## How Transformers Perform PCA and Sparse Recovery

Qinyan Liu

July 6, 2025

Qinyan Liu

How Transformers Perform PCA and Spa

July 6, 2<u>025</u>

#### Transformers

- Encoder-based Transformers
- Decoder-based Transformers

2 How Transformers Perform PCA

How Transformers Perform Sparse Recovery: A brief example of decoder-based Transformers

Aspect	Encoder-based Transformer	Decoder-based Transformer
Primary Use	Understanding	Generation
Attention Type	Bidirectional	Causal (unidirectional)
Masking	No mask (full context)	Causal mask (only past tokens visible)
Input	Full sequence at once	Left-to-right, token by token
Training Objec-	Masked Language Modeling	Autoregressive Language Mod-
tive	(MLM)	eling
Example Mod- els	BERT, RoBERTa	GPT, ChatGPT
Output	Contextual embeddings	Predicted next tokens

3)) J

(Attention layer). A (self-)attention layer with M heads is denoted as  $Attn_{\theta}(\cdot)$  with parameters  $\theta = \{(\mathbf{V}_m, \mathbf{Q}_m, \mathbf{K}_m)\}_{m \in [M]} \subset \mathbb{R}^{D \times D}$ . On any input sequence  $\mathbf{H} \in \mathbb{R}^{D \times N}$ ,

$$\widetilde{\mathbf{H}} = \operatorname{Attn}_{\boldsymbol{\theta}}(\mathbf{H}) := \mathbf{H} + \frac{1}{N} \sum_{m=1}^{M} \left( \mathbf{V}_m \mathbf{H} \cdot \sigma \left( (\mathbf{Q}_m \mathbf{H})^\top (\mathbf{K}_m \mathbf{H}) \right) \right) \in \mathbb{R}^{D \times N},$$

where  $\sigma: \mathbb{R} \to \mathbb{R}$  is the ReLU function. In vector form,

$$\widetilde{\mathbf{h}}_{i} = [\operatorname{Attn}_{\boldsymbol{\theta}}(\mathbf{H})]_{i} = \mathbf{h}_{i} + \sum_{m=1}^{M} \frac{1}{N} \sum_{j=1}^{N} \sigma\left(\langle \mathbf{Q}_{m} \mathbf{h}_{i}, \mathbf{K}_{m} \mathbf{h}_{j} \rangle\right) \cdot \mathbf{V}_{m} \mathbf{h}_{j}.$$

(MLP layer). A (token-wise) MLP layer with hidden dimension D' is denoted as

 $\mathrm{MLP}_{\boldsymbol{\theta}}(\cdot)$  with parameters  $\boldsymbol{\theta} = (\mathbf{W}_1, \mathbf{W}_2) \in \mathbb{R}^{D' \times D} \times \mathbb{R}^{D \times D'}$ .

On any input sequence  $\mathbf{H} \in \mathbb{R}^{D \times N}$ ,

$$\widetilde{\mathbf{H}} = \mathrm{MLP}_{\boldsymbol{\theta}}(\mathbf{H}) := \mathbf{H} + \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{H}),$$

where  $\ \sigma:\mathbb{R}\to\mathbb{R}$  is the ReLU function. In vector form, we have

 $\widetilde{\mathbf{h}}_i = \mathbf{h}_i + \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{h}_i).$ 

(Transformer). An L-layer Transformer, denoted as  $TF_{\theta}(\cdot)$ , is a composition of L self-attention layers each followed by an MLP layer:

$$\mathbf{H}^{(L)} = \mathrm{TF}_{\boldsymbol{\theta}}(\mathbf{H}^{(0)}), \text{ where } \mathbf{H} = \mathbf{H}^{(0)} \in \mathbb{R}^{D \times N}$$

is the input sequence, and

$$\mathbf{H}^{(\ell)} = \mathrm{MLP}_{\boldsymbol{\theta}_{\mathrm{mlp}}^{(\ell)}} \left( \mathrm{Attn}_{\boldsymbol{\theta}_{\mathrm{attn}}^{(\ell)}} (\mathbf{H}^{(\ell-1)}) \right), \quad \ell \in \{1, \dots, L\}.$$

Above, the parameter  $\boldsymbol{\theta} = (\boldsymbol{\theta}_{attn}^{(1:L)}, \boldsymbol{\theta}_{mlp}^{(1:L)})$  consists of the attention layers  $\boldsymbol{\theta}_{attn}^{(\ell)} = \{(\mathbf{V}_m^{(\ell)}, \mathbf{Q}_m^{(\ell)}, \mathbf{K}_m^{(\ell)})\}_{m \in [M]} \subset \mathbb{R}^{D \times D}$  and the MLP layers  $\boldsymbol{\theta}_{mlp}^{(\ell)} = (\mathbf{W}_1^{(\ell)}, \mathbf{W}_2^{(\ell)}) \in \mathbb{R}^{D \times D'} \times \mathbb{R}^{D' \times D}.$ 

$$TF_{\boldsymbol{\theta}}(\mathbf{H}) := \widetilde{\mathbf{W}_{\mathbf{0}}} \times \mathrm{MLP}_{\boldsymbol{\theta}_{\mathrm{mlp}}^{(\ell)}} \left( \mathrm{Attn}_{\boldsymbol{\theta}_{\mathrm{attn}}^{(\ell)}}(\mathbf{H}^{(\ell-1)}) \right) \times \widetilde{\mathbf{W}_{\mathbf{1}}}, \quad \ell \in \{1, \dots, L\}$$

The two additional matrices  $\widetilde{\mathbf{W}_0} \in \mathbb{R}^{d_1 \times D}$  and  $\widetilde{\mathbf{W}_1} \in \mathbb{R}^{N \times d_2}$  serve for the dimension adjustment purpose such that the output of  $TF_{\boldsymbol{\theta}}()$  will be of dimension  $\mathbb{R}^{d_1 \times d_2}$ .

#### We additionally define the following norm of a Transformer $\mathrm{TF}_{\theta}$ :

$$\begin{split} \|\boldsymbol{\theta}\|_{op} &:= \\ \max_{\ell \in [L]} \left\{ \max_{m \in [M]} \left\{ \|\mathbf{Q}_m^{(\ell)}\|_{op}, \|\mathbf{K}_m^{(\ell)}\|_{op} \right\} + \sum_{m=1}^M \|\mathbf{V}_m^{(\ell)}\|_{op} + \|\mathbf{W}_1^{(\ell)}\|_{op} + \|\mathbf{W}_2^{(\ell)}\|_{op} \right\} \end{split}$$

We can prove that Transformer is Lipschitz continous to this norm when inputs are bounded.

#### Input:

$$\mathbf{H} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_N & \mathbf{x}_{N+1} \\ y_1 & y_2 & \cdots & y_N & 0 \\ \mathbf{p}_1 & \mathbf{p}_2 & \cdots & \mathbf{p}_N & \mathbf{p}_{N+1} \end{bmatrix} \in \mathbb{R}^{D \times (N+1)}, \quad \mathbf{p}_i := \begin{bmatrix} \mathbf{0}_{D-(d+3)} \\ 1 \\ 1\{i < N+1\} \end{bmatrix} \in \mathbb{R}^{D \times (N+1)}$$

 $\{\mathbf{p}_i\}$  are fixed vectors consisting of ones, zeros, and indicator for being the trained token (similar to a positional encoding vector). We assume

 $\|\mathbf{x}_i\|_2 \le B_x, |y_i| \le B_y, a.s.$ 

#### Input:

$$\mathbf{H} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_N & \mathbf{x}_{N+1} \\ y_1 & y_2 & \cdots & y_N & 0 \\ \mathbf{p}_1 & \mathbf{p}_2 & \cdots & \mathbf{p}_N & \mathbf{p}_{N+1} \end{bmatrix} \in \mathbb{R}^{D \times (N+1)}, \quad \mathbf{p}_i := \begin{bmatrix} \mathbf{0}_{D-(d+3)} \\ 1 \\ 1\{i < N+1\} \end{bmatrix} \in \mathbb{R}^{D \times (N+1)}$$

 $\{\mathbf{p}_i\}$  are fixed vectors consisting of ones, zeros, and indicator for being the trained token (similar to a positional encoding vector). We assume

$$\|\mathbf{x}_i\|_2 \le B_x, |y_i| \le B_y, a.s.$$

Output:

$$\widetilde{\mathbf{H}} = \mathsf{TF}_{\theta}(\boldsymbol{H})$$

$$\hat{y}_{N+1} = \widetilde{\operatorname{read}}_y(\widetilde{\mathbf{H}}) := \operatorname{clip}_R\left(\left(\widetilde{\mathbf{h}}_{N+1}\right)_{d+1}\right)$$

3 🕨 🤅 3

Decoder TFs are the same as encoder TFs, except that the attention layers are replaced by masked attention layers with a specific decoder-based (causal) attention mask.

