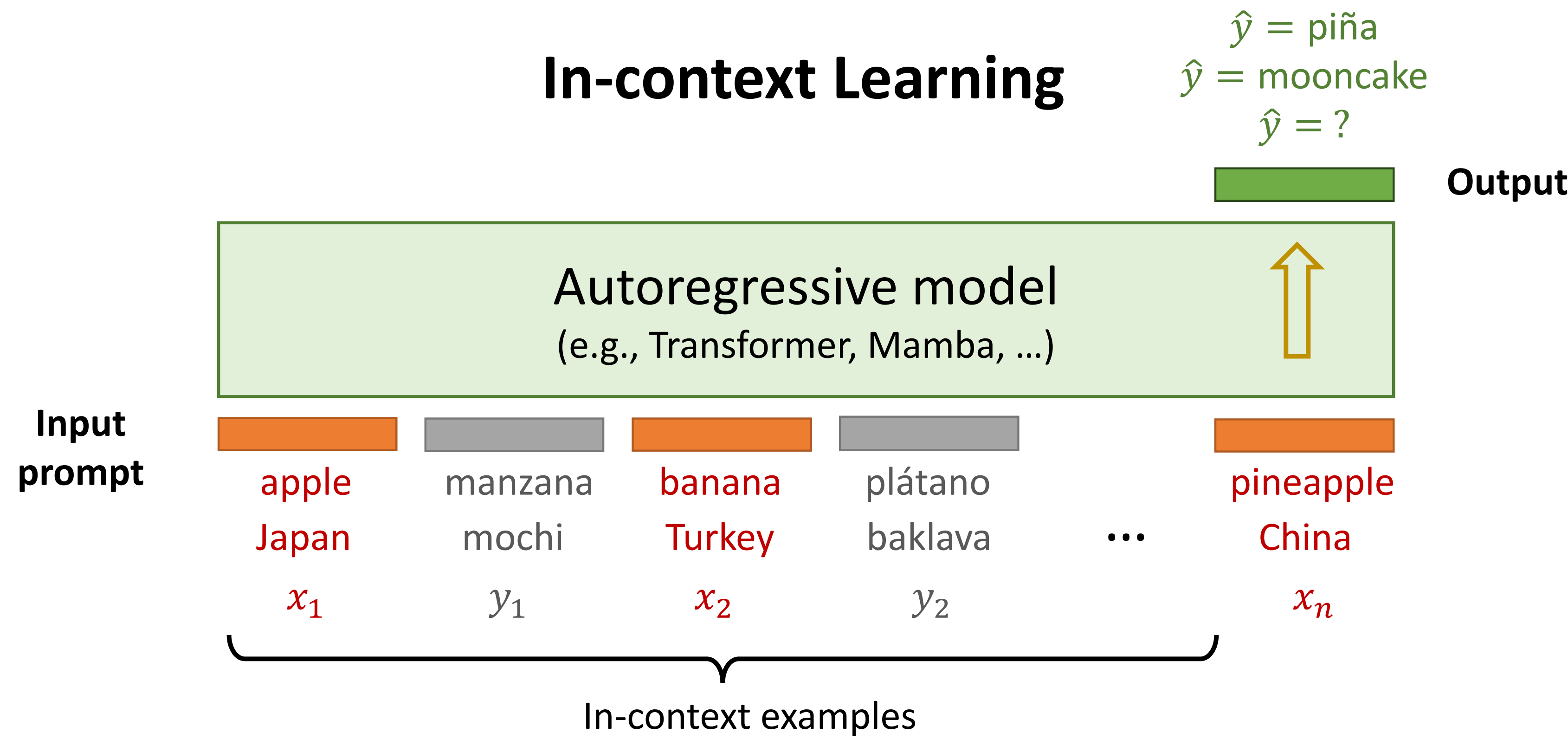


# Gating is Weighting: Understanding Gated Linear Attention through In-context Learning

Yingcong Li<sup>\*1,5</sup> Davoud Ataee Tarzanagh<sup>\*2</sup> Ankit Singh Rawat<sup>3</sup> Maryam Fazel<sup>4</sup> Samet Oymak<sup>1</sup>  
 Equal contribution<sup>\*</sup> University of Michigan<sup>1</sup> Samsung SDS Research America<sup>2</sup> Google Research NYC<sup>3</sup> University of Washington<sup>4</sup> New Jersey Institute of Technology<sup>5</sup>



