

## Tensor Attention Training:

# Provably Efficient Learning of Higher-order Transformers

Yingyu Liang, Zhenmei Shi, Zhao Song, Yufa Zhou



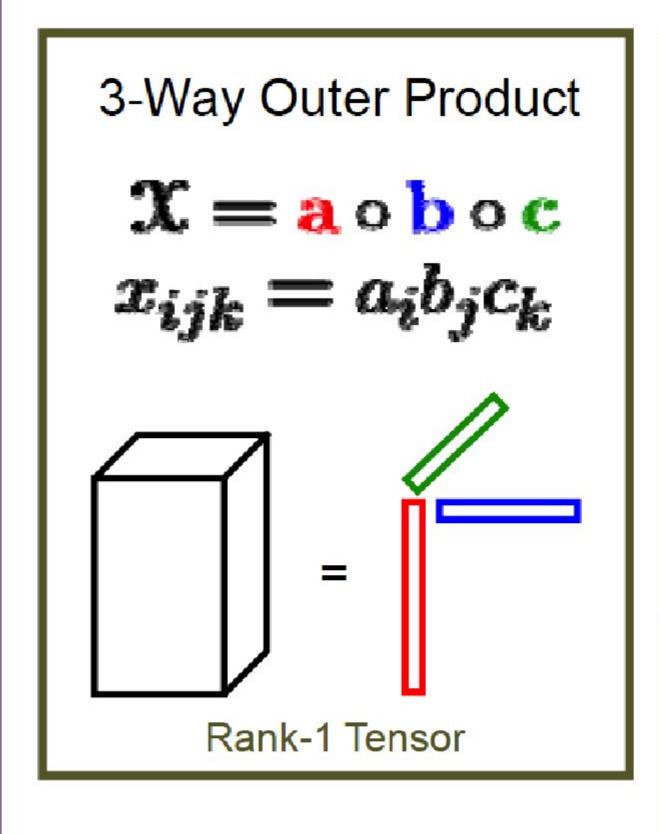


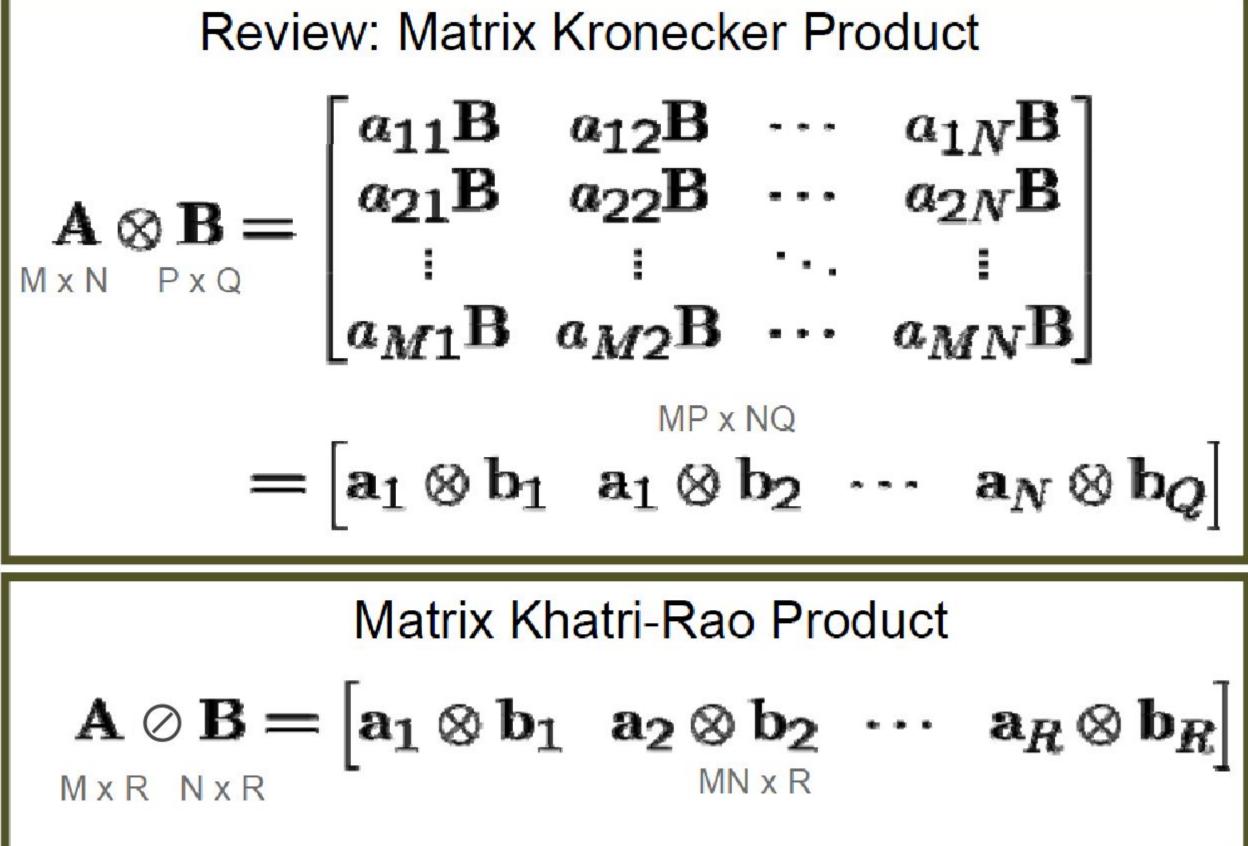


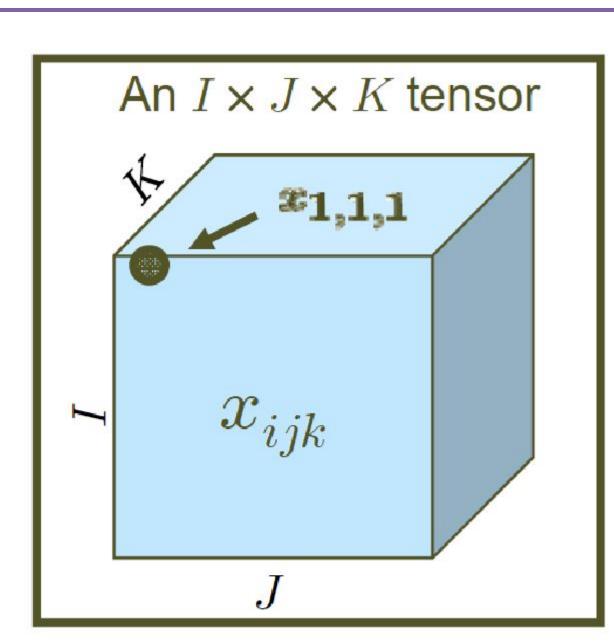


#### Background

⊗ Kronecker product and ⊘ Kathri-Rao product



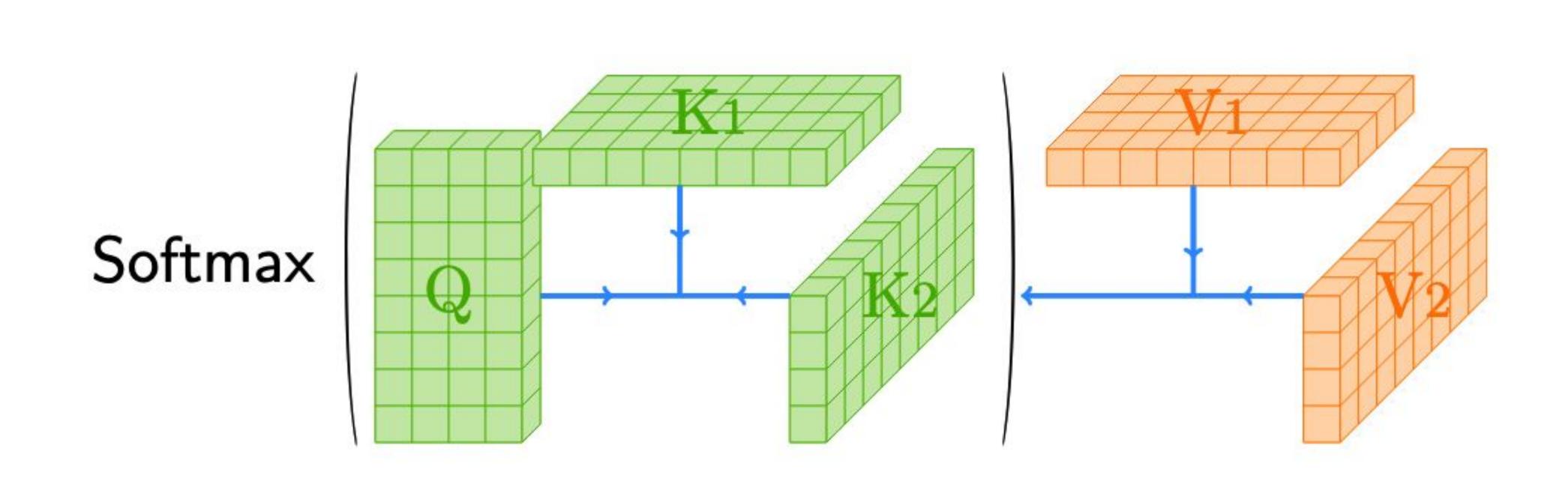




 $3^{\rm rd}$  order tensor mode 1 has dimension I mode 2 has dimension J mode 3 has dimension K

Source: CSE 6363 Machine Learning from UT Arlington

#### **Motivation**



Tensor Attention, defined as  $\mathsf{Softmax}(Q(K_1 \oslash K_2)^\top)(V_1 \oslash V_2)$ , is a higher-order generalization of matrix attention that can capture high-order/multi-view information intrinsically. Meanwhile, it faces a cubic computational complexity bottleneck. Therefore, in this work, we pose the following question:

Can we achieve almost linear time for gradient computation in Tensor Attention Training?

