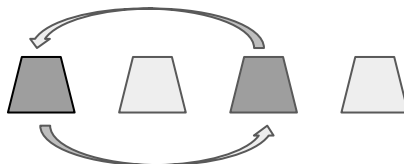
Products of Generalized Householder (GH) Matrices

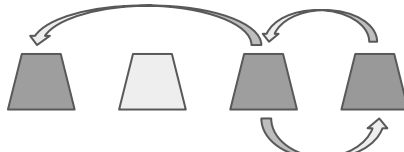
$$\mathcal{M}_k^n(\Omega) := \{ \mathbf{C}_1 \mathbf{C}_2 \cdots \mathbf{C}_k : \mathbf{C}_i = \mathbf{I} - \beta_i \mathbf{v}_i \mathbf{v}_i^\top, \quad (1 - \beta_i) \in \Omega, \quad \mathbf{v}_i \in \mathbb{R}^n, \|\mathbf{v}_i\| = 1 \}$$

Eigenvalue range

DeltaNet has state-transition matrices in $\mathcal{M}_1^n([0, 1])$

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

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


$$(\mathbf{I} - 2\mathbf{v}_1 \mathbf{v}_1^\top)(\mathbf{I} - 2\mathbf{v}_2 \mathbf{v}_2^\top) =$$


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 $\delta(\cdot, w) : Q \rightarrow Q$ are permutations of at most k elements.

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**

Recap of Theoretical Contributions



1. Any Linear RNN with state transition matrices having only **positive real eigenvalues cannot solve parity**.
2. **Diagonal and Triangular Linear RNNs cannot solve modular counting**, even with negative real eigenvalues.




1. Linear RNNs with **products of GH state transition matrices, each with negative eigenvalues**, can mimic any FSA and
2. Can do it with products of $k-1$ GH matrices and **one layer** if the transitions are only permutations of at most k elements.

Open question:

- What can be done with **a single GH matrix + multiple layers?**
Addition modulo m can be done with 2 layers! (See Appendix).

Eigenvalue Extension for Mamba and DeltaNet

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

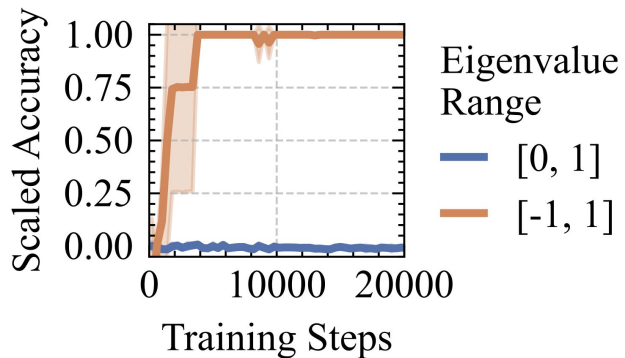
Change for DeltaNet is a *one-liner!*

```
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])
```

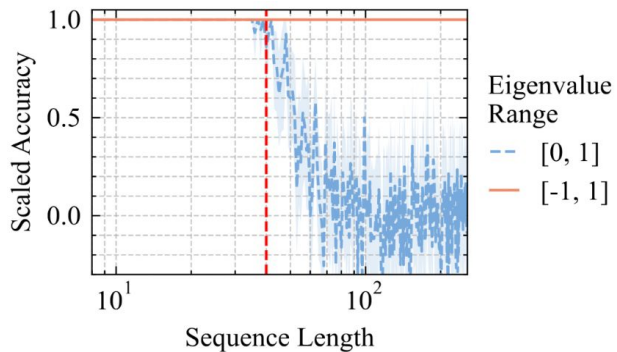
Code from [Flash Linear Attention](#) (Yang et al. 2024)

Experiments - Chomsky Hierarchy

→ Can we actually solve parity using linear RNNs?