## Motivation



**Q1:** What optimization algorithms are implemented by different model architectures in ICL?

### Key Findings on ICL Implementations:

Linear attention  $\implies$  Preconditioned Gradient Descent (PGD)<sup>[1,2]</sup>  
 State-space model/H3  $\implies$  Sample-weighted PGD (WPGD)<sup>[1]</sup>  
**Our work:** Gated linear attention  $\implies$  Data-dependent WPGD (DWPGD)<sup>[4]</sup>

**Q2:** What are **Gated Linear Attention (GLA)** architectures?

**Definition 1 (GLA):** Given a sequence of (query, key, value) embeddings  $(\mathbf{q}_i, \mathbf{k}_i, \mathbf{v}_i)_{i=1}^n$ , the GLA recurrence is given by

$$\mathbf{S}_i = \mathbf{G}_i \odot \mathbf{S}_{i-1} + \mathbf{v}_i \mathbf{k}_i^\top, \quad \text{and} \quad \mathbf{o}_i = \mathbf{S}_i \mathbf{q}_i, \quad i \in \{1, \dots, n\}. \quad (\text{GLA})$$

Here, the *gating variable*  $\mathbf{G}_i$  is applied to the state  $\mathbf{S}_{i-1}$  through the Hadamard product  $\odot$ .

**Note:** When  $\mathbf{G}_i \equiv \mathbf{1}$ , (GLA) reduces to causal linear attention.

## Popular GLA Architectures<sup>[3]</sup>:

Model <sup>[3]</sup>	Parameterization
Mamba (Gu & Dao, 2023)	$\mathbf{G}_t = \exp(-(\mathbf{1}^\top \boldsymbol{\alpha}_t) \odot \exp(A)), \boldsymbol{\alpha}_t = \text{softplus}(\mathbf{x}_t W_{\alpha_1} W_{\alpha_2})$
Mamba-2 (Dao & Gu, 2024)	$\mathbf{G}_t = \gamma_t \mathbf{1}^\top \mathbf{1}, \gamma_t = \exp(-\text{softplus}(\mathbf{x}_t W_\gamma) \exp(a))$
mLSTM (Beck et al., 2024; Peng et al., 2021)	$\mathbf{G}_t = \gamma_t \mathbf{1}^\top \mathbf{1}, \gamma_t = \sigma(\mathbf{x}_t W_\gamma)$
Gated Retention (Sun et al., 2024)	$\mathbf{G}_t = \gamma_t \mathbf{1}^\top \mathbf{1}, \gamma_t = \sigma(\mathbf{x}_t W_\gamma)^{\frac{1}{2}}$
DFW (Mao, 2022; Pramanik et al., 2023)	$\mathbf{G}_t = \boldsymbol{\alpha}_t \boldsymbol{\beta}_t^\top, \boldsymbol{\alpha}_t = \sigma(\mathbf{x}_t W_\alpha), \boldsymbol{\beta}_t = \sigma(\mathbf{x}_t W_\beta)$
GateLoop (Katsch, 2023)	$\mathbf{G}_t = \boldsymbol{\alpha}_t^\top \mathbf{1}, \boldsymbol{\alpha}_t = \sigma(\mathbf{x}_t W_{\alpha_1}) \exp(\mathbf{x}_t W_{\alpha_2} \mathbf{t})$
HGRN-2 (Qin et al., 2024b)	$\mathbf{G}_t = \boldsymbol{\alpha}_t^\top \mathbf{1}, \boldsymbol{\alpha}_t = \gamma + (1 - \gamma) \sigma(\mathbf{x}_t W_\alpha)$
RWKV-6 (Peng et al., 2024)	$\mathbf{G}_t = \boldsymbol{\alpha}_t^\top \mathbf{1}, \boldsymbol{\alpha}_t = \exp(-\exp(\mathbf{x}_t W_\alpha))$
<b>Gated Linear Attention (GLA)</b>	$\mathbf{G}_t = \boldsymbol{\alpha}_t^\top \mathbf{1}, \boldsymbol{\alpha}_t = \sigma(\mathbf{x}_t W_{\alpha_1} W_{\alpha_2})^{\frac{1}{2}}$

**Our Goal:** Develop a **mathematical understanding** of the GLA mechanism through the lens of in-context learning and optimization.