#### Definition 4

Masked Attention Layer A masked attention layer with M heads is denoted as  $\operatorname{MAttn}_{\boldsymbol{\theta}}(\cdot)$  with parameters  $\boldsymbol{\theta} = \left\{ (\boldsymbol{V}_m, \boldsymbol{Q}_m, \boldsymbol{K}_m) \in (\mathbb{R}^{D \times D})^3 \right\}_{m=1}^M$ . On any input sequence  $\boldsymbol{H} \in \mathbb{R}^{D \times N'}$  with  $N' \leq N$ ,

$$\begin{split} &\widetilde{\boldsymbol{H}} = \mathrm{MAttn}_{\boldsymbol{\theta}}(\boldsymbol{H}) := \\ &\boldsymbol{H} + \sum_{m=1}^{M} (\boldsymbol{V_m} \boldsymbol{H}) \times \left( (MSK_{1:N',1:N'}) \circ \sigma(\boldsymbol{Q_m} \boldsymbol{H})^\top (\boldsymbol{K_m} \boldsymbol{H}) \right) \in \mathbb{R}^{D \times N'}, \end{split}$$

9/60

< 4<sup>™</sup> > <

In the definition,  $\circ$  denotes the entry-wise (Hadamard) product of two matrices, and  $MSK \in \mathbb{R}^{N \times N}$  is the mask matrix given by

$$MSK = \begin{bmatrix} 1 & 1/2 & 1/3 & \cdots & 1/N \\ 0 & 1/2 & 1/3 & \cdots & 1/N \\ 0 & 0 & 1/3 & \cdots & 1/N \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & \cdots & 1/N \end{bmatrix}.$$

In vector form, we have

$$\widetilde{\boldsymbol{h}}_i = [\operatorname{Attn}_{\boldsymbol{\theta}}(\boldsymbol{H})]_i = \boldsymbol{h}_i + \sum_{m=1}^M \frac{1}{i} \sum_{j=1}^i \sigma(\langle \boldsymbol{Q}_m \boldsymbol{h}_i, \boldsymbol{K}_m \boldsymbol{h}_j \rangle) \cdot \boldsymbol{V}_m \boldsymbol{h}_j.$$

(Decoder-based Transformer). An *L*-layer decoder-based Transformer, denoted as  $DTF_{\theta}(\cdot)$ , is a composition of *L* self-attention layers each followed by an MLP layer:

 $\mathbf{H}^{(L)} = \mathrm{DTF}_{\boldsymbol{\theta}}(\mathbf{H}^{(0)}), \text{ where } \mathbf{H} = \mathbf{H}^{(0)} \in \mathbb{R}^{D \times N}$ 

is the input sequence, and

$$\mathbf{H}^{(\ell)} = \mathrm{MLP}_{\boldsymbol{\theta}_{\mathrm{mlp}}^{(\ell)}} \left( \mathrm{MAttn}_{\boldsymbol{\theta}_{\mathrm{mattn}}^{(\ell)}} (\mathbf{H}^{(\ell-1)}) \right), \quad \ell \in \{1, \dots, L\}.$$

Above, the parameter  $\boldsymbol{\theta} = (\boldsymbol{\theta}_{\text{mattn}}^{(1:L)}, \boldsymbol{\theta}_{\text{mlp}}^{(1:L)})$  consists of the attention layers  $\boldsymbol{\theta}_{\text{mattn}}^{(\ell)} = \{(\mathbf{V}_m^{(\ell)}, \mathbf{Q}_m^{(\ell)}, \mathbf{K}_m^{(\ell)})\}_{m \in [M]} \subset \mathbb{R}^{D \times D}$  and the MLP layers  $\boldsymbol{\theta}_{\text{mlp}}^{(\ell)} = (\mathbf{W}_1^{(\ell)}, \mathbf{W}_2^{(\ell)}) \in \mathbb{R}^{D \times D'} \times \mathbb{R}^{D' \times D}.$ 

### Input Matrix

Input:

$$\mathbf{H} = \begin{bmatrix} \boldsymbol{x}_{1} & \mathbf{0} & \dots & \boldsymbol{x}_{N} & \mathbf{0} & \boldsymbol{x}_{N+1} \\ 0 & y_{1} & \dots & 0 & y_{N} & 0 \\ \boldsymbol{p}_{1} & \boldsymbol{p}_{2} & \dots & \boldsymbol{p}_{2N-1} & \boldsymbol{p}_{2N} & \boldsymbol{p}_{2N+1} \end{bmatrix} \in \mathbb{R}^{D \times (2N+1)},$$
$$\boldsymbol{p}_{i} := \begin{bmatrix} \mathbf{0}_{D-(d+4)} \\ \lceil i/2 \rceil \\ 1 \\ \text{mod}(i+1,2) \end{bmatrix} \in \mathbb{R}^{D-(d+1)}$$

July 6, 2025

< 3 >

æ

### Input Matrix

Input:

$$\mathbf{H} = \begin{bmatrix} \boldsymbol{x}_{1} & \mathbf{0} & \dots & \boldsymbol{x}_{N} & \mathbf{0} & \boldsymbol{x}_{N+1} \\ 0 & y_{1} & \dots & 0 & y_{N} & 0 \\ \boldsymbol{p}_{1} & \boldsymbol{p}_{2} & \dots & \boldsymbol{p}_{2N-1} & \boldsymbol{p}_{2N} & \boldsymbol{p}_{2N+1} \end{bmatrix} \in \mathbb{R}^{D \times (2N+1)},$$
$$\boldsymbol{p}_{i} := \begin{bmatrix} \mathbf{0}_{D-(d+4)} \\ \lceil i/2 \rceil \\ 1 \\ \text{mod}(i+1,2) \end{bmatrix} \in \mathbb{R}^{D-(d+1)}$$

Output:

 $\widetilde{\mathbf{H}} = \mathsf{DTF}_{\theta}(\boldsymbol{\mathit{H}})$ 

$$\hat{y}_{N+1} = \widetilde{\operatorname{read}}_y(\widetilde{\mathbf{H}}) := \operatorname{clip}_R\left(\left(\widetilde{\mathbf{h}}_{2N+1}\right)_{d+1}\right)$$

æ

- The input format for decoder-based Transformer is different from the input format for encoder-based Transformers.
- The main difference is that  $(x_i, y_i)$  are in different tokens in the former, whereas  $(i, y_i)$  are in the same token in the latter.
- The reason for the former (i.e., different tokens in decoder) is that we want to avoid every  $[x_i; 0]$  token seeing the information of  $y_i$ , since we will evaluate the loss at every token.
- The reason for the latter (i.e., the same token in encoder) is for presentation convenience: since we only evaluate the loss at the last token, it is not necessary to alternate between  $[x_i; 0]$  and  $[0; y_i]$  to avoid information leakage.

#### (Proof to be added)

#### Proposition 6

Input format conversion There exists a 2-layer decoder TF with 3 heads per layer, hidden dimension 2 and  $\|\boldsymbol{\theta}\|_2 \leq 12$  such that upon taking input  $\boldsymbol{H}$  of format for decoder-based Transformers, it outputs  $\widetilde{\boldsymbol{H}} = \text{DTF}(\boldsymbol{H})$  with

$$\widetilde{\mathbf{H}} = egin{bmatrix} m{x}_1 & m{0} & \dots & m{x}_N & m{0} & m{x}_{N+1} \ 0 & y_1 & \dots & 0 & y_N & 0 \ m{p}_1 & m{p}_2 & \dots & m{p}_{2N-1} & m{p}_{2N} & m{p}_{2N+1} \end{bmatrix} \in \mathbb{R}^{D imes (2N+1)},$$

In particular, this format contains the format for encoders as a submatrix, by restricting to the  $\{1, 2, ..., D - 1, D - 2, D\}$  rows and  $\{2, 4, ..., 2N - 2, 2N, 2N + 1\}$  columns.

#### Transformers

- Encoder-based Transformers
- Decoder-based Transformers

### 2 How Transformers Perform PCA

How Transformers Perform Sparse Recovery: A brief example of decoder-based Transformers

Algorithm 1: Power Method for the Left Singular Vectors

**Data:** Matrix  $X \in \mathbb{R}^{d \times N}$ , Number of Iterations  $\tau$ Symmertize  $A = XX^{\top} \in \mathbb{R}^{d \times d}$ ; Let the set of eigenvectors be  $\mathcal{V} = \{\}$ . Initialize  $A_1 \leftarrow A$ ; for  $\ell \leftarrow 1$  to k do Sample a random vector  $\mathbf{v}_{0,\ell} \in \mathbb{S}^{N-1}$ . Initialize  $\mathbf{v}_{\ell}^{(0)} \leftarrow \mathbf{v}_{0,\ell}$ ; for  $t \leftarrow 1$  to  $\tau$  do Apply the procedure to obtain the principle eigenvector  $\mathbf{v}_{\ell}^{(t)} = \frac{A_{\ell}\mathbf{v}_{\ell}^{(t-1)}}{\|A_{\ell}\mathbf{v}_{\ell}^{(t-1)}\|_{2}}$ ; Let  $\mathcal{V} \leftarrow \mathcal{V} \cup \{\mathbf{v}_{\ell}^{(\tau)}\}$ ; Compute the eigenvalue estimate  $\hat{\lambda}_{\ell} \leftarrow \|A_{\ell}\mathbf{v}_{\ell}^{(\tau)}\|_{2}$ ; Update the matrix by  $A_{\ell+1} = A_{\ell} - \hat{\lambda}_{\ell}\mathbf{v}_{\ell}^{(\tau)}\mathbf{v}_{\ell}^{(\tau),\top}$ ; return  $\mathcal{V}$ ;

э.