	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



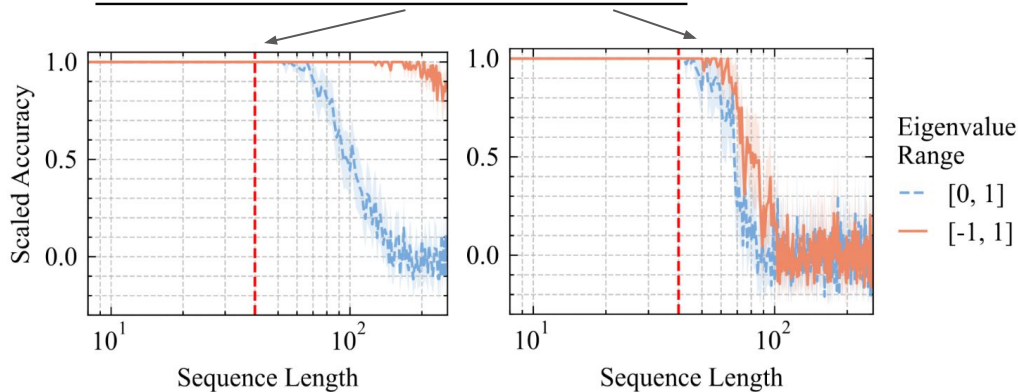
Experiments - Chomsky Hierarchy

Mod. Arithm. (w/o brackets): $2 - 3 - 3 * 2 \bmod 5 = 3$

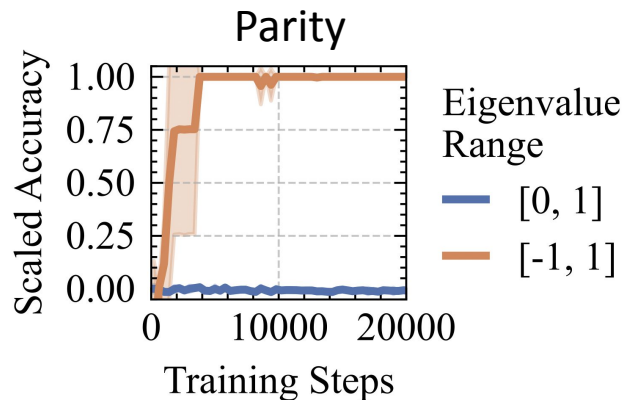
(w/ brackets): $((((3+3)+-1)+-2)-((3-(-3))+((1)+4))) \bmod 5 = 2$

	Mod. Arithm. (w/o brackets)	Mod. Arithm. w/ brackets)
Transformer	0.031	0.025
mLSTM	0.040 (0.04)	0.034 (0.03)
sLSTM	0.787 (1.00)	0.173 (0.57)
Mamba [0, 1]	0.095	0.092
Mamba [-1, 1]	0.241	0.136
DeltaNet [0, 1]	0.314	0.137
DeltaNet [-1, 1]	0.971	0.200