## Main Results

• **Input prompt:** Construct input prompt as follows:

$$\mathbf{Z} = [\mathbf{z}_1 \dots \mathbf{z}_n \mathbf{z}_{n+1}]^\top = \begin{bmatrix} \mathbf{x}_1 \dots \mathbf{x}_n & \mathbf{x}_{n+1} \\ y_1 \dots y_n & 0 \end{bmatrix}^\top.$$

### Theorem 1: GLA $\Leftrightarrow$ DWPGD

Consider model construction  $\mathbf{W}_k = \begin{bmatrix} \mathbf{P}_k & 0 \\ 0 & 0 \end{bmatrix}, \mathbf{W}_q = \begin{bmatrix} \mathbf{P}_q & 0 \\ 0 & 0 \end{bmatrix}, \mathbf{W}_v = \begin{bmatrix} \mathbf{0}_{d \times d} & 0 \\ 0 & 1 \end{bmatrix}$ , and take the last coordinate of the last token output denoted by  $(\mathbf{o}_{n+1})_{d+1}$  as a prediction. Then, we have

$$f_{\text{GLA}}(\mathbf{Z}) := (\mathbf{o}_{n+1})_{d+1} = \mathbf{x}^\top \hat{\boldsymbol{\beta}}, \quad \text{where} \quad \hat{\boldsymbol{\beta}} = \mathbf{P}_q (\mathbf{X} \mathbf{P}_k \odot \boldsymbol{\Omega})^\top \mathbf{y}.$$

Here,  $\boldsymbol{\Omega} = [\mathbf{g}_{1:n+1} \dots \mathbf{g}_{n:n+1}]^\top \in \mathbb{R}^{n \times d}$ , where  $\mathbf{g}_{i:n+1}, i \in [n]$  is given by

$$\mathbf{g}_{i:n+1} := \mathbf{g}_{i+1} \odot \mathbf{g}_{i+2} \dots \mathbf{g}_{n+1} \in \mathbb{R}^d, \quad \text{and} \quad \mathbf{G}_i = \begin{bmatrix} * & * \\ \mathbf{g}_i^\top & * \end{bmatrix}.$$

**Note:** When  $\mathbf{G}_i \equiv \mathbf{1}, \hat{\boldsymbol{\beta}} = \mathbf{P}_q \mathbf{P}_k^\top \mathbf{X}^\top \mathbf{y}$  reduces to *preconditioned gradient descent (PGD)*.

## Optimization Landscape of WPGD

• **Learning problem:** Given  $(\mathbf{x}, y, \mathbf{X}, \mathbf{y}) \sim \mathcal{D}$ , learn optimal WPGD:

$$\mathcal{L}_{\text{WPGD}}^* := \min_{\mathbf{P} \in \mathbb{R}^{d \times d}, \boldsymbol{\omega} \in \mathbb{R}^n} \mathbb{E} \left[ \left( y - \mathbf{x}^\top \mathbf{P} \mathbf{X}^\top (\boldsymbol{\omega} \odot \mathbf{y}) \right)^2 \right].$$

• **Data model:** Correlated tasks  $\boldsymbol{\beta}_i \sim \mathcal{N}(0, \mathbf{I})$  jointly Gaussian,  $\mathbf{x}_i \sim \mathcal{N}(0, \boldsymbol{\Sigma}), y_i \sim \mathcal{N}(\mathbf{x}_i^\top \boldsymbol{\beta}_i, \sigma^2)$ .

### Theorem 2: Stationary Points

Define  $h_1 : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  and  $h_2 : [1, \infty) \rightarrow \mathbb{R}_+$  as

$$h_1(\bar{\gamma}) := \left( \sum_{i=1}^n \frac{\lambda_i a_i^2}{(1 + \lambda_i \bar{\gamma})^2} \right) \left( \sum_{i=1}^n \frac{a_i^2}{(1 + \lambda_i \bar{\gamma})^2} \right)^{-1},$$

$$h_2(\gamma) := \left( 1 + M \left( \sum_{i=1}^d \frac{s_i^2}{(M + s_i \gamma)^2} \right) \left( \sum_{i=1}^d \frac{s_i^3}{(M + s_i \gamma)^2} \right)^{-1} \right)^{-1},$$

where  $\{s_i\}, \{\lambda_i\}$  are eigenvalues of  $\boldsymbol{\Sigma}, \mathbf{R}; \{a_i\}$  from  $\mathbf{r} = \mathbf{E} \mathbf{a}$ ;  $M = \sigma^2 + \sum_i s_i$ .