æ

### Pretraining via Supervised Learning

We construct the input of the Transformer as a *context-augmented matrix* given by the following:

$$\mathbf{H} = \begin{bmatrix} \mathbf{X} \\ \mathbf{P} \end{bmatrix} \in \mathbb{R}^{D \times N}, \quad \mathbf{P} = \begin{bmatrix} \widetilde{p}_{1,1}, \dots, \widetilde{p}_{1,N} \\ \widetilde{p}_{2,1}, \dots, \widetilde{p}_{2,N} \\ \vdots \\ \widetilde{p}_{D-d,1}, \dots, \widetilde{p}_{D-d,N} \end{bmatrix} \in \mathbb{R}^{(D-d) \times N},$$

The auxillary matrix  $\mathbf{P}$  contains contextual information; the design also maked sure P is unrelated to X. The experiments show that P is not necessary for the pre-trained Transformer to performer PCA with high accuracy.

We construct the input of the Transformer as a *context-augmented matrix* given by the following:

$$\mathbf{H} = \begin{bmatrix} \mathbf{X} \\ \mathbf{P} \end{bmatrix} \in \mathbb{R}^{D \times N}, \quad \mathbf{P} = \begin{bmatrix} \widetilde{p}_{1,1}, \dots, \widetilde{p}_{1,N} \\ \widetilde{p}_{2,1}, \dots, \widetilde{p}_{2,N} \\ \vdots \\ \widetilde{p}_{D-d,1}, \dots, \widetilde{p}_{D-d,N} \end{bmatrix} \in \mathbb{R}^{(D-d) \times N},$$

The auxillary matrix  $\mathbf{P}$  contains contextual information; the design also maked sure P is unrelated to X. The experiments show that P is not necessary for the pre-trained Transformer to performer PCA with high accuracy. **Output:** 

$$TF_{\theta}(\mathbf{H}) = \begin{bmatrix} \hat{v}_1^{\top} & \dots & \hat{v}_k^{\top} \end{bmatrix}^{\top} \in \mathbb{R}^{d \times k}$$

which corresponds to the estimated principal eigenvectors of the matrix  ${f X}.$ 

Consider a set of samples  $\{\mathbf{X}^{(i)}\}_{i\in[n]}$  i.i.d. sampled from some distribution  $p_{\mathbf{X}}$ , we construct their oracle top-k principal components as  $\mathbf{V}^{(i)} = \begin{bmatrix} \mathbf{v}_1^{i,\top} & \dots & \mathbf{v}_k^{i,\top} \end{bmatrix}^{\top}$  and the context-augmented input matrix as  $\mathbf{H}_i$  for each  $\mathbf{X}_i$ . Then, the pretraining procedure is given by minimizing the following objective for some convex loss function  $L(\cdot, \cdot) : \mathbb{R}^{d \times k} \times \mathbb{R}^{d \times k} \to \mathbb{R}$ ,

$$\hat{\theta} = \operatorname*{arg\,min}_{\theta \in \Theta(B_{\theta}, B_M)} \sum_{i=1}^n L(TF_{\theta}(\mathbf{H}^{(i)}), \mathbf{V}^{(i)}).$$
(1)

Here we consider  $\Theta(B_{\theta}) := \{\theta : \|\theta\| \le B_{\theta}, \max_{\ell} M^{\ell} \le B_M\}$  to be the space of parameters. We also consider guarantees when  $L(\mathbf{x}_1, \mathbf{x}_2) := \|\mathbf{x}_1 - \mathbf{x}_2\|_2$  in the theory.

July 6, 2025

Our design of the matrix  ${\bf P}$  consists of three parts:

- Place Holder. For  $\ell \in \{1\} \cup [4:k+3]$  and  $i \in [N]$ , we let  $\widetilde{\mathbf{p}}_{\ell,i} = \mathbf{0} \in \mathbb{R}^{d \times 1}$ . The placeholders in  $\mathbf{P}$  record the intermediate results in the forward propagation. Recall that k is the number of eigenvectors we hope to recover.
- **2** Identity Matrix. We let  $[\widetilde{\mathbf{p}}_{2,1} \ldots \widetilde{\mathbf{p}}_{2,N}] = [\mathbf{I}_d \quad \mathbf{0}_{d \times (N-d)}]$ . The identity matrix in  $\mathbf{P}$  helps us screen out all the covariates  $\mathbf{X}$  in the forward propagation.
- Random Samples on the Hypersphere. We let  $\tilde{\mathbf{p}}_{3,1}, \ldots \tilde{\mathbf{p}}_{3,k}$  be the i.i.d. samples uniformly distributed on  $\mathbb{S}^{d-1}$ . The random samples on the sphere correspond to the initial vectors  $\mathbf{v}_{0,\ell}$  for  $\ell \in [k]$  in algorithm 1.

The auxiliary matrix is designed for purely technical reasons. Moreover, our experiments suggest that such an auxiliary matrix is not necessary for the task.

- ロ ト ・ 同 ト ・ 三 ト ・ 三 ト - -

#### Theorem 7

Assume that the eigenvalues of  $\mathbf{X}\mathbf{X}^{\top}$  to be  $\lambda_1 > \lambda_2 > \ldots > \lambda_k > \ldots$ . Let  $\Delta := \min_{1 \le i < j \le k} |\lambda_i - \lambda_j|$ . Assume that the initialized vectors  $\widetilde{\mathbf{p}}_{3,1}, \ldots \widetilde{\mathbf{p}}_{3,N}$  satisfy  $\widetilde{\mathbf{p}}_{3,i}^{\top}\mathbf{v}_i \ge \delta$  for all  $i \in [k]$  and make the rest of the vectors  $\mathbf{0}$ . Then, there exists a Transformer model with the number of layers  $L = 2\tau + 4k + 1$  and the number of heads  $B_M \le \lambda_1^d \frac{C}{\epsilon^2}$  with  $\tau \le \frac{\log(1/\epsilon_0 \delta)}{\epsilon_0}$  such that for all  $\epsilon_0, \epsilon > 0$ , the final output  $[v_1, \ldots, v_k]$  given by the Transformer model achieves

$$\|v_{\eta+1} - \mathbf{v}_{\eta+1}\|_2 \le C\tau\epsilon\lambda_1^2 + \frac{C\lambda_1\sqrt{\epsilon_0}}{\Delta}\prod_{i=1}^k \frac{5\lambda_{i+1}}{\Delta}$$

Moreover, consider the accuracy of multiple vs as a whole. There exists  $\theta \in \Theta(B_{\lambda_1}, B_M)$  that satisfies

$$L(TF_{\theta}(\mathbf{H}), \mathbf{V}) \leq C\tau\epsilon k\lambda_1^2 + C\frac{\epsilon_0\lambda_1^2}{\Delta^2}\sum_{\eta=1}^{k-1}\prod_{i=1}^{\eta}\frac{25\lambda_{i+1}^2}{\Delta^2}^{1/2}$$

- The first error term comes from the approximation of the Power Method iterations by transformers.
- The second error term comes from finite iteration.

The eigenvalues  $\lambda_1 \asymp \lambda_2 \asymp \ldots \asymp \lambda_k \asymp \Delta$ . Then our results boil down to

$$|TF_{\theta}(\mathbf{H}) - [\mathbf{v}_1^{\top}, \mathbf{v}_2^{\top}, \dots, \mathbf{v}_k^{\top}]^{\top}|_2 \lesssim \tau \epsilon k \lambda_1^2 + \frac{\lambda_1}{\Delta} \sqrt{k\epsilon_0}.$$

These results hide dimension d in the universal constant. We note that the dimension significantly affects the approximation bound of Transformers. This is mainly due to the limitations given by approximating high dimensional functions by ReLU neural networks.

In the above theorem, our results rely on the random initialization of  $\mathbf{P}$ . We show that the conditions on  $\widetilde{\mathbf{p}}_{3,1}, \ldots, \widetilde{\mathbf{p}}_{3,N}$  can be achieved through sampling from isotropic Gaussians, given by the following lemma.

#### Lemma 8

Consider  $\mathbf{x}$  to be a random vector sampled uniformly at random on  $\mathbb{S}^{d-1}$ . Let  $\mathbf{v}$  be any unit length vector, then we have for all  $\delta < \frac{1}{2}d^{-1}$ ,  $\mathbb{P}(|\mathbf{v}^{\top}\mathbf{x}| \leq \delta) \leq \frac{1}{\sqrt{\pi}}\sqrt{\delta} + \exp(-C\delta^{-\frac{1}{2}})$ . Therefore, for all  $\delta < \frac{1}{2}d^{-1}$ , the event in theorem 7 is achieved with

$$\mathbb{P}\exists i \in [k] \text{ such that } \mathbf{x}_i^\top \mathbf{v}_i \leq \frac{\delta}{\sqrt{d}} \leq \frac{k\sqrt{\delta}}{\sqrt{\pi}} + k\exp(-C\delta^{-1}).$$

Given the approximation error provided by theorem 7, we further provide the generalization error bound for the ERM defined by  $\hat{\theta} = \arg \min_{\theta \in \Theta(B_{\theta}, B_M)} \sum_{i=1}^{n} L(TF_{\theta}(\mathbf{H}^{(i)}), \mathbf{V}^{(i)})$ . This requires us to consider the following regularity conditions on the underlying distribution of  $\mathbf{X}\mathbf{X}^{\top}$  (which also translates to the distribution for the pre-training instances  $\mathbf{X}$ ).

・ロト ・ 四ト ・ ヨト ・ ヨト … ヨ

The distribution of  $\mathbf{X}\mathbf{X}^{\top}$  supports on the following space

$$\mathbb{X} := \{ \mathbf{A} : \mathbf{A} \in \mathbf{S}_{++}^d, B_X \ge \lambda_1(\mathbf{A}) > \lambda_2(\mathbf{A}) > \dots > \lambda_k(\mathbf{A}), \\ \inf_{1 \le i < j \le k} \lambda_i(\mathbf{A}) - \lambda_j(\mathbf{A}) \ge \Delta \}.$$

The above assumption can be generalized to a distribution that supports X with high probability. Examples of such distribution include the Wishart distribution under the Gaussian design.