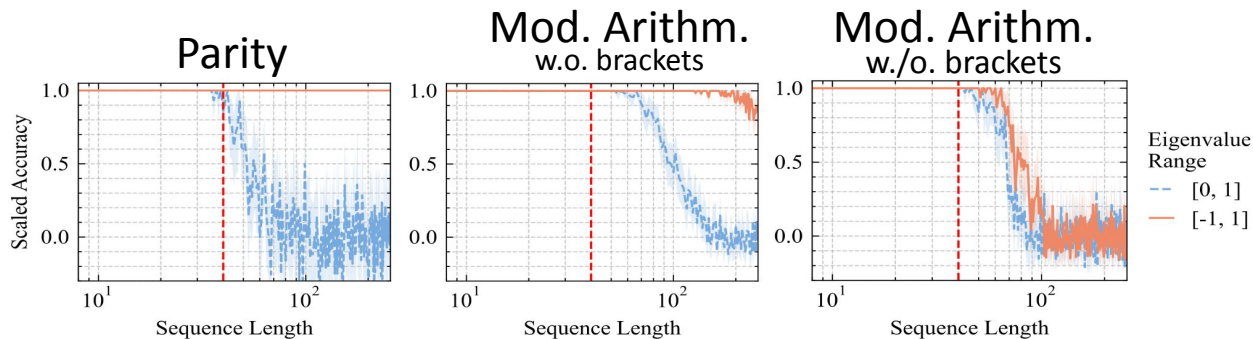
DeltaNet



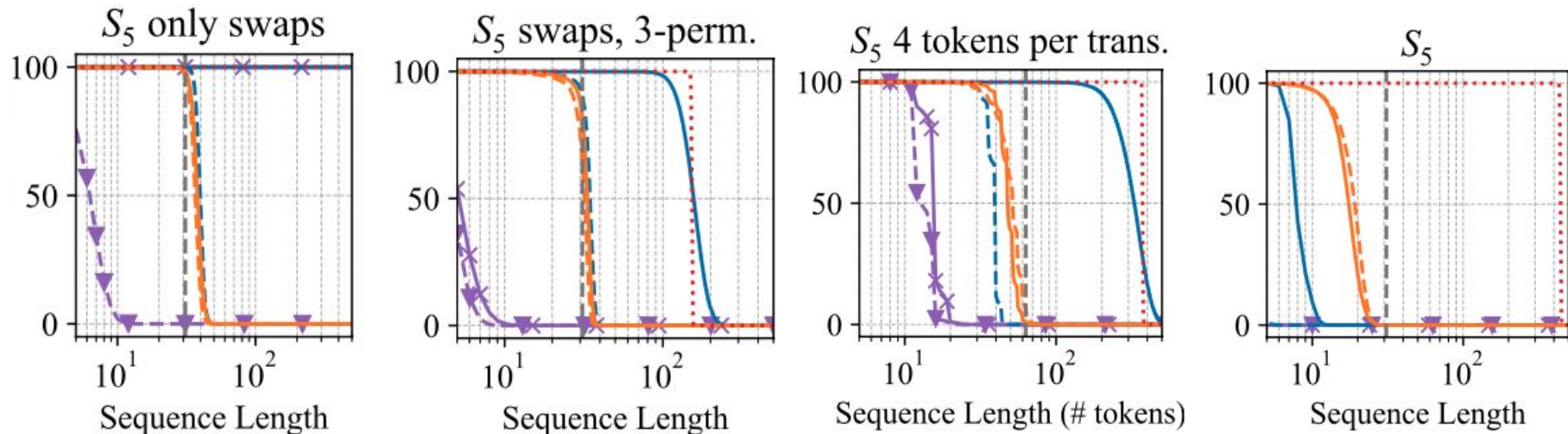
Synthetic Experiments - Part 1



	Parity	Mod. Arithm. (w/o brackets)	Mod. Arithm. w/ brackets
Transformer	0.022	0.031	0.025
mLSTM	0.087 (0.04)	0.040 (0.04)	0.034 (0.03)
sLSTM	1.000 (1.00)	0.787 (1.00)	0.173 (0.57)
Mamba [0, 1]	0.000	0.095	0.092
Mamba [-1, 1]	1.000	0.241	0.136
DeltaNet [0, 1]	0.017	0.314	0.137
DeltaNet [-1, 1]	1.000	0.971	0.200



Synthetic Experiments - Part 2



- ▼-- DeltaNet [0,1] (1L)
- ×— DeltaNet [-1,1] (1L)
- - - DeltaNet [0,1] (5L)
- DeltaNet [-1,1] (5L)
- - - Mamba [0,1] (5L)
- Mamba [-1,1] (5L)
- Full matrix simple

