

Trained Mamba Emulates Online Gradient Descent in In-Context Linear Regression



Jiarui Jiang^{*1}, Wei Huang^{*2}, Miao Zhang^{†1}, Taiji Suzuki^{3,2}, Liqiang Nie¹

¹Harbin Institute of Technology, Shenzhen, ²RIKEN AIP, ³University of Tokyo



Introduction

What is ICL?

In-context learning (ICL) is a powerful paradigm that enables models to generalize to unseen tasks by dynamically leveraging contextual examples (such as input-output pairs) without task-specific fine-tuning.

e.g.

Task: Translate English to French

Input: I love you

Output: Je t'aime

→ **prompt (context)**

Input: Good morning

Output: Bonjour

The model needs to learn from the context

Input: How are you?

Output:

→ **and then answer this question**

Why study ICL on Mamba?

- ICL is a **critical capacity** for large foundation models.
- Mamba is an **efficient** Transformer alternative with linear complexity,
- but **theoretical understanding** of Mamba's ICL remains limited.

Highlight

- We reveal that trained Mamba can emulate **online gradient descent** in In-Context Linear Regression.
- We establish a **convergence guarantee** for Mamba from random initialization to ICL solution, and further derive the **loss bound** after convergence. The loss matches that of Transformers.

Problem Setup

Data Model

$$\mathbf{x}_i, \mathbf{x}_q, \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$$

For each sequence, we have a new function $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$.

Given prompt $\{\mathbf{x}_i, y_i\}_{i=1}^N$, where $y_i = f(\mathbf{x}_i) = \mathbf{w}^\top \mathbf{x}_i$,

we expect the model to learn the latent function $f(\cdot)$,

and predict $y_q = f(\mathbf{x}_q) = \mathbf{w}^\top \mathbf{x}_q$ for query \mathbf{x}_q

$$\begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_N & \mathbf{x}_q \\ y_1 & y_2 & \cdots & y_N & 0 \end{bmatrix} \quad \text{target: predict } y_q$$

prompt (context)

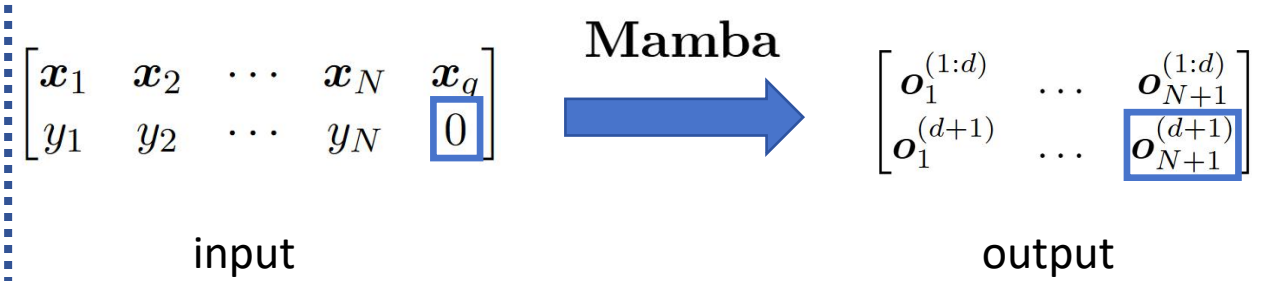
Mamba Model

Mamba Model. We consider a S6 layer of Mamba $\mathbf{o}_{1:L} = \text{Mamba}(\theta; \mathbf{u}_{1:L})$ with selection discretization, and linear recurrence components, where $\mathbf{u}_i, \mathbf{o}_i \in \mathbb{R}^{d_s}$. It can be described as follows:

$$\mathbf{h}_i^{(i)} = \bar{\mathbf{A}}_i \mathbf{h}_{i-1}^{(i)} + \bar{\mathbf{B}}_i \mathbf{u}_i^{(i)} \in \mathbb{R}^{d_h \times 1} \quad (1a) \quad \bar{\mathbf{A}}_i = \exp(\Delta_i \mathbf{A}) \in \mathbb{R}^{d_h \times d_h} \quad (2a)$$

$$\mathbf{o}_i^{(i)} = \mathbf{C}_i^\top \mathbf{h}_i^{(i)}, \quad \mathbf{C}_i \in \mathbb{R}^{d_h \times 1} \quad (1b) \quad \bar{\mathbf{B}}_i = (\Delta_i \mathbf{A})^{-1} (\exp(\Delta_i \mathbf{A}) - \mathbf{I}) \Delta_i \mathbf{B}_i \in \mathbb{R}^{d_h \times 1} \quad (2b)$$

$$\mathbf{B}_i = \mathbf{W}_B \mathbf{u}_i + \mathbf{b}_B \quad (3) \quad \mathbf{C}_i = \mathbf{W}_C \mathbf{u}_i + \mathbf{b}_C \quad (4) \quad \Delta_i = \text{softplus}(\mathbf{w}_\Delta^\top \mathbf{u}_i + \mathbf{b}_\Delta) \quad (5)$$