Given the above assumption, we are ready to state the generalization bound.

#### Proposition 9

Under assumption 1 and using the notations given by theorem 7, with probability at least  $1 - \xi$ , the ERM solution  $\theta$  satisfies

$$\mathbb{E}[L(TF_{\theta}(\mathbf{H}), \mathbf{V}) | \theta] \leq \inf_{\theta \in \Theta(B_{\theta}, B_M)} \mathbb{E}[L(TF_{\theta}(\mathbf{H}), \mathbf{V})] + C\sqrt{\frac{k^3 L B_M d^2 \log(B_X + k) + \log(1/\xi)}{n}},$$

where the expectation is taken over the new sample  $\mathbf{X}$ .

< 4<sup>™</sup> > <

### Corollary 10

Under assumption 1, with probability at least  $1 - \xi - \frac{k\sqrt{\delta}}{\sqrt{\pi}} - k \exp(-C\delta^{-1/2}) \text{ for all } \delta < d^{-1} \text{ we have for all } \epsilon, \epsilon_0 > 0,$   $L_{PCA}(\theta, \mathbf{P}) := \mathbb{E}[L(TF_{\theta}(\mathbf{H}), \mathbf{V}) \Big| \theta, \mathbf{P}] \lesssim \tau \epsilon k \lambda_1^2 + \frac{\epsilon_0 \lambda_1^2}{\Delta^2} \sum_{\eta=1}^{k-1} \prod_{i=1}^{\eta} \frac{25\lambda_{i+1}^2}{\Delta^2}^{1/2} + \sqrt{\frac{k^3 \log(\delta/\epsilon_0) B_X^d d^2 \log(B_X + k) + \log(1/\xi)}{n\epsilon_0 \epsilon^2}}.$ 

If we consider optimizing the bound w.r.t.  $\epsilon_0$  and  $\epsilon$ , we obtain that  $L_{PCA}(\theta,\mathbf{P}) \lesssim n^{-1/5}$  with high probability, given that the parameters and the dimension d are of constant scales.

July 6, 2025

・ロト ・ 母 ト ・ ヨ ト ・ ヨ ト …… ヨ

Our proof can be disected into the following setps: We construct a Transformer with fixed parameters that performs:

- Interpretation of the symmetrized covariate matrix;
- One approximation of the power method;
- The removal of the principal eigenvectors;
- Adjust the dimension of the output through multiplying the two matrices  $\widetilde{\mathbf{W}}_0$  and  $\widetilde{\mathbf{W}}_1$  on the left and right.

### 1. The Covariate Matrix

Construct 
$$\mathbf{H} = \begin{bmatrix} \mathbf{X}_1, \dots, \mathbf{X}_N \\ \widetilde{\mathbf{p}}_{1,1}, \dots, \widetilde{\mathbf{p}}_{1,N} \\ \widetilde{\mathbf{p}}_{2,1}, \dots, \widetilde{\mathbf{p}}_{2,N} \\ \vdots \\ \widetilde{\mathbf{p}}_{\ell,1}, \dots, \widetilde{\mathbf{p}}_{\ell,N} \end{bmatrix} = \begin{bmatrix} \mathbf{X} \\ \mathbf{P} \end{bmatrix}$$
, we let the number of heads

m = 2, and construct the first covariate layer as follows,

$$\mathbf{V}_{1}^{cov} = I_{D} = -\mathbf{V}_{2}^{cov}, \quad \mathbf{Q}_{1}^{cov,\top}\mathbf{K}_{1}^{cov} = -\mathbf{Q}_{2}^{\top}\mathbf{K}_{2} = \begin{bmatrix} \mathbf{0}_{N+1\times d} & I_{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{D\times}$$
(2)

$$\widetilde{\mathbf{p}}_{1,\ell,j} = \mathbf{0}, \qquad \widetilde{\mathbf{p}}_{2,\ell,j} = \begin{cases} 1_{\ell=j} & \text{when } \ell \leq d \\ 0 & \text{when } \ell > d \end{cases}.$$
(3)

æ

28 / 60

∃ ► < ∃ ►</p>

### 1. The Covariate Matrix

Under the above constructions, we obtain that

$$\begin{aligned} \mathbf{Q}_1^{\top} \mathbf{K}_1 \mathbf{H} &= \begin{bmatrix} I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{D \times N}, \quad \mathbf{Q}_2^{\top} \mathbf{K}_2 \mathbf{H} = \begin{bmatrix} -I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{D \times N}, \\ \sigma(\mathbf{H}^{\top} \mathbf{Q}_1^{\top} \mathbf{K}_1 \mathbf{H}) + \sigma(\mathbf{H}^{\top} \mathbf{Q}_2^{\top} \mathbf{K}_2 \mathbf{H}) &= \begin{bmatrix} \mathbf{X}^{\top}, \mathbf{0} \end{bmatrix} \in \mathbb{R}^{N \times N}. \end{aligned}$$

We further obtain that

$$\frac{1}{N}\sum_{m=1}^{M} (\mathbf{V}_m \mathbf{H}) \times \sigma((\mathbf{Q}_m \mathbf{H})^{\top}(\mathbf{K}_m \mathbf{H})) = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{X} \mathbf{X}^{\top} \in \mathbb{R}^{d \times d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{D \times D}$$

Г

Therefore, the output is given by  $\widetilde{\mathbf{H}}^{cov} = \Big|_{\widetilde{\mathbf{H}}}$ 

$$\begin{bmatrix} \mathbf{X} \\ \mathbf{X} \mathbf{X}^\top, \mathbf{0} \\ \widetilde{\mathbf{p}}_{2,1}, \dots, \widetilde{\mathbf{p}}_{2,N} \\ \widetilde{\mathbf{p}}_{\ell,1}, \dots, \widetilde{\mathbf{p}}_{\ell,N} \end{bmatrix}$$

V

п

Step 1: Obtaining the vector given by  $\mathbf{X}\mathbf{X}^{\top}\mathbf{v}$ .

Step 2: Approximation of the value of the inverse norm given by  $1/\|\mathbf{X}\mathbf{X}^{\top}\mathbf{v}\|_{2}$ . We show that one can use the multihead ReLU Transformer to achieve both goals simulatenously, whose parameters are given by

$$\begin{split} \mathbf{V}_{1}^{pow,1} &= -\mathbf{V}_{2}^{pow,1} = \begin{bmatrix} \mathbf{0}_{(3d+1)\times(2d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{(d)\times(2d+1)} & I_{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{Q}_{1}^{pow,1} &= -\mathbf{Q}_{1}^{pow,1} = \begin{bmatrix} \mathbf{0}_{(d+1)\times(d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d\times(d+1)} & I_{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{K}_{1}^{pow,1} &= \mathbf{K}_{2}^{pow,1} = \begin{bmatrix} \mathbf{0}_{(3d+1)\times(3d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d\times(3d+1)} & I_{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \qquad \widetilde{\mathbf{p}}_{4,j} = \mathbf{0} \text{ for all } j \in [N]. \end{split}$$

We can calculate that the output of the first power iteration layer is given by

$$\widetilde{\mathbf{H}}^{pow,1} = \begin{bmatrix} \mathbf{X} \\ \widetilde{\mathbf{y}}^{\top} \\ \mathbf{X}\mathbf{X}^{\top}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{2,1}, \dots, \widetilde{\mathbf{p}}_{2,N} \\ \widetilde{\mathbf{p}}_{3,1}, \dots, \widetilde{\mathbf{p}}_{3,N} \\ \mathbf{X}\mathbf{X}^{\top} \widetilde{\mathbf{p}}_{3,1}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{5,1}, \dots, \widetilde{\mathbf{p}}_{5,N} \\ \vdots \\ \widetilde{\mathbf{p}}_{\ell,1}, \dots, \widetilde{\mathbf{p}}_{\ell,N} \end{bmatrix}.$$

July 6, 2025

э

Then, using lemma, we design an extra attention layer that performs the normalizing procedure, with the following parameters for all  $m \in [M]$ ,

$$\begin{split} \mathbf{V}_{m}^{pow,2} &= \begin{bmatrix} \mathbf{0}_{d \times (4d+1)} & c_{m}\mathbf{I}_{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \qquad \mathbf{Q}_{m}^{pow,2} = \begin{bmatrix} \mathbf{0}_{d \times (2d+1)} & \mathbf{I}_{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{K}_{m}^{pow,2} &= \begin{bmatrix} \mathbf{0}_{1 \times (3d+1)} & \mathbf{a}_{m}^{\top} & \mathbf{0} \\ \vdots & & \\ \mathbf{0}_{1 \times (3d+1)} & \mathbf{a}_{m}^{\top} & \mathbf{0} \\ \mathbf{0}_{(D-d) \times (3d+1)} & \mathbf{0} & \mathbf{0} \end{bmatrix}. \end{split}$$