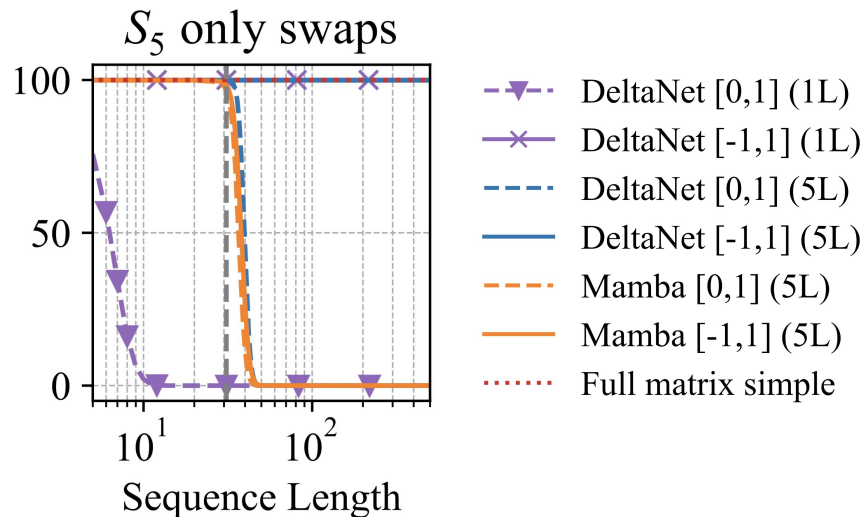
Modified DeltaNet (DeltaNet [-1,1]) can learn with only swap transitions or, with multiple layers, when a transition is encoded with multiple tokens.

Experiments - Permutation Groups

S_5 (Permutation group of 5 elements) (Only swaps)

Example: S_5 only swaps:

$$(1, 2, 3, 4, 5) \circ (1 \rightarrow 2, 2 \rightarrow 1) \\ = (2, 1, 3, 4, 5)$$



Results:

→ DeltaNet [-1, 1] can solve S_5 only swaps (even with 1 layer).

→ Mamba [-1, 1] can't.

Experiments - Permutation Groups

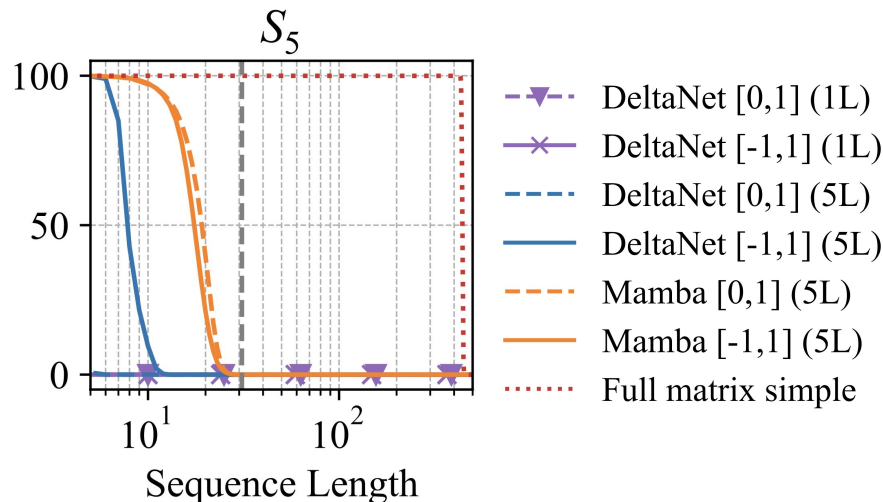
S_5 (Permutation group of 5 elements)

Example: S_5

$(1, 2, 3, 4, 5) \circ$

$(1 \rightarrow 2, 2 \rightarrow 4, 3 \rightarrow 2, 4 \rightarrow 3, 5 \rightarrow 1)$

$= (2, 4, 2, 3, 1)$



Results:

→ DeltaNet and Mamba can't solve S_5 (expected)

