Implicit Language Models are RNNs: Balancing Parallelization and Expressivity

Babak Rahmani*,

M Schöne*, H Kremer, F Falck, H Ballani, J Gladrow (*eaual contribution)



Microsoft Research



Implicit Language Modelling

Expressivity

➢ Parallel Training

➤Scalable



Reasoning with LLMs

Puzzles

"swap ball 1 and 3, swap ball 3 and 5, swap ball 4 and 2, .."

Where does ball 5 end up?

Math



Can LLMs learn the algorithm to calculate the next element?

Code Evaluation

1	<pre># Assing variables</pre>
2	x = 3
3	y = 5
4	
5	# Manipulate variables
6	x = x + 2
7	y = y * 2
8	
9	# Loop
9 10	<pre># Loop ~ for i in range(3):</pre>
9 10 11	<pre># Loop </pre> <pre> for i in range(3): x += i </pre>
9 10 11 12	<pre># Loop for i in range(3): x += i y -= 1</pre>
9 10 11 12 13	<pre># Loop for i in range(3):</pre>
9 10 11 12 13 14	<pre># Loop for i in range(3):</pre>
9 10 11 12 13 14 15	<pre># Loop for i in range(3): x += i y -= 1 print("Final x:", x) print("Final y:", y)</pre>

S_5 Permutation composition

Example

Consider the sequence in S_5 : $\sigma_1 = (5, 4, 1, 3, 2), \ \sigma_2 = (3, 4, 1, 2, 5), \ \text{and} \ \sigma_3 = (5, 3, 1, 2, 4).$ The prefix product is computed as follows:

$$\sigma_{1,2} = \sigma_2 \circ \sigma_1 = (1, 3, 5, 4, 2),$$

$$\sigma_{1,2,3} = \sigma_3 \circ \sigma_{1,2} = (2, 5, 1, 3, 4).$$

Hence, the final permutation is $\sigma_{1,2,3} = (2, 5, 1, 3, 4)$.

Task: Given the sequence: **SEQ**>, predict the next element in the sequence.





S_5 Permutation composition

Example

Consider the sequence in S_5 : $\sigma_1 = (5, 4, 1, 3, 2), \ \sigma_2 = (3, 4, 1, 2, 5), \ \text{and} \ \sigma_3 = (5, 3, 1, 2, 4).$ The prefix product is computed as follows:

$$\sigma_{1,2} = \sigma_2 \circ \sigma_1 = (1, 3, 5, 4, 2),$$

$$\sigma_{1,2,3} = \sigma_3 \circ \sigma_{1,2} = (2, 5, 1, 3, 4)$$

Hence, the final permutation is $\sigma_{1,2,3} = (2, 5, 1, 3, 4)$. **Task**: Given the sequence: $\langle SEQ \rangle$, predict the next element in the sequence.





GPT-like models cannot track state

Transformer models are fixed-depth

> This **provably limits** what they can express

Notably, they cannot (really) track state*

*They often learn shortcuts



A Formal Hierarchy of RNN Architectures Merrill, W., et al 2020 The Illusion of State in State-Space Models Merrill, W., et.al 2024

What does it take to solve this task?



Merrill, W., et.al 2024

What does it take to solve this task?



A Formal Hierarchy of RNN Architectures Merrill, W., et al 2020 The Illusion of State in State-Space Models Merrill, W., et.al 2024



Implicit State Space Models

Legend

Iteration steps

Implicit State Space Models







11

Duality between parallel and sequential mode



Training Implicit Models

One training step

- > 1- Fixed-point search:
 - $z^s = f(z^{s-1}, x)$
- 2- Phantom gradients:

Fixed-point search



• $z^{s+1} = (1-\lambda)z^s + \lambda f(z^s, x)$

On Training Implicit Models Geng z et. al. 2022

Why Implicit LLMs are RNNs?

