

# Transformers as Support Vector Machines

Davoud Ataee Tarzanagh<sup>1</sup> Yingcong Li<sup>2,4</sup> Christos Thrampoulidis<sup>3</sup> Samet Oymak<sup>4</sup>

University of Pennsylvania<sup>1</sup> University of California, Riverside<sup>2</sup> University of British Columbia<sup>3</sup> University of Michigan<sup>4</sup>

## Question

Can we characterize the optimization landscape and implicit bias of Transformers' attention mechanism?

## Optimization Methods

- **W-parameterization**: **Gradient Descent** with stepsize  $\eta > 0$ :

$$\mathbf{W}(k+1) = \mathbf{W}(k) - \eta \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}(k)), \quad (\text{GD-W})$$

- **(K, Q)-parameterization**: **Regularization Path** with radius  $R > 0$ :

$$(\mathbf{K}_R, \mathbf{Q}_R) = \arg \min_{\|\mathbf{K}\|_F^2 + \|\mathbf{Q}\|_F^2 \leq 2R} \mathcal{L}(\mathbf{K}, \mathbf{Q}). \quad (\text{RP-KQ})$$

## Softmax-Attention

$$f(\mathbf{X}) = h(\mathbf{X}^\top \mathbb{S}(\mathbf{X} \mathbf{Q} \mathbf{K}^\top \mathbf{X}^\top))$$

- $\mathbf{K}, \mathbf{Q} \in \mathbb{R}^{d \times m}$ ,  $\mathbf{W} := \mathbf{K} \mathbf{Q}^\top$ : attention weights,
- $\mathbb{S}(\cdot)$ : softmax function,  $h(\cdot)$ : prediction head.

**Problem Description**: Given training dataset  $(Y_i, \mathbf{X}_i, \mathbf{z}_i)_{i=1}^n$  where  $Y_i \in \{-1, 1\}$ ,  $\mathbf{X}_i \in \mathbb{R}^{T \times d}$  and  $\mathbf{z}_i \in \mathbb{R}^d$ , we explore the training risk with a loss  $\ell$  as follows:

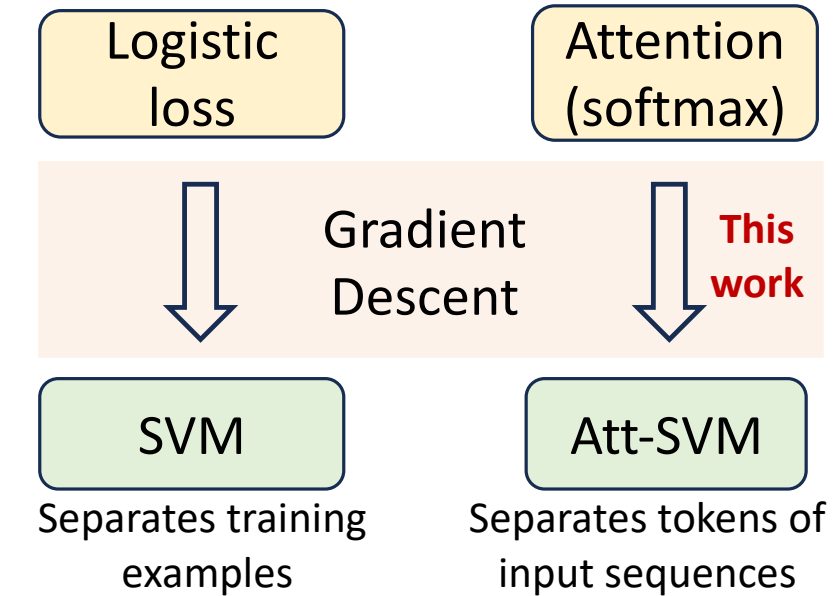
$$\mathcal{L}(\mathbf{K}, \mathbf{Q}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i \cdot f(\mathbf{X}_i)), \quad (\text{ERM})$$

$$\text{where } f(\mathbf{X}_i) = h(\mathbf{X}_i^\top \mathbb{S}(\mathbf{X}_i \mathbf{K} \mathbf{Q}^\top \mathbf{z}_i)).$$

## Motivation

Exploring implicit bias is a key step in unraveling the generalization of the (softmax-)attention mechanism.