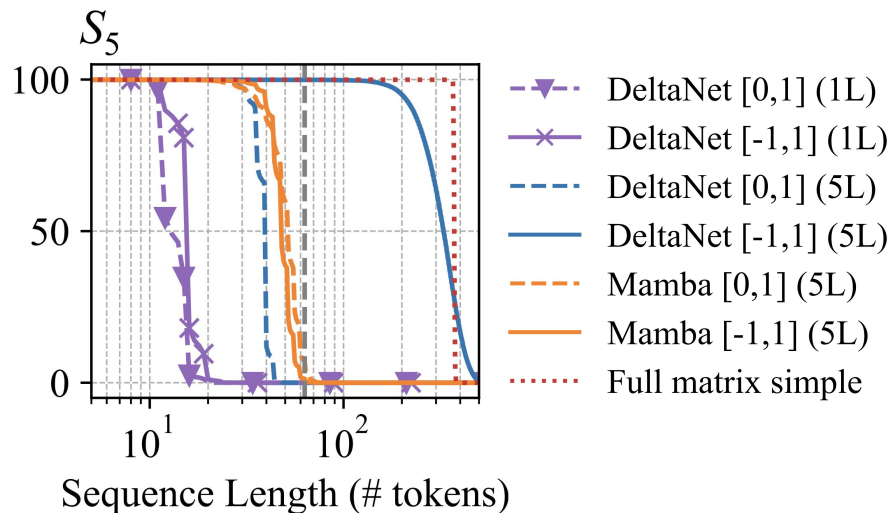
→ Linear RNN with full state-transition matrix can learn to solve S_5

Experiments - Permutation Groups

S_5 (Permutation group of 5 elements)

How can we solve S_5 using DeltaNet?

→ **Products of Householders**



Experiments - Permutation Groups

S_5 (Permutation group of 5 elements)

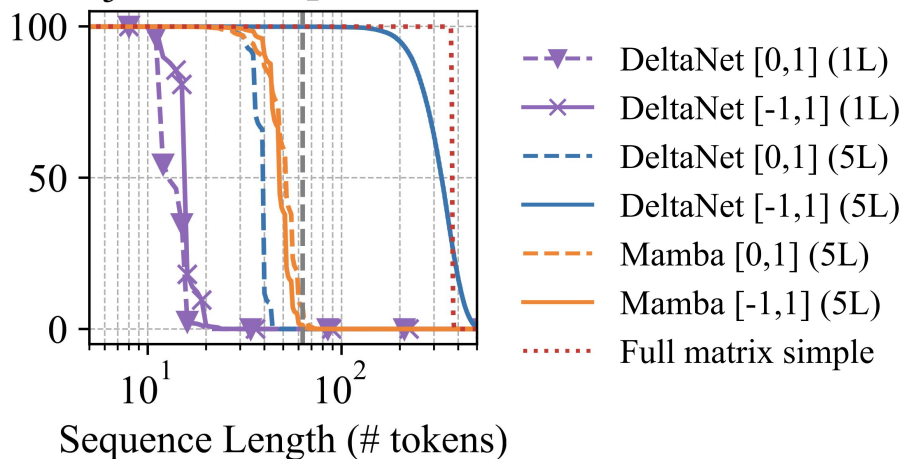
Example: S_5

$(1, 2, 3, 4, 5) \circ$

$(1 \rightarrow 2, 2 \rightarrow 4, 3 \rightarrow 2, 4 \rightarrow 3, 5 \rightarrow 1)$

$= (2, 4, 2, 3, 1)$

S_5 4 tokens per trans.