July 6, 2025

Then, given  $\mathbf{V}_m^{pow,2}$  we can show that under the condition given by lemma, we have

$$\begin{aligned} & \left\| \sum_{m=1}^{M} \mathbf{V}_{m}^{pow,2} \widetilde{\mathbf{H}}^{pow,1} \sigma((\mathbf{Q}_{m}^{pow,2} \widetilde{\mathbf{H}}^{pow,1})^{\top} (\mathbf{K}_{m}^{pow,2} \widetilde{\mathbf{H}}^{pow,1})) - \left| \begin{array}{c} \mathbf{0}_{4d+1} \\ \mathbf{X} \mathbf{X}^{\top} \widetilde{\mathbf{p}}_{3,1} \\ \frac{\|\mathbf{X} \mathbf{X}^{\top} \widetilde{\mathbf{p}}_{3,1}\|_{2}}{0} \\ \mathbf{0} \\ \leq \epsilon, \end{aligned} \right. \end{aligned}$$

Moreover, we can further achieve that

$$\left\| \sum_{m=1}^{M} \mathbf{V}_{m}^{pow,2j} \widetilde{\mathbf{H}}^{pow,1} \sigma((\mathbf{Q}_{m}^{pow,2} \widetilde{\mathbf{H}}^{pow,1})^{\top} (\mathbf{K}_{m}^{pow,2} \widetilde{\mathbf{H}}^{pow,1})) - \begin{bmatrix} \mathbf{0}_{4d} \\ \frac{\mathbf{X} \mathbf{X}^{\top} \widetilde{\mathbf{p}}_{3,1}}{\|\mathbf{X} \mathbf{X}^{\top} \widetilde{\mathbf{p}}_{3,1}\|_{2}} - \mathbf{0} \\ \mathbf{0} \end{bmatrix} \right\| = \mathbf{0}^{\mathbf{U}_{4d}} \mathbf{U}_{m}^{\mathbf{U}_{4d}} \mathbf{U}_{m}^{\mathbf{$$

 $< \epsilon \| \mathbf{X} \mathbf{X}^{\top} \widetilde{\mathbf{p}}_{3,1} \|_2.$ 

Hence, using the fact that  $\widetilde{\mathbf{H}}^{pow,2} = \widetilde{\mathbf{H}}^{pow,1} + \sum_{i=1}^{m} \mathbf{V}_{m}^{pow,2} \widetilde{\mathbf{H}}^{pow,1} \sigma((\mathbf{Q}_{m}^{pow,2} \widetilde{\mathbf{H}}^{pow,1})^{\top}(\mathbf{K}_{m}^{pow,2} \widetilde{\mathbf{H}}^{pow,1}))$ , we obtain that

$$\left\| \widetilde{\mathbf{H}}^{pow,2} - \begin{bmatrix} \mathbf{X} \\ \widetilde{\mathbf{y}} \\ \mathbf{X}\mathbf{X}^{\top}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{2,1}, \dots, \widetilde{\mathbf{p}}_{2,N} \\ \widetilde{\mathbf{p}}_{3,1}, \dots, \widetilde{\mathbf{p}}_{3,N} \\ \frac{\mathbf{X}\mathbf{X}^{\top}\widetilde{\mathbf{p}}_{3,1}}{\|\mathbf{X}\mathbf{X}^{\top}\widetilde{\mathbf{p}}_{3,1}\|_{2}}, \dots \mathbf{0} \\ \vdots \end{bmatrix} \right\|_{2} < \epsilon \| \mathbf{X}\mathbf{X}^{\top}\widetilde{\mathbf{p}}_{3,1} \|_{2}.$$

Then we construct another attention layer, which performs similar calculations as that of pow, 1 but switch the rows of  $\tilde{\mathbf{p}}_{3,1}$  with that of  $\frac{\mathbf{X}\mathbf{X}^{\top}\tilde{\mathbf{p}}_{3,1}}{\|\mathbf{X}\mathbf{X}^{\top}\tilde{\mathbf{p}}_{3,1}\|_{2}}$ . Our construction for the third layer is given by

July 6, 2025

34 / 60

$$\begin{split} \mathbf{V}_1^{pow,3} &= -\mathbf{V}_2^{pow,3} = \begin{bmatrix} \mathbf{0}_{(3d+1)\times(2d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d\times(2d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{Q}_1^{pow,3} &= -\mathbf{Q}_2^{pow,3} = \begin{bmatrix} \mathbf{0}_{(3d+1)\times(d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d\times(d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{K}_1^{pow,3} &= \mathbf{K}_2^{pow,3} = \begin{bmatrix} \mathbf{0}_{(4d+1)\times(4d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d\times(4d+1)} & I_d & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \qquad \widetilde{\mathbf{p}}_{4,j} = \mathbf{0} \text{ for all } j \in [N]. \end{split}$$

Consider we are doing in total of  $\tau$  power iterations, we can set for all  $\tau \in \mathbb{N}^*\text{,}$ 

$$\mathbf{V}_{m}^{pow,2\tau+1} = \mathbf{V}_{m}^{pow,3}, \quad \mathbf{Q}_{m}^{pow,2\tau+1} = \mathbf{Q}_{m}^{pow,3}, \quad \mathbf{K}_{m}^{pow,2\tau+1} = \mathbf{K}_{m}^{pow,3}, \\ \mathbf{V}_{m}^{pow,2\tau+2} = \mathbf{V}_{m}^{pow,4}, \quad \mathbf{Q}_{m}^{pow,2\tau+2} = \mathbf{Q}_{m}^{pow,4}, \quad \mathbf{K}_{m}^{pow,2\tau+2} = \mathbf{K}_{m}^{pow,4}.$$
Qinyan Liu How Transformers Perform PCA and Spa July 6, 2025 35/60

Therefore, taking another layer of normalization, we can show that

$$\left\| \widetilde{\mathbf{H}}^{pow,3} - \begin{bmatrix} \mathbf{X} \\ \widetilde{\mathbf{y}} \\ \mathbf{X}\mathbf{X}^{\top}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{2,1}, \dots, \widetilde{\mathbf{p}}_{2,N} \\ \widetilde{\mathbf{p}}_{3,1}, \dots, \widetilde{\mathbf{p}}_{3,N} \\ \frac{(\mathbf{X}\mathbf{X}^{\top})^2 \widetilde{\mathbf{p}}_{3,1}}{\|\mathbf{X}\mathbf{X}^{\top} \widetilde{\mathbf{p}}_{3,1}\|_2^2}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{5,1}, \dots, \widetilde{\mathbf{p}}_{5,N} \\ \vdots \\ \widetilde{\mathbf{p}}_{\ell,1}, \dots, \widetilde{\mathbf{p}}_{\ell,N} \end{bmatrix} \right\|_2 \le 2\epsilon \|\mathbf{X}\mathbf{X}^{\top}\|_2.$$

July 6, 2025

э

Then, using the sublinearity of errors, we can show that for  $au \in \mathbb{N}$ ,

$$\left\| \widetilde{\mathbf{H}}^{pow,2\tau+2} - \begin{bmatrix} \mathbf{X} \\ \widetilde{\mathbf{y}} \\ \mathbf{X}\mathbf{X}^{\top}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{3,1}, \dots, \widetilde{\mathbf{p}}_{3,N} \\ \widetilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{5,1}, \dots, \widetilde{\mathbf{p}}_{5,N} \\ \vdots \\ \widetilde{\mathbf{p}}_{\ell,1}, \dots, \widetilde{\mathbf{p}}_{\ell,N} \end{bmatrix} \right\|_{\infty} \leq \tau \epsilon \| \mathbf{X}\mathbf{X}^{\top} \|_{2}, \quad \widetilde{\mathbf{p}}_{3,1}^{(\tau)} = \frac{\mathbf{X}\mathbf{X}^{\top} \widetilde{\mathbf{p}}_{3,1}^{(\tau-1)}}{\| \mathbf{X}\mathbf{X}^{\top} \widetilde{\mathbf{p}}_{3,1}^{(\tau-1)} \|_{2}}$$

< ≣ > < ≣ >July 6, 2025

< 4<sup>™</sup> > <

37 / 60

3

If we denote  $\mathbf{v}_i$  as the eigenvector corresponds to the i th largest eigenvalue of  $\mathbf{X}\mathbf{X}^{\top}$ . Let the eigenvalues of  $\mathbf{X}\mathbf{X}^{\top}$  be denoted by  $\lambda_1 > \lambda_2 > \cdots > \lambda_n$ . Given  $|\widetilde{\mathbf{p}}_{3,1}^{\top}\mathbf{v}_1| > \delta$  and  $|\sqrt{\lambda_1} - \sqrt{\lambda_2}| = \Omega(1)$ . Theorem 3.11 in blum2020foundations page 53 shows that given  $k = \frac{\log(1/\epsilon_0\delta)}{2\epsilon_0}$  and  $\|\widetilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2 = \|\mathbf{v}_1\|_2 = 1$ , one immediately obtains that  $\widetilde{\mathbf{p}}_{3,1}^{(\tau),\top}\mathbf{v}_1 \ge 1 - \epsilon_0$ ,  $\|\widetilde{\mathbf{p}}_{3,1}^{(\tau)} - \mathbf{v}_1\|_2 = \sqrt{2 - 2\mathbf{v}_1^{\top}\widetilde{\mathbf{p}}_{3,1}^{(\tau)}} = \sqrt{2\epsilon_0}$ .