## Conclusion



**Transformers are SVMs!**

## Attention SVM

For given indices of **selected tokens**  $\alpha = (\alpha_i)_{i=1}^n$ , define

- **SVM for W-parameterization**:

$$\mathbf{W}^{mm} = \arg \min_{\mathbf{W}} \|\mathbf{W}\|_F$$

$$\text{s.t. } (\mathbf{x}_{i\alpha_i} - \mathbf{x}_{it})^\top \mathbf{W} \mathbf{z}_i \geq 1, \quad \forall i, t \ (t \neq \alpha_i)$$

- **SVM for (K, Q)-parameterization** ( $\mathbf{W} := \mathbf{K} \mathbf{Q}^\top$ ):

$$\mathbf{W}_*^{mm} \in \arg \min_{\text{rank}(\mathbf{W}) \leq m} \|\mathbf{W}\|_*$$

$$\text{s.t. } (\mathbf{x}_{i\alpha_i} - \mathbf{x}_{it})^\top \mathbf{W} \mathbf{z}_i \geq 1, \quad \forall i, t \ (t \neq \alpha_i)$$

## Convergence of Attention Weights with Linear Head

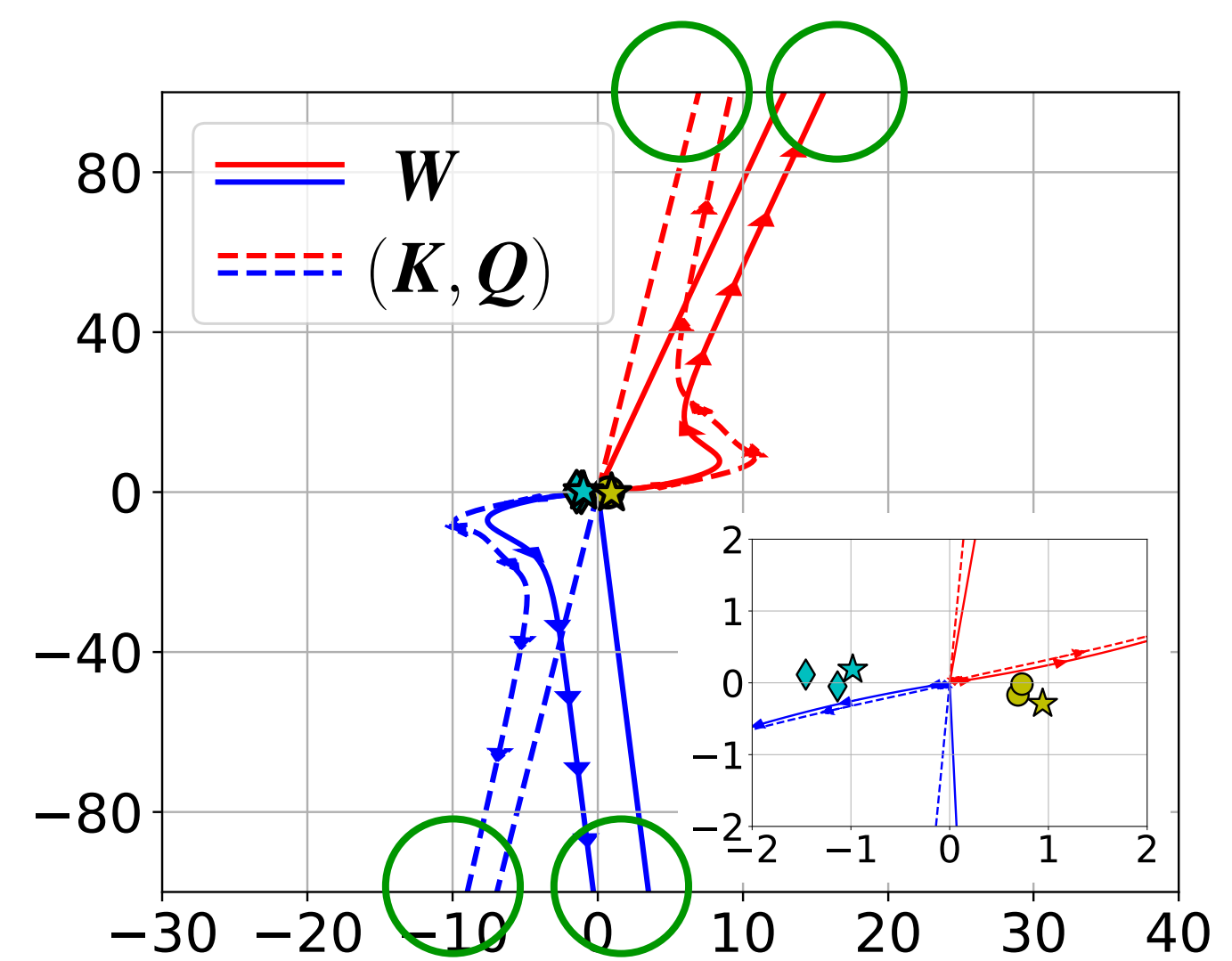
**Assumptions**: Over any bounded interval  $[a, b]$ :