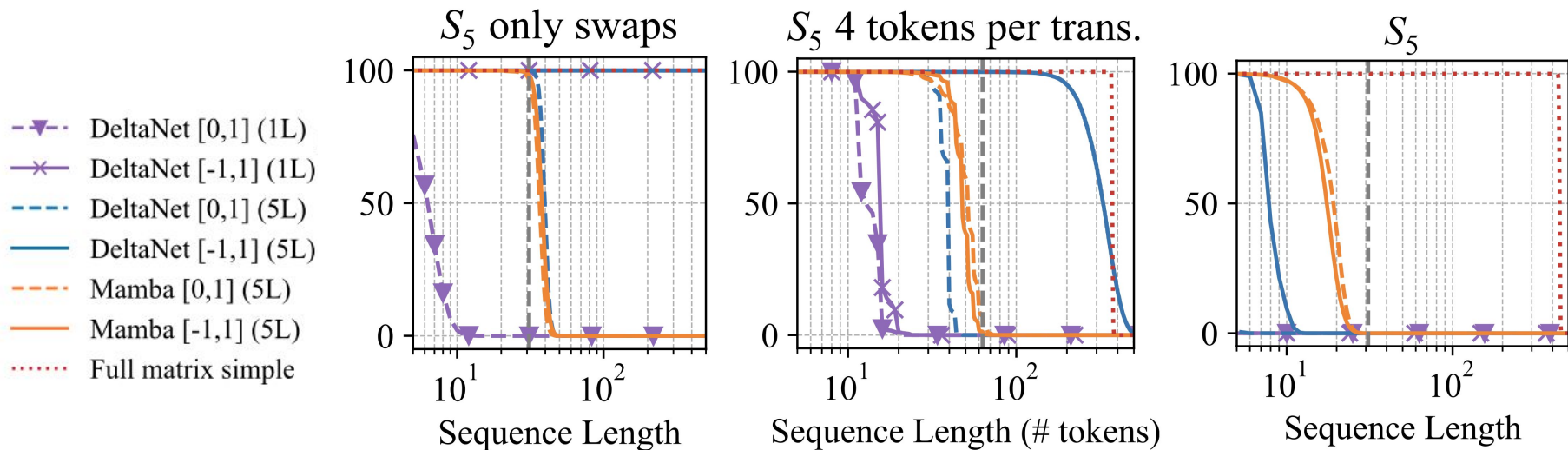


Results:

→ DeltaNet and Mamba can't solve S_5 (expected)

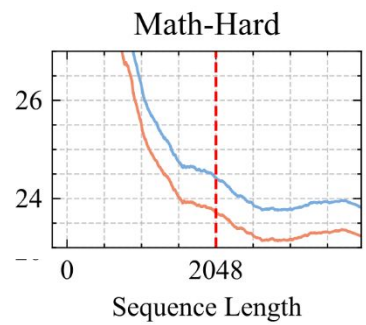
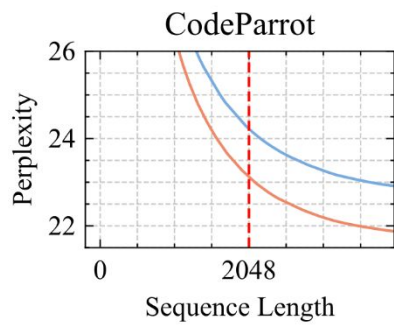
→ Linear RNN with full state-transition matrix can learn to solve S_5