#### **Problem Setup**

#### **Definition 1** (Tensor attention optimization)

Suppose  $A_1,A_2,A_3,A_4,A_5,E\in\mathbb{R}^{n\times d}$  and  $Y_1,Y_2\in\mathbb{R}^{d\times d}$  are given. Let  $D(X)=\mathrm{diag}(\exp(A_1X(A_2\otimes A_3)^\top/d)\mathbf{1}_{n^2})\in\mathbb{R}^{n\times n}$  and  $Y=Y_1\otimes Y_2\in\mathbb{R}^{d^2\times d}$ . We formulate the attention optimization problem as:

 $\min_{X \in \mathbb{R}^{d \times d^2}} \mathsf{Loss}(X) := 0.5 \|D(X)^{-1} \exp(A_1 X (A_2 \otimes A_3)^\top / d) (A_4 \otimes A_5) Y - E\|_F^2.$ 

**Definition 2** (Approximate Tensor Attention Loss Gradient Computation (ATAttLGC $(n,d,B,\epsilon)$ )

Suppose  $A_1, A_2, A_3, A_4, A_5, E \in \mathbb{R}^{n \times d}$  and  $X_1, X_2, X_3, Y_1, Y_2 \in \mathbb{R}^{d \times d}$ . Let  $X = X_1 \cdot (X_2 \oslash X_3)^{\top} \in \mathbb{R}^{d \times d^2}$ . Let  $\epsilon, B > 0$ . Assume that  $\max\{\|A_1X_1\|_{\infty}, \|A_2X_2\|_{\infty}, \|A_3X_3\|_{\infty}, \|A_4Y_1\|_{\infty}, \|A_5Y_2\|_{\infty}\} \leq B$ . Let us assume that any numbers in the previous matrices are in the  $\log(n)$  bits model. Then, our target is to output a matrix  $\widetilde{g} \in \mathbb{R}^{d \times d^2}$  to approximate the gradient of the loss function in **Definition 1**, satisfying  $\|\widetilde{g} - \frac{\mathrm{dLoss}(X)}{\mathrm{d} X}\|_{\infty} \leq \epsilon$ .

$$\min_{X \in \mathbb{R}^{d \times d^2}} 0.5 \left\| \left( n \right\|_{X}^{n} \right) \times \exp \left( n \left\| \frac{d}{A_1} \right\|_{X}^{d} \times d^2 \right\|_{X}^{d^2} \times d^2 \left\| \frac{d^2}{(A_2 \otimes A_3)^\top} \right\|_{X}^{n} \times d^2 \left\| \frac{d}{A_1} \right\|_{X}^{d} \times d^2 \left\| \frac{d}{A_1} \right\|_{X}^{n} \times d^2 \left\| \frac{d}{A_2} \right\|_{X}^{n} \times d^2 \left\| \frac{$$

### **Main Results**

#### Theorem 1 (Fast gradient computation)

Assume that any numbers in the matrices are in the  $\log(n)$  bits model. Then, there exist an algorithm that runs in almost linear time  $n^{1+o(1)}$  to solve

$$\mathsf{ATAttLGC}(n, d = O(\log n), B = o(\sqrt[3]{\log n}), \epsilon = 1/\operatorname{poly}(n)).$$

#### **Theorem 2** (Hardness)

Assume Strong Exponential Time Hypothesis (**SETH**). Let  $\gamma:\mathbb{N}\to\mathbb{N}$  be any function with  $\gamma(n)=o(\log n)$  and  $\gamma(n)=\omega(1)$ . For any constant  $\delta>0$ , when  $E=0,\ \mathbf{Y}=\mathbf{I}_d,\ \mathbf{X}=\lambda\mathbf{I}_d$  for some scalar  $\lambda\in[0,1]$ , it is impossible in  $O(n^{3-\delta})$  time to solve

$$\mathsf{ATAttLGC}(n, d = \Theta(\log n), B = \Theta(\sqrt[3]{\gamma(n) \cdot \log n}), \epsilon = O(1/(\log n)^4)).$$