Use $\hat{y}_q = \mathbf{o}_{N+1}^{(d+1)}$ as the prediction for y_q

Assumptions

Assumption 4.1 (1) Matrix $\mathbf{A} = -\mathbf{I}_{d_h}$. (2) The vector \mathbf{w}_Δ is fixed as zero $\mathbf{0}$, and \mathbf{b}_Δ is fixed as $\ln(\exp((\ln 2)/N) - 1)$. (3) Matrices $\mathbf{W}_B, \mathbf{W}_C$ are initialized with entries drawn i.i.d. from the standard Gaussian distribution $\mathcal{N}(0, 1)$. (4) The hidden state dimension satisfies: $d_h = \Omega(d^2)$. (5) The learning rate satisfies: $\eta = O(d^{-2} d_h^{-1})$. (6) Bias vectors $\mathbf{b}_B, \mathbf{b}_C$ are initialized as zero $\mathbf{0}$. (7) Token length $N = \Omega(d)$.

B, C are initialized with Gaussian distribution

Training Algorithm

Mse Loss + Gradient Descent

$$\mathcal{L}(\theta) = \mathbb{E}_{\mathbf{x}_{1:N}, \mathbf{x}_q, \mathbf{w}} \left[\frac{1}{2} (\hat{y}_q - y_q)^2 \right]$$

$$\theta'(t+1) = \theta'(t) - \eta \nabla_{\theta'} \mathcal{L}(\theta(t))$$

Main Results

Theorem 4.1 Under Assumption 4.1 if the Mamba is trained with gradient descent, and given a new prompt $(\mathbf{e}_1, \dots, \mathbf{e}_N, \mathbf{e}_q)$, then with probability at least $1 - \delta$ for some $\delta \in (0, 1)$, the trainable parameters $\theta'(t) = \{\mathbf{W}_B(t), \mathbf{W}_C(t), \mathbf{b}_B(t), \mathbf{b}_C(t)\}$ converge as $t \rightarrow \infty$ to parameters that satisfies:

$$(a) \text{ Projected hidden state: } (\mathbf{W}_C^\top(t))_{[1:d,:]}(t) \mathbf{h}_i^{(d+1)} = \alpha (\mathbf{W}_C^\top(t))_{[1:d,:]} \mathbf{h}_{i-1}^{(d+1)} + (1 - \alpha) \beta y_i \mathbf{x}_i,$$

$$(b) \text{ Prediction for target: } \hat{y}_q = \mathbf{x}_q^\top \sum_{i=0}^{N-1} (1 - \alpha) \alpha^{i+1} \beta y_{N-i} \mathbf{x}_{N-i},$$

$$(c) \text{ Population loss: } \mathcal{L}(\theta(t)) \leq \frac{3d(d+1)}{2N},$$

$$\text{where } \alpha = \exp((-\ln 2)/N), \beta = \frac{2(1+\alpha)}{\alpha(3(1-\alpha)d+4-2\alpha)}.$$

Theorem 4.1 (a)

Mamba emulates **online GD** for ICL

Theorem 4.1 (b)

The prediction for y_q is a **linear combination** of $y_l \mathbf{x}_q^\top \mathbf{x}_l$

Theorem 4.1 (c)

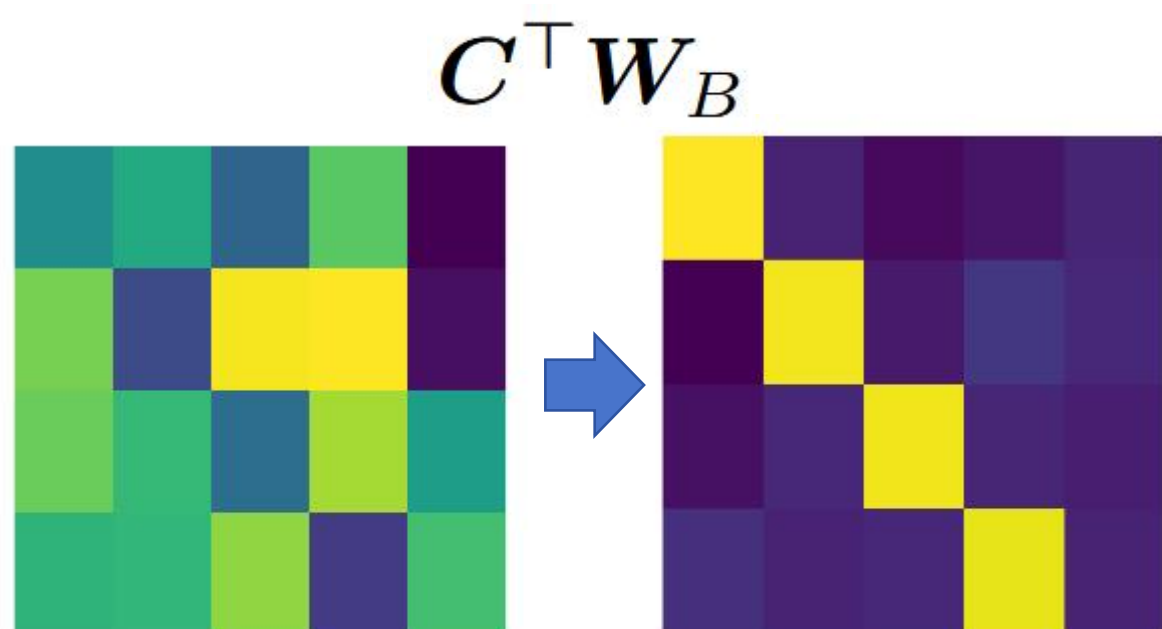
The Loss bound is **linearly** related to the context length N