- 1  $\ell : \mathbb{R} \rightarrow \mathbb{R}$  is strictly decreasing; and
- 2  $\ell'$  is Lipschitz continuous and bounded.

### Theorem I: Convergence of Gradient Descent

Under Assumptions 1&2 and  $h(\mathbf{x}) = \mathbf{v}^\top \mathbf{x}$ , **GD-W** with a fixed  $\eta$  and proper starting point satisfy  $\lim_{k \rightarrow \infty} \|\mathbf{W}(k)\|_F = \infty$  and  $\lim_{k \rightarrow \infty} \mathbf{W}(k) / \|\mathbf{W}(k)\|_F = \mathbf{W}^{mm} / \|\mathbf{W}^{mm}\|_F$ .

- **Regularized path**: Under Assumptions 1&2 and  $h(\mathbf{x}) = \mathbf{v}^\top \mathbf{x}$ , **RP-KQ** satisfies  $\lim_{R \rightarrow \infty} \frac{\mathbf{K}_R \mathbf{Q}_R^\top}{R} = \frac{\mathbf{W}_*^{mm}}{\|\mathbf{W}_*^{mm}\|_F}$ .



- Arrows: GD trajectories.
- Lines: the SVM directions mapped to  $\mathbf{z}$ , e.g.,  $\mathbf{W} \mathbf{z}$ .

## Implicit Bias of Attention with Nonlinear Head

**Q**: What is the implicit bias and the form of  $\mathbf{W}(k)$  when the GD solution is composed by multiple tokens?

### General SVM

$$\mathbf{W}(k) \approx \mathbf{W}^{fin} + \|\mathbf{W}(k)\|_F \cdot \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F}$$