The stationary point  $(\mathbf{P}^*, \boldsymbol{\omega}^*)$  (up to rescaling) is:

$$\mathbf{P}^* = \boldsymbol{\Sigma}^{-\frac{1}{2}} \left( \frac{\gamma^*}{M} \cdot \boldsymbol{\Sigma} + \mathbf{I} \right)^{-1} \boldsymbol{\Sigma}^{-\frac{1}{2}}, \quad \boldsymbol{\omega}^* = (h_2(\gamma^*) \cdot \mathbf{R} + \mathbf{I})^{-1} \mathbf{r},$$

where  $\gamma^*$  is a fixed point of  $h_1(h_2(\gamma)) + 1$ .

### Theorem 3: Global Uniqueness

Under mild spectral gap condition  $\Delta_{\boldsymbol{\Sigma}} \cdot \Delta_{\mathbf{R}} < M + s_{\min}$ :

**T1.** Mapping  $h_1(h_2(\gamma)) + 1$  is a contraction with unique fixed point  $\gamma^*$ .

**T2.** Loss has unique global minimum  $(\mathbf{P}^*, \boldsymbol{\omega}^*)$  up to rescaling.

**Key insight:** Optimal  $\boldsymbol{\omega}^* = (h_2(\gamma^*) \mathbf{R} + \mathbf{I})^{-1} \mathbf{r}$  depends on task correlations  $\mathbf{R}, \mathbf{r}$ , enabling context-aware weighting.

## Optimization Landscape of GLA

**Definition 2(Multi-task Prompt and GLA Objective):** Consider prompt with  $K$  correlated tasks  $(\boldsymbol{\beta}_k)_{k=1}^K$  and one query task  $\boldsymbol{\beta}$ . For each task  $k \in [K]$ , a prompt of length  $n_k$  is drawn, consisting of IID input-label pairs  $\{(\mathbf{x}_i^k, y_i^k)\}_{i=1}^{n_k}$ . Let  $n := \sum_{k=1}^K n_k$ .

$$\mathbf{Z} = \begin{bmatrix} \underbrace{\mathbf{x}_1^1 \dots \mathbf{x}_{n_1}^1 0}_{\text{task 1}} & \underbrace{\mathbf{x}_1^K \dots \mathbf{x}_{n_K}^K 0}_{\text{task K}} & \mathbf{x} \\ \underbrace{y_1^1 \dots y_{n_1}^1 0}_{\text{task 1}} & \underbrace{y_1^K \dots y_{n_K}^K 0}_{\text{task K}} & 0 \\ 0 \dots 0 & \mathbf{c}^1 & 0 \dots 0 & \mathbf{c}^K & 0 \end{bmatrix}^\top. \quad (1)$$

Here,  $\{\mathbf{c}^1, \dots, \mathbf{c}^K\}$  are  $K$  **linearly independent** contextual features. The GLA optimization problem is described as:

$$\mathcal{L}_{\text{GLA}}^* := \min_{\mathbf{P}_k, \mathbf{P}_q \in \mathbb{R}^{d \times d}, G \in \mathcal{G}} \mathcal{L}_{\text{GLA}}(\mathbf{P}_k, \mathbf{P}_q, G) \quad \text{where} \quad \mathcal{L}_{\text{GLA}}(\mathbf{P}_k, \mathbf{P}_q, G) = \mathbb{E} \left[ (y - f_{\text{GLA}}(\mathbf{Z}))^2 \right].$$

Here,  $G(\cdot)$  represents the gating function and  $\mathcal{G}$  denotes the function search space.