Compare with Transformer

$$\text{Transformer}(\mathbf{e}_1, \dots, \mathbf{e}_N, \mathbf{e}_q) \approx \mathbf{x}_q^\top \left(\frac{1}{N} \sum_{i=1}^N y_i \mathbf{x}_i \right) \approx \mathbf{x}_q^\top \mathbf{w}.$$

Transformer simulates a **single step** of GD

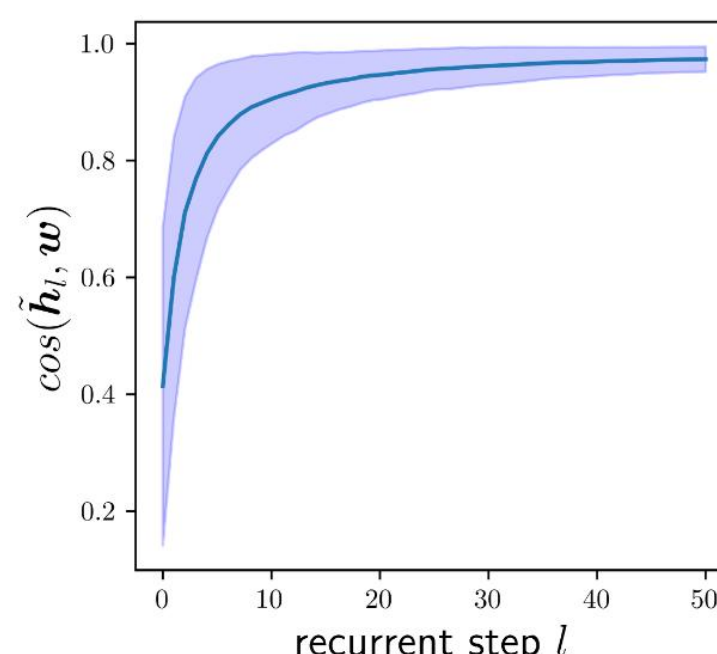
Experimental Results



(a) Initialization

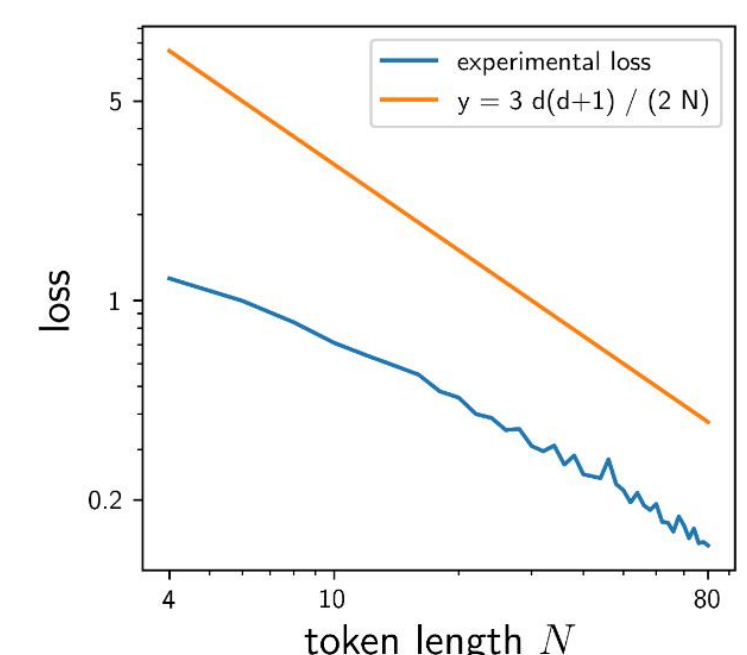
(b) Trained parameter

A Gaussian initialized Mamba can converge to a solution that can perform ICL



$$\text{define } \tilde{\mathbf{h}}_l := (\mathbf{W}_C^\top)_{[1:d,:]} \mathbf{h}_l^{(d+1)}$$

The direction of $\tilde{\mathbf{h}}_l$ converges toward \mathbf{w} as mamba processes multiple prompts



The loss has a linear upper bound, which is comparable to Transformer

Insights & takeaways & future work

- The idea of “**online GD**” can provide insights into designing new architectures
- The theoretical analysis can be extended to **multi-layer Mamba** containing nonlinear layers
- The theoretical analysis can be extended to more **complex and diverse ICL tasks** beyond linear regression

Contact



paper



LinkedIn



Wechat

email: jiaruij@outlook.com