And we also consider the approximation of the maximum eigenvalue. Note that using  $\|\mathbf{v}_1\|_2=1,$  we have

$$\begin{split} \|\mathbf{X}\mathbf{X}^{\top}\|_{2} &= \|\mathbf{X}\mathbf{X}^{\top}\mathbf{v}_{1}\|_{2} = \|\mathbf{X}\mathbf{X}^{\top}\widetilde{\mathbf{p}}_{3,1}^{(\tau)} + \mathbf{X}\mathbf{X}^{\top}(\mathbf{v}_{1} - \widetilde{\mathbf{p}}_{3,1}^{(\tau)})\|_{2} \\ &\leq \|\mathbf{X}\mathbf{X}^{\top}\widetilde{\mathbf{p}}_{3,1}^{(\tau)}\|_{2} + \|\mathbf{X}\mathbf{X}^{\top}(\mathbf{v}_{1} - \widetilde{\mathbf{p}}_{3,1}^{(\tau)})\|_{2} \\ &\leq \|\mathbf{X}\mathbf{X}^{\top}\widetilde{\mathbf{p}}_{3,1}^{(\tau)}\|_{2} + \|\mathbf{X}\mathbf{X}^{\top}\|_{2}\|\mathbf{v}_{1} - \widetilde{\mathbf{p}}_{3,1}^{(\tau)}\|_{2}. \end{split}$$

Similarly we can also derive that

 $\|\mathbf{X}\mathbf{X}^{\top}\|_2 \geq \|\mathbf{X}\mathbf{X}^{\top}\widetilde{\mathbf{p}}_{3,1}^{(\tau)}\|_2 - \|\mathbf{X}\mathbf{X}^{\top}\|_2\|\mathbf{v}_1 - \widetilde{\mathbf{p}}_{3,1}\|_2. \text{ Then we show that}$ 

38 / 60

Step 1: The computation of the estimated eigenvalue  $\|\mathbf{X}\mathbf{X}^{\top}\widetilde{\mathbf{p}}_{3,1}\|_2$ . Step 2: The construction of the low rank update  $\widetilde{\mathbf{p}}_{3,1}\widetilde{\mathbf{p}}_{3,1}^{\top}$ . For step (1), we consider the following construction:

$$\begin{split} \mathbf{V}_{1}^{rpe,1} &= -\mathbf{V}_{2}^{rpe,1} = \begin{bmatrix} \mathbf{0}_{(3d+1)\times(2d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d\times(2d+1)} & I_{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{Q}_{1}^{rpe,1} &= -\mathbf{Q}_{2}^{rpe,1} = \begin{bmatrix} \mathbf{0}_{(d+1)\times(d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d\times(d+1)} & I_{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{K}_{1}^{rpe,1} &= \mathbf{K}_{2}^{rpe,1} = \begin{bmatrix} \mathbf{0}_{(4d+1)\times(4d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0}_{d\times(4d+1)} & I_{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}. \end{split}$$

July 6, 2025

Note that the above construction is similar to the first layer of the power method. Under this construction, we can show that

Then, we construct the next layer, using the notations in lemma, for  $M \geq \|\mathbf{X}\mathbf{X}^{\top}\|_{2}^{d} \frac{C(d)}{\epsilon^{2}}$  for all  $m \in [M]$  we have

$$\begin{split} \mathbf{V}_{m}^{rpe,2} &= \begin{bmatrix} \mathbf{0}_{d \times (4d+1)} & d_{m} \mathbf{I}_{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \qquad \mathbf{Q}_{m}^{rpe,2} = \begin{bmatrix} \mathbf{0}_{d \times (2d+1)} & \mathbf{I}_{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{K}_{m}^{rpe,2} &= \begin{bmatrix} \mathbf{0}_{1 \times (5d+1)} & \mathbf{b}_{m}^{\top} & \mathbf{0} \\ \vdots & & \\ \mathbf{0}_{1 \times (5d+1)} & \mathbf{b}_{m}^{\top} & \mathbf{0} \\ \mathbf{0}_{(D-d) \times (5d+1)} & \mathbf{0} & \mathbf{0} \end{bmatrix}. \end{split}$$

Given the above construction, we subsequently show that

$$(\mathbf{Q}_m^{rpe,2} \widetilde{\mathbf{H}}^{rpe,1})^\top = \begin{bmatrix} I_{d \times d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \qquad \mathbf{K}_m^{rpe,2} \widetilde{\mathbf{H}}^{rpe,1} = \begin{bmatrix} \mathbf{b}_m^\top \mathbf{X} \mathbf{X}^\top \widetilde{\mathbf{p}}_{3,1}^{(\tau)} & \mathbf{0} \\ \vdots & \vdots \\ \mathbf{b}_m^\top \mathbf{X} \mathbf{X}^\top \widetilde{\mathbf{p}}_{3,1}^{(\tau)} & \mathbf{0} \\ \mathbf{0}_{(D-d) \times 1} & \mathbf{0} \end{bmatrix}.$$

July 6, 2025

41

Hence, given the construction of  $\mathbf{V}_m^{rpe,2}$  , we can show that  $\widetilde{\mathbf{H}}^{rpe,2}$  satisfies

$$\begin{split} \widetilde{\mathbf{H}}^{rpe,2} &= \widetilde{\mathbf{H}}^{rpe,1} + \sum_{m \in [M]} \mathbf{V}_m^{rpe,2} \widetilde{\mathbf{H}}^{rpe,1} \times \sigma((\mathbf{K}_m^{rpe,2} \widetilde{\mathbf{H}}^{rpe,1})^\top (\mathbf{Q}_m^{rpe} \widetilde{\mathbf{H}}^{rpe,1})) \\ &= \mathbf{H}^{rpe,1} + \sum_{m \in [M]} \mathbf{V}_m^{rpe,2} \widetilde{\mathbf{H}}^{rpe,1} \times \sigma((\mathbf{K}_m^{rpe,2} \widetilde{\mathbf{H}}^{rpe,1})^\top (\mathbf{Q}_m^{rpe,2} \mathbf{H}^{rpe,1})) \\ &+ (\widetilde{\mathbf{H}}^{rpe,1} - \mathbf{H}^{rpe,1}) + \sum_{m \in [M]} \mathbf{V}_m^{rpe,2} \widetilde{\mathbf{H}}^{rpe,1} \times \sigma((\mathbf{K}_m^{rpe,2} \widetilde{\mathbf{H}}^{rpe,1})^\top \mathbf{Q}_m^{rpe,2}) \\ &- \sum_{m \in [M]} \mathbf{V}_m^{rpe,2} \widetilde{\mathbf{H}}^{rpe,1} \times \sigma((\mathbf{K}_m^{rpe,2} \widetilde{\mathbf{H}}^{rpe,1})^\top \mathbf{Q}_m^{rpe,2} \widetilde{\mathbf{H}}^{rpe,1}). \end{split}$$

Denote the sum of the first two terms as  $\hat{\mathbf{H}^{rpe,1}}$ .

42 / 60

We note that by lemma we can show that

$$\left\| H^{rpe,1} - \begin{bmatrix} \mathbf{X} \\ \widetilde{\mathbf{y}} \\ \mathbf{X}\mathbf{X}^{\top}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{2,1}, \dots, \widetilde{\mathbf{p}}_{2,N} \\ \widetilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \| \mathbf{X}\mathbf{X}^{\top}\widetilde{\mathbf{p}}_{3,N}^{(\tau)} \|_{2}^{\frac{1}{2}}\widetilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \vdots \\ \widetilde{\mathbf{p}}_{\ell,1}, \dots \widetilde{\mathbf{p}}_{\ell,N} \end{bmatrix} \right\|_{2} \leq C\tau\epsilon \| \mathbf{X}\mathbf{X}^{\top} \|_{2}^{2}.$$

Then the rest of the proof focuses on showing that the rest of the terms are small. Already we have

$$\|\widetilde{\mathbf{H}}^{rpe,1} - \mathbf{H}^{rpe,1}\|_2 \le \tau \epsilon \|\mathbf{X}\mathbf{X}^\top\|_2^2.$$

And for the last term, we can show that

Qinyan Liu

43 / 60

$$\Big\|\sum_{m\in[M]}\mathbf{V}_m^{rpe,2}\widetilde{\mathbf{H}}^{rpe,1}\times\sigma((\mathbf{K}_m^{rpe,2}\widetilde{\mathbf{H}}^{rpe,1})^{\top}(\mathbf{Q}_m^{rpe,2}\widetilde{\mathbf{H}}^{rpe,1}))-$$

$$\sum_{m \in [M]} \mathbf{V}_m^{rpe,2} \mathbf{H}^{rpe,1} \times \sigma((\mathbf{K}_m^{rpe,2} \mathbf{H}^{rpe,1})^\top (\mathbf{Q}_m^{rpe,2} \mathbf{H}^{rpe,1})) \Big\|_2 \le C \tau \epsilon \|\mathbf{X}\mathbf{X}^\top\|_2^2.$$

Collecting the above pieces, we finally show that

$$\left\|\widetilde{\mathbf{H}}^{rpe,2} - \begin{bmatrix} \mathbf{X} \\ \mathbf{X}\mathbf{X}^{\top}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{2,1}, \dots, \widetilde{\mathbf{p}}_{2,N} \\ \widetilde{\mathbf{p}}_{3,1}, \dots, \widetilde{\mathbf{p}}_{3,N} \\ \widetilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \|\mathbf{X}\mathbf{X}^{\top}\widetilde{\mathbf{p}}_{3,1}^{(\tau)}\|_{2}^{\frac{1}{2}}\widetilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \|\mathbf{X}\mathbf{X}^{\top}\widetilde{\mathbf{p}}_{3,1}^{(\tau)}\|_{2}^{\frac{1}{2}}\widetilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{6,1}, \dots, \widetilde{\mathbf{p}}_{6,N} \\ \vdots \\ \widetilde{\mathbf{p}}_{\ell,1}, \dots, \widetilde{\mathbf{p}}_{\ell,N} \end{bmatrix} \right\|_{2} \leq C\tau\epsilon \|\mathbf{X}\mathbf{X}^{\top}\|_{2}^{2}.$$