Given context examples  $\{(\mathbf{X}_k, \mathbf{y}_k) := (\mathbf{x}_i^k, y_i^k)_{i=1}^{n_k}\}_{k=1}^K$ , define the concatenated data  $(\mathbf{X}, \mathbf{y})$ :

$$\mathbf{X} = [\mathbf{X}_1^\top \dots \mathbf{X}_K^\top]^\top \in \mathbb{R}^{n \times d}, \quad \text{and} \quad \mathbf{y} = [\mathbf{y}_1^\top \dots \mathbf{y}_K^\top]^\top \in \mathbb{R}^n.$$

### Theorem 4 (Optimization Equivalence): GLA $\Leftrightarrow$ DWPGD

Consider GLA with input prompt  $\mathbf{Z}$  defined in (1). There exists a gating function  $G(\cdot)$  such that the optimal risk  $\mathcal{L}_{\text{GLA}}^*$  obeys

$$\mathcal{L}_{\text{GLA}}^* = \mathcal{L}_{\text{WPGD}}^*.$$

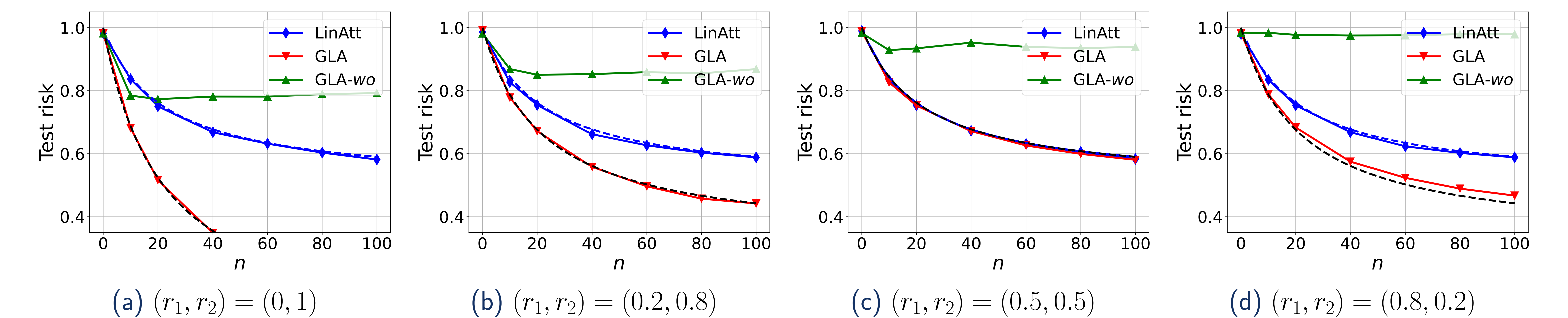
**Corollary (Loss landscape of one-layer linear attention):** Consider a single-layer linear attention following the same model constructions. Let  $\mathbf{R}$  and  $\mathbf{r}$  be the corresponding correlation matrix and vector. Suppose  $\boldsymbol{\Sigma} = \mathbf{I}$ . Then, the optimal risk obeys

$$\mathcal{L}_{\text{ATT}}^* := \min_{\mathbf{P} \in \mathbb{R}^{d \times d}} \mathcal{L}_{\text{WPGD}}(\mathbf{P}, \boldsymbol{\omega} = \mathbf{1}) = d + \sigma^2 - \frac{d(\mathbf{1}^\top \mathbf{r})^2}{n(d + \sigma^2 + 1) + \mathbf{1}^\top \mathbf{R} \mathbf{1}}.$$

## Experiments

**Data setting:**

- $K = 2, n_1 = n_2 = n/2, \boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \boldsymbol{\beta} \sim \mathcal{N}(0, \mathbf{I})$
- $(r_1, r_2) = (\text{corr\_coef}(\boldsymbol{\beta}_1, \boldsymbol{\beta}), \text{corr\_coef}(\boldsymbol{\beta}_2, \boldsymbol{\beta}))$  and  $\text{corr\_coef}(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) = 0$ .
- $G(\mathbf{z}) = \sigma(\mathbf{W}_g \mathbf{z}) \mathbf{1}^\top$  where  $\sigma(x) = (1 + e^{-x})^{-1}$  is the activation function.



• LinAtt: Linear attention; GLA: Gated linear attention; GLA-wo: GLA without contextual features

[1] Li Y, Rawat AS, Oymak S. Fine-grained analysis of in-context linear estimation: Data, architecture, and beyond. NeurIPS 2024.

[2] Ahn K, Cheng X, Daneshmand H, Sra S. Transformers learn to implement preconditioned gradient descent for in-context learning. NeurIPS 2023.

[3] Yang S, Wang B, Shen Y, Panda R, Kim Y. Gated linear attention transformers with hardware-efficient training. ICML 2024.

[4] Li Y, Tarzanagh DA, Fazel M, Oymak S. Gating is weighting: Understanding gated linear attention through in-context learning. COLM 2025.