**Finite component ( $\mathbf{W}^{fin}$ )**:

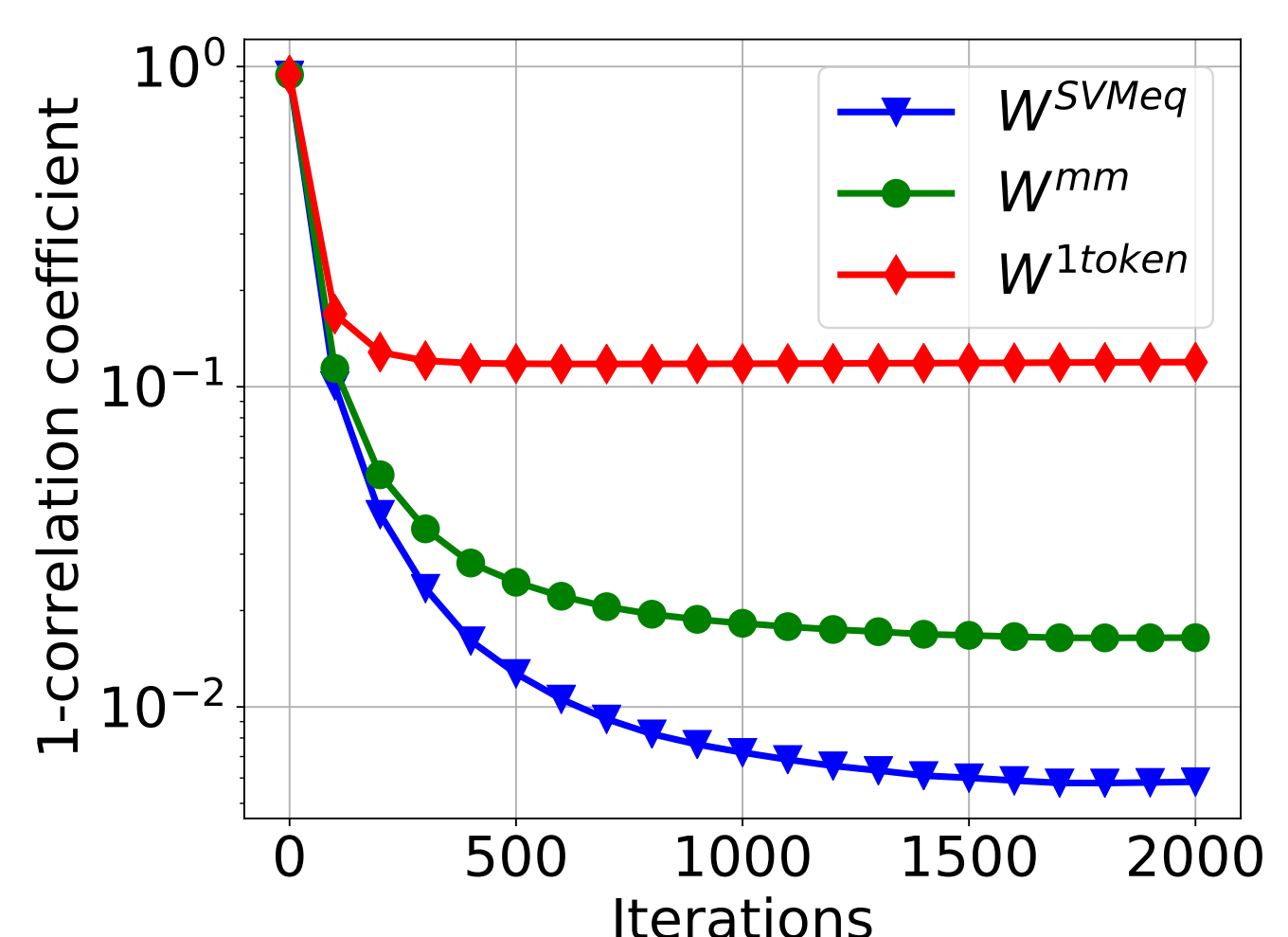
$$(\mathbf{x}_{it} - \mathbf{x}_{i\tau})^\top \mathbf{W}^{fin} \mathbf{z}_i = \log(s_{it}/s_{i\tau}) \quad \forall t, \tau \in \mathcal{O}_i, \ i \in [n].$$

**Directional component ( $\mathbf{W}^{mm}$ )**:

$$\mathbf{W}^{mm} = \arg \min_{\mathbf{W}} \|\mathbf{W}\|_F$$

$$\text{s.t. } \begin{cases} \forall t \in \mathcal{O}_i, \tau \in \bar{\mathcal{O}}_i : (\mathbf{x}_{it} - \mathbf{x}_{i\tau})^\top \mathbf{W} \mathbf{z}_i \geq 1, \\ \forall t, \tau \in \mathcal{O}_i : (\mathbf{x}_{it} - \mathbf{x}_{i\tau})^\top \mathbf{W} \mathbf{z}_i = 0. \end{cases}$$

- $\mathcal{O}_i, i \in [n]$ : Sets of relevant tokens.
- $s_{it}, t \in \mathcal{O}_i$ : Assigned softmax probabilities.



$$\mathbf{W}^{SVMeq} = \mathbf{W}^{fin} + C \cdot \mathbf{W}^{mm}$$

$$\text{where } C = \arg \max \langle \mathbf{W}^{SVMeq}, \mathbf{W}^{GD} \rangle.$$