Qinyan Liu

Then we construct another layer to remove the principal components from the matrix  ${\bf X}{\bf X}^\top$ , given by

$$\begin{split} -\mathbf{V}_{1}^{rpe,3} &= \mathbf{V}_{2}^{rpe,3} = \begin{bmatrix} \mathbf{0}_{(d+1)\times(4d+1)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & I_{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{Q}_{1}^{rpe,3} &= -\mathbf{Q}_{2}^{rpe,3} = \begin{bmatrix} \mathbf{0}_{d\times(4d+1)} & I_{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \\ \mathbf{K}_{1}^{rpe,3} &= \mathbf{K}_{2}^{rpe,3} = \begin{bmatrix} \mathbf{0}_{d\times(4d+1)} & I_{d} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}. \end{split}$$

Therefore, we can further show that

$$\widetilde{\mathbf{H}}^{rpe,3} = \widetilde{\mathbf{H}}^{rpe,2} + \sum_{m=1}^{2} \mathbf{V}_{m}^{rpe,3} \widetilde{\mathbf{H}}^{rpe,2} \times \sigma((\mathbf{Q}_{m}^{rpe,3} \widetilde{\mathbf{H}}^{rpe,2})^{\top} \mathbf{K}_{m}^{rpe,3} \widetilde{\mathbf{H}}^{rpe,2})$$

satisfies

$$\left\| \widetilde{\mathbf{H}}^{rpe,3} - \begin{bmatrix} \mathbf{X} \\ \mathbf{X}\mathbf{X}^{\top} - \| \mathbf{X}\mathbf{X}^{\top} \widetilde{\mathbf{p}}_{3,1}^{(\tau)} \|_{2} \widetilde{\mathbf{p}}_{3,1}^{(\tau)} \widetilde{\mathbf{p}}_{3,1}^{(\tau),\top}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{2,1}, \dots, \widetilde{\mathbf{p}}_{2,N} \\ \widetilde{\mathbf{p}}_{3,1}, \dots, \widetilde{\mathbf{p}}_{3,N} \\ \widetilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{5,1}, \dots, \widetilde{\mathbf{p}}_{5,N} \\ \vdots \\ \widetilde{\mathbf{p}}_{\ell,1}, \dots, \widetilde{\mathbf{p}}_{\ell,N} \end{bmatrix} \right\|_{2} \leq C\tau\epsilon \| \mathbf{X}\mathbf{X}^{\top} \|_{2}^{2}.$$

And then we proceed to recover the rest of the k principal eigenvectors using similar model architecture given by the ones used by the Power Iterations. For the computation over the  $\tau$ -th eigenvector, we denote  $\widetilde{\mathbf{H}}^{pow,\eta,1}$  till  $\widetilde{\mathbf{H}}^{pow,\eta,\tau}$  to be the intermediate states corresponding to the  $\eta$ -th power iteration. We denote  $\widetilde{\mathbf{H}}^{rpe,\eta,\tau_0}$  to be the output of  $\eta$ -th removal of principal eigenvector layers for the  $\tau$ -th eigenvector. Furthermore, we iteratively define

$$\mathbf{A}_1 = \mathbf{X}\mathbf{X}^\top - \|\mathbf{X}\mathbf{X}^\top \widetilde{\mathbf{p}}_{3,i}^{(\tau)}\|_2 \widetilde{\mathbf{p}}_{3,i}^{(\tau)} \widetilde{\mathbf{p}}_{3,i}^{(\tau),\top}, \qquad \mathbf{A}_{i+1} = \mathbf{A}_i - \|\mathbf{A}_i \widetilde{\mathbf{p}}_{3,i}^{(\tau)}\|_2 \widetilde{\mathbf{p}}_{3,i}^{(\tau)} \widetilde{\mathbf{p}}_{3,i}^{(\tau)}$$

Then, applying the subadditivity of the 2-norm, we can show that

47 / 60

$$\left\| \widetilde{\mathbf{H}}^{rpe,4,k} - \begin{bmatrix} \mathbf{X} \\ \mathbf{A}_{k+1}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{2,1}, \dots, \widetilde{\mathbf{p}}_{2,N} \\ \widetilde{\mathbf{p}}_{3,1}^{(\tau)}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{3,2}^{(\tau)}, \mathbf{0} \\ \vdots \\ \widetilde{\mathbf{p}}_{3,k}^{(\tau)}, \mathbf{0} \end{bmatrix} \right\|_{2} \leq C\tau k\epsilon \| \mathbf{X} \mathbf{X}^{\top} \|_{2}^{2}.$$

For simplicity, we denote  $\widetilde{\mathbf{A}} =$ 

$$\begin{bmatrix} \mathbf{X} \\ \mathbf{A}_{k+1}, \mathbf{0} \\ \widetilde{\mathbf{p}}_{2,1}, \dots, \widetilde{\mathbf{p}}_{2,N} \\ \widetilde{\mathbf{p}}_{3,1}, \dots, \widetilde{\mathbf{p}}_{3,N} \end{bmatrix} \text{ and } \widetilde{\mathbf{P}} = \begin{bmatrix} \widetilde{\mathbf{p}}_{3,1} \\ \widetilde{\mathbf{p}}_{3,2} \\ \vdots \\ \widetilde{\mathbf{p}}_{3,k} \\ \widetilde{\mathbf{p}}_{3,k} \end{bmatrix} \text{ from }$$

here.

 $\Gamma_{\alpha}(\tau) \mathsf{T}$ 

# 4. Finishing Up

Our construction gives the following:

$$\widetilde{\mathbf{W}}_0 = \begin{bmatrix} \mathbf{0}, I_{kd} \end{bmatrix}, \qquad \widetilde{\mathbf{W}}_1 = \begin{bmatrix} 1 \\ \mathbf{0}_{N-1} \end{bmatrix}.$$

And we can show that

$$\left\| \widetilde{\mathbf{W}}_{0} \widetilde{\mathbf{H}}^{rpe,4,k} \widetilde{\mathbf{W}}_{1} - \begin{bmatrix} \widetilde{\mathbf{p}}_{3,1}^{(\tau)} \\ \widetilde{\mathbf{p}}_{3,2}^{(\tau)} \\ \vdots \\ \widetilde{\mathbf{p}}_{3,k}^{(\tau)} \end{bmatrix} \right\|_{2} \leq C \tau k \epsilon \| \mathbf{X} \mathbf{X}^{\top} \|_{2}^{2}$$

We further use the result given by lemma, denote  $a_{\eta} := \|\mathbf{v}_{\eta} - \widetilde{\mathbf{p}}_{3,\eta}^{(\tau)}\|_2$ ,  $\widehat{\lambda}_{\eta} = \|\mathbf{A}_{\eta}\widetilde{\mathbf{p}}_{3,\eta}^{(\tau)}\|_2$ , and  $b_{\eta} := |\lambda_{\eta} - \widehat{\lambda}_{\eta}|$  for  $\eta \in [k]$ , we obtain that for all  $\eta \geq 1$ , given the number of iterations  $\tau \geq C \frac{\log(1/\epsilon_0 \delta)}{2\epsilon_0}$  where the constant value C depends on d,

Qinyan Liu

$$a_{\eta+1} \le \frac{\max_{i \in [\eta]} b_i + \sum_{i=1}^{\eta} 2\lambda_i a_i}{\Delta}, \qquad b_{\eta+1} \le \frac{2\lambda_{\eta+1}}{\Delta} \max_{i \in [\eta]} b_{\eta} + \sum_{i=1}^{\eta} 2\lambda_i a_i + \lambda_i a_i$$

Further note that the starting point is given by  $a_1 \leq \sqrt{2\epsilon_0}$ ,  $b_1 \leq \lambda_1 \sqrt{2\epsilon_0}$ . Introducing  $A_\eta = \sum_{i=1}^{\eta} 2\lambda_i a_i$ , we obtain that  $A_{\eta+1} = \sum_{i=1}^{\eta+1} 2\lambda_i a_i = A_\eta + 2\lambda_{\eta+1}a_{\eta+1}$  which alternatively implies that

$$\frac{1}{2\lambda_{\eta+1}}(A_{\eta+1} - A_{\eta}) \le \frac{\max_{i \in [\eta]} b_i + A_{\eta}}{\Delta},$$
$$b_{\eta+1} \le \frac{2\lambda_{\eta+1}}{\Delta} \max_{i \in [\eta]} b_{\eta} + A_{\eta} + \lambda_{\eta+1}\sqrt{2\epsilon_0}.$$