$$\begin{aligned} h_{t} &= \Lambda(x_{t})h_{t-1} + u(x_{t}) \\ y_{t} &= f_{\theta}(h_{t-1}, x_{t}) \end{aligned} \begin{cases} h_{t}^{(s)} &= \Lambda\left(z_{t-1}^{(s-1)}, x_{t}\right)h_{t-1}^{(s)} + u\left(z_{t-1}^{(s-1)}, x_{t}\right) & \text{Eq.1} \\ z_{t-1}^{s} &= f_{\theta}\left(z_{t-1}^{(s-1)}, h_{t}^{(s)}, x_{t}\right) & \text{Eq.2} \end{cases} \\ \begin{cases} h_{t}^{*} &= \Lambda(z_{t-1}^{*}, x_{t})h_{t-1}^{*} + u(z_{t-1}^{*}, x_{t}) & z_{t}^{*} &= \lim_{s \to \infty} z_{t}^{s} \\ z_{t-1}^{*} &= f_{\theta}(z_{t-1}^{*}, h_{t}^{*}, x_{t}) & h_{t}^{*} &= \lim_{s \to \infty} h_{t}^{s} \end{aligned}$$

Theorem 1. Consider an implicit SSM defined by equations (1) and (2). Then the transition function $h_{t-1}^{(*)} \rightarrow h_t^{(*)}$ is **non-linear** and **non-diagonal**. Consequently, the state-to-state Jacobian is a **non-diagonal** operator.

Why Implicit LLMs are RNNs?

Theorem 1. Consider an implicit SSM defined by equations (1) and (2). Then the transition function $h_{t-1}^{(*)} \rightarrow h_t^{(*)}$ is **non-linear** and **non-diagonal**. Consequently, the state-to-state Jacobian is a **non-diagonal** operator.



Lifting the "Illusion of state": Word-problem



Catbabi dataset: a bag of reasoning tasks

Basic Deduction

mice are afraid of cats. cats are afraid of wolves. emily is a cat. wolves are afraid of cats. jessica is a mouse. gertrude is a wolf. sheep are afraid of mice. winona is a wolf. what is winona afraid of ? cat

Path finding

the bathroom is south of the garden. the kitchen is north of the bedroom. the hallway is north of the garden. the bedroom is west of the garden. the office is west of the hallway. how do you go from the hallway to the bedroom? s,w



Implicit SSMs outperform explicit models in Catbabi reasoning



Implicit Large Language Models

Implicit Llama and Mamba2

> Size: 1.3B, 760M, 350M, 130M

Dataset: Pile (207B tokens)

Curriculum Training



Implicit Large Language Models are Parameter Optimal



DEQ LLMs outperform baselines in HellaSwag

A woman is outside with a bucket and a dog. The dog is running around trying to avoid a bath. She...

A. rinses the bucket off with soap and blow dry the dog's head.

B. uses a hose to keep it from getting soapy.

C. gets the dog wet, then it runs away again.

D. gets into a bath tub with the dog.



Downstream task performances

Llama

Model	Dataset/Tokens (B)	D-Pile ppl↓	LAMBADA ppl↓	LAMBADA acc↑	HellaSwag acc↑	PIQA acc↑	Arc-E acc↑	Arc-C acc↑	WinoGrande acc↑	OpenbookQA acc↑	Average acc↑	
Llama	SlimPajama/300	-	9.90	0.5141	0.5216	0.7095	0.5648	0.2875	0.5667	-	0.5274	
Llama [†]	D-Pile/207	8.88	5.77	0.6375	0.5448	0.7171	0.5905	0.2816	0.6054	0.338	0.5307	
Llama (4+1)-ours	D-Pile/207	<u>8.27</u>	<u>5.15</u>	<u>0.6524</u>	0.5853	<u>0.7312</u>	<u>0.6052</u>	0.3097	0.5967	0.356	0.5481	l
Llama (32+4)-ours	D-Pile/207	7.90	4.82	0.6703	0.5995	0.7416	0.6187	0.3012	<u>0.5991</u>	<u>0.344</u>	0.5535	
Llama	SlimPajama/300	-	7.23	0.5744	0.5781	0.7312	0.6279	0.3174	0.5904	-	0.5699	
Llama [†]	D-Pile/207	7.99	4.95	0.6569	0.5936	0.7432	<u>0.6385</u>	0.3217	0.6062	0.352	0.5589	
Llama (4+1)-ours	D-Pile/207	7.66	4.40	0.6852	0.6397	<u>0.7448</u>	0.6338	<u>0.3396</u>	0.6575	0.360	0.5801	
Llama (32+4)-ours	D-Pile/207	7.24	4.24	0.6901	0.6583	0.7465	0.6654	0.3601	<u>0.6401</u>	0.364	0.5892	

