# Unraveling Attention via Convex Duality: Analysis and Interpretations of Vision Transformers

Arda Sahiner[1], Tolga Ergen[1], Batu Ozturkler[1], Burak Bartan[1], John Pauly[1], Morteza Mardani[2], Mert Pilanci[1]

[1]Department of Electrical Engineering, Stanford University, Stanford, CA
[2]NVIDIA Corporation, Santa Clara, CA

# Introduction

- Vision transformer architectures very successful at image tasks
- A variety of attention mechanisms have been proposed: self-attention, MLP-Mixer, Fourier Neural Operator (FNO), and more!
- However, these architectures are not theoretically understood

# Contributions

- ▶ Prove that self-attention, MLP-Mixer, and FNO with linear and ReLU activation can be solved to their global optima by demonstrating their equivalence to convex optimization problems.
- ▶ Provide interpretability to the optimization objectives of these attention modules.
- ▶ Validate the (convex) vision transformers perform better than baseline convex methods in a transfer learning task (see paper).

# Preliminaries: Background

▶ Training data $\{X_i \in \mathbb{R}^{s \times d}\}_{i=1}^n$, corresponding labels of arbitrary size $\{Y_i \in \mathbb{R}^{r \times c}\}_{i=1}^n$

▶ Solve the optimization problem

$$p^* := \min_\theta \sum_{i=1}^n \mathcal{L}\left(f_\theta(X_i), Y_i\right) + \mathcal{R}(\theta) \tag{1}$$

▶ Can include classification, where $r = 1$, or for regression, where $r = s$, one can directly use squared loss or other convex loss functions.

▶ One may also use this formulation to apply to both supervised and self-supervised learning.

# Single Block of Multi-Head Self-Attention

▶ $j$th self-attention head given by

$$f_j(X_i) := \sigma \left( \frac{X_i Q_j K_j^\top X_i^\top}{\sqrt{d}} \right) X_i V_j, \qquad (2)$$

▶ Multi-head self-attention

$$f_{MHSA}(X_i) := \begin{bmatrix} f_1(X_i) & \cdots & f_m(X_i) \end{bmatrix} W$$

$$= \sum_{j=1}^{m} \sigma \left( \frac{X_i Q_j K_j^\top X_i^\top}{\sqrt{d}} \right) X_i V_j W_j \qquad (3)$$

▶ Can simplify as

$$f_{MHSA}(X_i) := \sum_{j=1}^{m} \sigma \left( \frac{X_i W_{1j} X_i^\top}{\sqrt{d}} \right) X_i W_{2j}. \qquad (4)$$

# Convexity of Linear Multi-Head Self-Attention

## Theorem

*Pose the non-convex weight-decay linear-activation multi-head self-attention training problem*
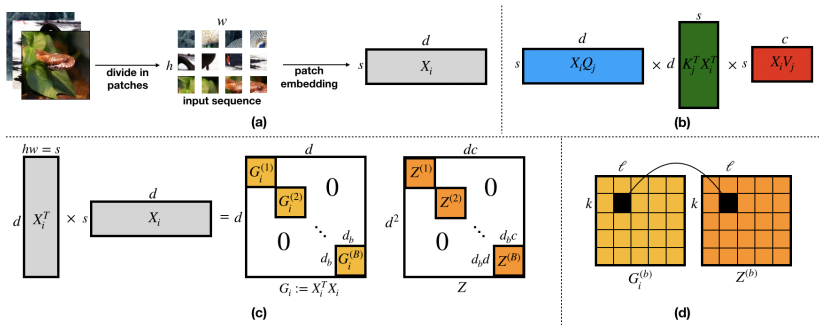
$$p_{SA}^* := \min_{W_{1j}, W_{2j}} \sum_{i=1}^{n} \mathcal{L}(\sum_{j=1}^{m} X_i W_{1j} X_i^\top X_i W_{2j}, Y_i)$$
$$+ \frac{\beta}{2} \sum_{j=1}^{m} \|W_{1j}\|_F^2 + \|W_{2j}\|_F^2. \tag{5}$$

*Then, for $\beta > 0$ and $m \geq m^*$ where $m^* \leq \min\{d^2, dc\}$, this is equivalent to a convex optimization problem*

$$p_{SA}^* = \min_{Z \in \mathbb{R}^{d^2 \times dc}} \sum_{i=1}^{n} \mathcal{L}\left(\sum_{k=1}^{d} \sum_{\ell=1}^{d} G_i[k, \ell] X_i Z^{(k,\ell)}, Y_i\right) + \beta \|Z\|_* \tag{6}$$

*where $G_i := X_i^\top X_i$ and $Z^{(k,\ell)} \in \mathbb{R}^{d \times c}$.*

# Interpretation of Linear Multi-Head Self-Attention



Figure 1: (a) Input image is first divided into $hw = s$ patches, where each patch is represented by a latent vector of dimension $d$. (b) The (non-convex) scaled dot-product self-attention applies learnable weights $Q_j$, $K_j$, $V_j$ to the patch embeddings $X_i$. (c) In the equivalent convex optimization problem for the self-attention training objective, the Gram matrix $G_i$ is formed that groups latent features in $B$ different blocks, (d) and accordingly the nuclear norm regularization is imposed on the dual variables $Z$ based on the similarity scores $G_i[k, l]$.

# Conclusion

A similar procedure can be used to analyze self-attention with ReLU, as well as

- ▶ MLP-Mixer
- ▶ FNO
- ▶ A modification of FNO called block-FNO (BFNO)

with both ReLU and linear activations (see paper). In summary,

- ▶ Studied the vision transformer problem by finding convex equivalents to single attention blocks.
- ▶ These dual forms provide new interpretations, and provide better convex solvers than previously formulated