Qinyan Liu

July 6, 2025

### 4. Finishing Up

We use the fact  $\frac{\lambda_{\eta}}{\Delta} > 1$  for all  $\eta \in [k]$  to show the following

$$A_{\eta+1} + \max_{i \in [\eta+1]} b_i \le \frac{5\lambda_{\eta+1}}{\Delta} (A_\eta + \max_{i \in [\eta]} b_i) + \lambda_1 \sqrt{2\epsilon_0}, \qquad A_1 + b_1 = 2\lambda_1 \sqrt{2\epsilon_0}$$

which implies that

$$\begin{aligned} A_{\eta+1} + \max_{i \in [\eta+1]} b_i + \frac{\lambda_1 \sqrt{2\epsilon_0}}{\Delta} &\leq \frac{5\lambda_1}{\Delta} A_\eta + \max_{i \in [\eta]} b_i + \frac{\lambda_1 \sqrt{2\epsilon_0}}{\frac{5\lambda_1}{\Delta} - 1}, \\ A_{\eta+1} + \max_{i \in [\eta+1]} b_i + \frac{\lambda_1 \sqrt{2\epsilon_0}}{\frac{5\lambda_1}{\Delta} - 1} &\leq A_1 + b_1 + \frac{\lambda_1 \sqrt{2\epsilon_0}}{\frac{5\lambda_1}{\Delta} - 1} \prod_{i=1}^{\eta} \frac{5\lambda_{i+1}}{\Delta} \\ &= \lambda_1 \sqrt{2\epsilon_0} 2 + \frac{1}{\frac{5\lambda_1}{\Delta} - 1} \prod_{i=1}^{\eta} \frac{5\lambda_{i+1}}{\Delta}. \end{aligned}$$

< A > <

э

51 / 60

# 4. Finishing Up

Therefore, for  $\eta \leq k$ , we have for all  $\eta \in [k-1]$ ,

$$a_{\eta+1} \leq \frac{1}{\Delta} \lambda_1 \sqrt{2\epsilon_0} 2 + \frac{1}{\frac{5\lambda_1}{\Delta} - 1} \prod_{i=1}^{\eta} \frac{5\lambda_{i+1}}{\Delta} - \frac{\lambda_1 \sqrt{2\epsilon_0}}{\frac{5\lambda_1}{\Delta} - 1},$$
$$b_{\eta+1} \leq \frac{2\lambda_\eta \lambda_1 \sqrt{2\epsilon_0}}{\Delta} 2 + \frac{1}{\frac{5\lambda_1}{\Delta} - 1} \prod_{i=1}^{\eta} \frac{5\lambda_{i+1}}{\Delta} + \lambda_{\eta+1} \sqrt{2\epsilon_0}.$$

Therefore collecting pieces, we conclude that there exists a transformer with number of layers  $2\tau + 4k + 1$  and number of heads  $M \leq \lambda_1^d \frac{C(d)}{\epsilon^2}$  such that the final output  $v_1, \ldots, v_k$  given by the Transformer model satisfy  $\forall \eta \in [k-1]$ ,

$$\|v_{\eta+1} - \mathbf{v}_{\eta+1}\|_2 \le C\tau\epsilon\lambda_1^2 + \frac{1}{\Delta}\lambda_1\sqrt{2\epsilon_0}2 + \frac{1}{\frac{5\lambda_1}{\Delta} - 1}\prod_{i=1}^{\eta}\frac{5\lambda_{i+1}}{\Delta} - \frac{\lambda_1\sqrt{2\epsilon_0}}{\frac{5\lambda_1}{\Delta} - 1}$$

And the rest of the result directly follows.

Qinyan Liu

#### Transformers

- Encoder-based Transformers
- Decoder-based Transformers

#### 2 How Transformers Perform PCA

B How Transformers Perform Sparse Recovery: A brief example of decoder-based Transformers

## Sparse Recovery: Formulated as LASSO Problems

$$\beta^{(k+1)} = \mathcal{S}_{\alpha/L}(\beta^{(k)} - \frac{1}{L}\mathbf{X}^{\top}(\mathbf{X}\beta^{(k)} - \mathbf{y})).$$

Here  $\alpha$  is a coefficient controlling the sparsity penalty. We denote the transpose of the *i*-th row in **X** by  $\mathbf{x}_i$ , i.e.,  $\mathbf{x}_i = [\mathbf{X}^\top]_{:,i}$ . Algorithms:

- Gradient Descent: inefficent
- Learning Iterative Shrinkage Thresholding Algorithm (LISTA): feed-forward neural network
- Classical variants of LISTA, e.g. LISTA-CP The update rule in the k-th iteration:

$$\beta^{(k+1)} = \mathcal{S}_{\theta^{(k)}} \left( \beta^{(k)} - (\mathbf{D}^{(k)})^{\top} (\mathbf{X} \beta^{(k)} - \mathbf{y}) \right), \tag{4}$$

where  $\{\theta^{(k)}, \mathbf{D}^{(k)}\}$  are learnable parameters.

• Transformer: LISTA-VM(LISTA with varying measurements)

In this work, we set the activation function as the element-wise ReLU function. Next, we define the one-layer decoder-based Transformer structure.

#### Definition 11 (Transformer layer)

A one-layer decoder-based Transformer is parameterized by  $\Theta := {\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}, (\mathbf{V}_m, \mathbf{Q}_m, \mathbf{K}_m)_{m \in [M]}}$ , denoted as  $_{\theta}$ . Therefore, give input sequence  $\mathbf{H} \in \mathbb{R}^{d \times N}$ , the output sequence is:

$$TF_{\theta}(\mathbf{H}) = MLP_{\{\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}\}} \Big( Attn_{\{(\mathbf{V}_m, \mathbf{Q}_m, \mathbf{K}_m)\}}(\mathbf{H}) \Big).$$

Given an in-context sparse recovery instance  $\mathcal{I} = (\mathbf{X}, \mathbf{y}, \mathbf{x}_{N+1})$  we embed the instance into an input sequence  $\mathbf{H}^{(1)} \in \mathbb{R}^{(2d+2) \times (2N+1)}$  as follows:

$$\mathbf{H}^{(1)}(\mathcal{I}) = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_1 & \cdots & \mathbf{x}_N & \mathbf{x}_N & \mathbf{x}_{N+1} \\ 0 & y_1 & \cdots & 0 & y_N & 0 \\ \beta_1^{(1)} & \beta_2^{(1)} & \cdots & \beta_{2N-1}^{(1)} & \beta_{2N}^{(1)} & \beta_{2N+1}^{(1)} \\ 1 & 0 & \cdots & 1 & 0 & 1 \end{bmatrix},$$
(5)

where  $\{\beta_i^{(1)}\}_{i\in[2N+1]} \in \mathbb{R}^d$  are implicit parameter vectors initialized as  $\mathbf{0}_d$ , and  $\mathbf{x}_i$  is the *i*-th column of the transposed measurement matrix, i.e,  $[\mathbf{X}^{\top}]_{:,i}$ . We note that a similar embedding structure is adopted in bai2023transformers.

For  $\mathbf{x} \sim P_{\mathbf{x}}$  and  $\beta^* \sim P_{\beta}$ , we assume  $\|\mathbf{x}\| \leq b_{\mathbf{x}}$  and  $\|\beta^*\|_1 \leq b_{\beta}$  almost surely. Besides, we consider the noiseless scenario where  $\epsilon = \mathbf{0}$ .

< 47 ▶ <

э

#### Theorem 12 (Equivalence between ICL and LISTA-VM)

With the Transformer structure described above, under Assumption 2, there exists a set of parameters in the Transformer so that for any  $k \in [1:K]$ ,  $n \in [N]$ , we have

$$\beta_{2n+1}^{(k+1)} = \mathcal{S}_{\theta^{(k)}} \Big( \beta_{2n+1}^{(k)} - \frac{1}{2n+1} \mathbf{M}^{(k)} [\mathbf{X}]_{1:n,:}^{\top} ([\mathbf{X}]_{1:n,:} \beta_{2n+1}^{(k)} - \mathbf{y}_{1:n}) \Big), \quad (6)$$

where  $\mathbf{M}^{(k)} \in \mathbb{R}^{d \times d}$  is embedded in the k-th Transformer layer.

57 / 60

#### Theorem 13 (Convergence of ICL)

Let  $\delta \in (0,1)$ ,  $N_0 = 8(4S-2)^2 \frac{\log d + \log S - \log \delta}{c}$ ,  $\alpha_n = -\log \left(1 - \frac{2}{3}\gamma + \gamma(2S-1)\sqrt{\frac{\log d - \log \delta}{nc}} + \sqrt{\frac{\log S - \log \delta}{nc}}\right)$ , where c is a positive constant and  $\gamma$  is a positive constant satisfies  $\gamma \leq \frac{3}{2}$ . For a K-layer Transformer model with the structure described previously, under Assumption 2, there exists a set of parameters such that for any randomly generated sparse recovery instance and any  $n \in [N_0 : N]$ , with probability at least  $1 - \delta$ , we have

$$\left\|\beta_{2n+1}^{(K+1)} - \beta^*\right\| \le b_{\beta} e^{-\alpha_n K}.$$

#### Theorem 14

Under the same setting, for any  $n \in [N_0 + 1 : N + 1]$ , with probability at least  $1 - n\delta - \delta'$ , we have

$$||y_n - \hat{y}_n|| \le b_{\mathbf{x}}(1 - \frac{2}{3}\gamma)^K + \frac{c_4K}{\sqrt{n}}(1 - \frac{2}{3}\gamma)^{K-1}$$

for a linear read-out function, and with probability at least  $1 - \delta$ , we have  $||y_n - \hat{y}_n|| \le c_5 e^{-\alpha_n K}$  for a query read-out function, where  $c_4$ ,  $c_5$  are constants.

- Bai(2023) "Transformers as statisticians: Provable in-context learning with in-context algorithm selection."
- Vaswani(2017) "Attention is all you need."
- He(2024) "Can transforeners perform PCA?"
- Liu(2025) "On the learn-to-optimize capabilities of Transformer in in-context sparse recovery"