Downstream task performances

Llama

	Model	Dataset/Tokens (B)	D-Pile ppl↓	LAMBADA ppl↓	LAMBADA acc↑	HellaSwag acc↑	PIQA acc↑	Arc-E acc↑	Arc-C acc↑	WinoGrande acc↑	OpenbookQA acc↑	Average acc↑	
	Llama	SlimPajama/300	-	9.90	0.5141	0.5216	0.7095	0.5648	0.2875	0.5667	-	0.5274	7
Σ	Llama [†]	D-Pile/207	8.88	5.77	0.6375	0.5448	0.7171	0.5905	0.2816	0.6054	0.338	0.5307	i
90	Llama (4+1)-ours	D-Pile/207	8.27	<u>5.15</u>	<u>0.6524</u>	0.5853	<u>0.7312</u>	<u>0.6052</u>	0.3097	0.5967	0.356	<u>0.5481</u>	ì
\sim	Llama (32+4)-ours	D-Pile/207	7.90	4.82	0.6703	0.5995	0.7416	0.6187	0.3012	0.5991	<u>0.344</u>	0.5535	
	Llama	SlimPajama/300	-	7.23	0.5744	0.5781	0.7312	0.6279	0.3174	0.5904	-	0.5699	Ì
ш	Llama [†]	D-Pile/207	7.99	4.95	0.6569	0.5936	0.7432	<u>0.6385</u>	0.3217	0.6062	0.352	0.5589	
<u>.</u>	Llama (4+1)-ours	D-Pile/207	<u>7.66</u>	<u>4.40</u>	0.6852	0.6397	<u>0.7448</u>	0.6338	<u>0.3396</u>	0.6575	0.360	0.5801	
-	Llama (32+4)-ours	D-Pile/207	7.24	4.24	0.6901	0.6583	0.7465	0.6654	0.3601	<u>0.6401</u>	0.364	0.5892	

Mamba2

Σ	Mamba2	Pile/300	9.23	5.86	0.6167	0.5492	0.7198	0.6103	0.2850	0.6030	0.362	0.5351
00	Mamba2*	D-Pile/207	8.98	6.24	0.6125	0.5418	0.7231	0.6044	0.2858	0.5777	0.338	0.5262
76	Mamba2(4+1)-ours	D-Pile/207	8.60	6.15	0.6117	0.5569	0.7296	0.6077	0.3140	0.5509	0.336	0.5295
	Mamba2(24+4)-ours	D-Pile/207	8.35	<u>5.90</u>	0.6191	0.5698	0.7334	<u>0.6090</u>	0.3131	0.5730	0.338	0.5365
	Mamba2	Pile/300	8.40	5.02	0.6559	0.5995	0.7378	0.6418	0.3319	0.6117	0.378	0.5652
B	Mamba2*	D-Pile/207	8.28	5.12	0.6456	0.5939	<u>0.741</u> 6	0.6145	0.3123	0.6117	0.352	0.5531
<u> </u>	Mamba2(4+1)-ours	D-Pile/207	7.97	5.21	0.6383	0.6136	0.7437	0.6343	0.3302	0.5746	0.354	0.5555
`	Mamba2(24+4)-ours	D-Pile/207	7.70	4.99	0.6489	0.6267	<u>0.7416</u>	<u>0.6423</u>	0.3336	0.5888	0.352	0.5620

Test-time compute performance



Implicit Mamba2 1.3B

Implicit Llama 1.3B



Summary & Outlook

What we did:

- 1. Theoretically showed that Implicit SSMs are RNNs
- 2. Lifting the illusion of state space models
- 3. Scalable Implicit training up to 1.3B
- **4. Improved state tracking** and broader down stream task in llama and Mamba2