Experiments - Permutation Groups



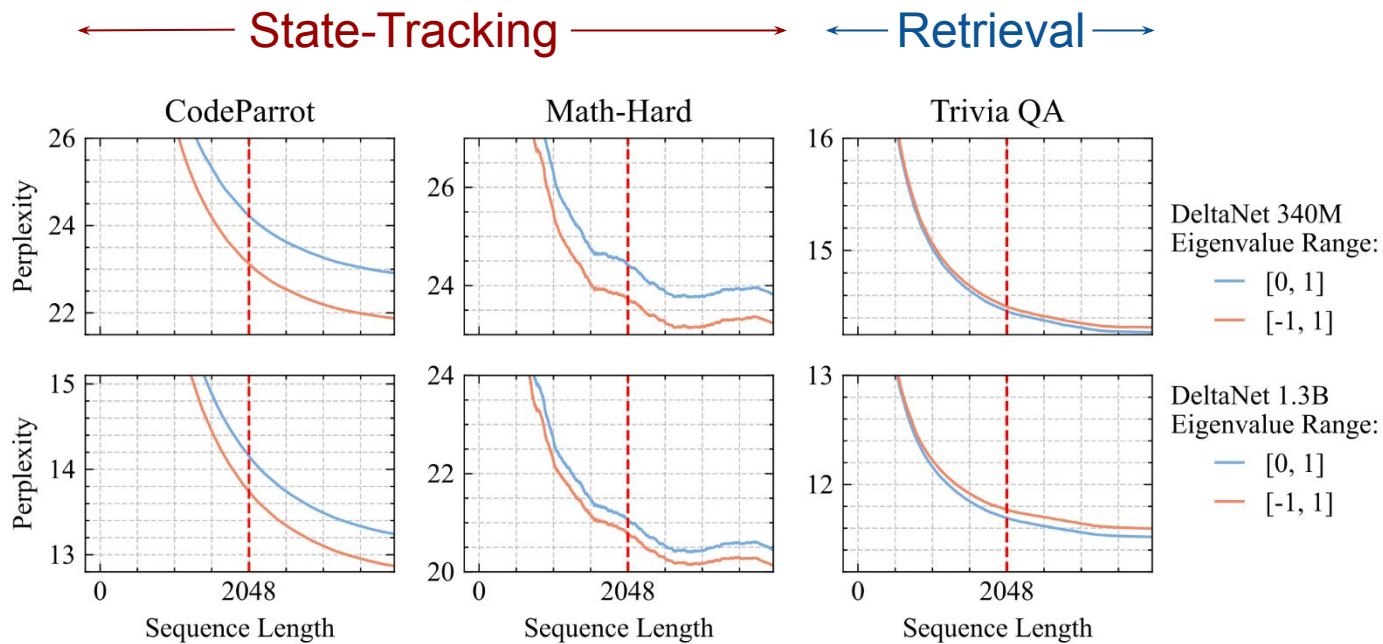
Experiments - Language Modelling

← State-Tracking → ← Retrieval →



DeltaNet 340M
Eigenvalue Range:
— [0, 1]
— [-1, 1]

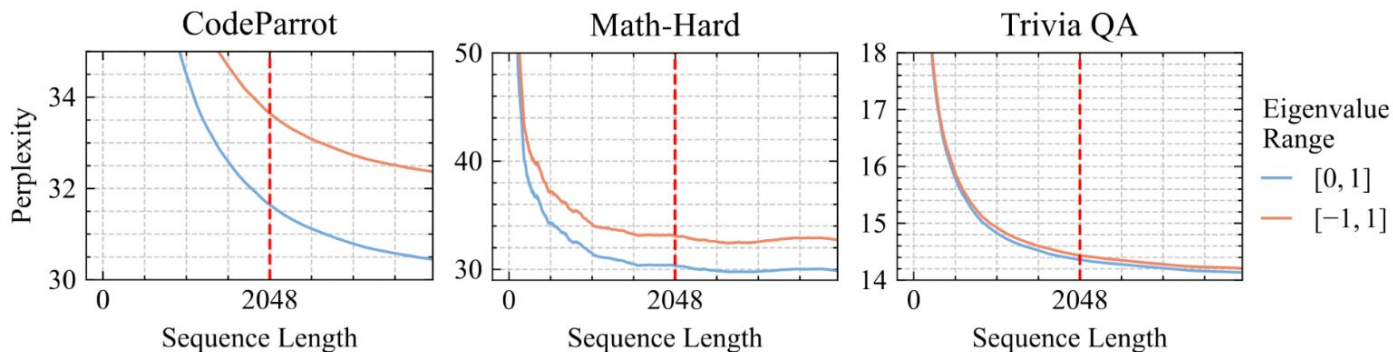
Experiments - Language Modelling



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

Experiments - Language Modelling

Results for **Mamba 370M**:



- *Mamba* doesn't benefit from extended eigenvalue range in language modelling.

Conclusion

- Inclusion of negative eigenvalues expands the expressivity of linear RNNs.
- DeltaNet is promising due to its superior expressivity compared to Mamba.

- Future Directions:
 - Assess real-world improvements in language modeling.
 - Increase expressivity of linear RNNs through more complex state-transition matrices.
 - Understanding the trade-off between associative recall and state-tracking.