

# LEARNING TO RECALL WITH TRANSFORMERS BEYOND ORTHOGONAL EMBEDDINGS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Modern large language models (LLMs) excel at tasks that require storing and retrieving knowledge, such as factual recall and question answering. Transformers are central to this capability because they can encode information during training and retrieve it at inference. Existing theoretical analyses typically study transformers under idealized assumptions such as infinite data or orthogonal embeddings. In realistic settings, however, models are trained on finite datasets with non-orthogonal (random) embeddings. We address this gap by analyzing a single-layer transformer with random embeddings trained with (empirical) gradient descent on a simple token-retrieval task, where the model must identify an informative token within a length- $L$  sequence and learn a one-to-one mapping from tokens to labels. Our analysis tracks the “early phase” of gradient descent and yields explicit formulas for the model’s storage capacity—revealing a multiplicative dependence between sample size  $N$ , embedding dimension  $d$ , and sequence length  $L$ . We validate these scalings numerically and further complement them with a lower bound for the underlying statistical problem, demonstrating that this multiplicative scaling is intrinsic under non-orthogonal embeddings.

## 1 INTRODUCTION

Large language models (LLMs) routinely answer knowledge questions with little or no external context, indicating that substantial factual information is stored in parameters and can be retrieved by suitable prompts (Petroni et al., 2019; Jiang et al., 2020; Roberts et al., 2020). A sharper theoretical account of how such parametric memories are learned and accessed is increasingly important: it can guide scaling choices (e.g., trading off memory capacity against compute budgets, Carlini et al., 2022; Allen-Zhu & Li, 2024) and illuminate failure modes (e.g., hallucination, Zucchet et al., 2025; Huang et al., 2025). Motivated by empirical results documenting the prevalence of parametric factual recall and its scaling with model size (Allen-Zhu & Li, 2024; Morris et al., 2025), recent theoretical works have begun to analyze the capacity and learning dynamics of transformers on controlled factual-recall tasks (Cabannes et al., 2024a; Nichani et al., 2025).

Many theoretical studies of transformer optimization work in population-dynamics settings and adopt simplifying assumptions such as treating token embeddings as orthogonal or one-hot vectors (see, e.g., Tian et al. 2023b; Chen et al. 2024; Ghosal et al. 2024). These choices do not always reflect practical applications, but make the math—particularly gradient calculations—more manageable. Furthermore, such population analyses do not characterize the statistical and computational complexity of gradient-based learning. Moreover, in factual-recall setups, it is known that strictly orthogonal embeddings are not capacity-optimal, whereas random/non-orthogonal embeddings (i.e., *superposition*) enable near-optimal factual storage (Nichani et al., 2025). At the same time, abandoning the orthogonality assumption introduces token interference that leads to intricate optimization behavior (e.g., oscillatory trajectories Cabannes et al., 2024b); in practice, superposition-based, memory-efficient solutions can also be more challenging to train (Elhage et al., 2022), highlighting a fundamental trade-off between optimization/statistical efficiency and optimal storage capacity.

Motivated by the above gaps, we aim to address the following question.

*Can we characterize the optimization and sample complexity of a transformer with non-orthogonal embeddings trained by gradient descent in the learning of a factual recall task?*

## 1.1 OUR CONTRIBUTIONS

In this paper, we analyze gradient-based learning of a single-layer transformer with an attention+MLP block and random embeddings on a synthetic task inspired by Nichani et al. (2025): the model must retrieve an informative token from a context containing many noisy tokens via attention, then map it to the correct label via factual recall. To mitigate the complex optimization dynamics arising from non-orthogonal embeddings, we follow Bietti et al. (2023); Oymak et al. (2023) and consider a simplified training regime involving only a few gradient steps with finite samples on the attention and value matrices. This perspective effectively zooms in on the “early phase” dynamics of gradient descent, a common focus in the feature-learning literature (Ba et al., 2022; Damian et al., 2022; Dandi et al., 2023; Wang et al., 2025).

Our analysis provides a fine-grained characterization of how vocabulary size  $V$ , sample size  $N$ , embedding dimension  $d$ , sequence length  $L$ , and MLP width  $m$  interact to permit successful gradient-based learning of the recall mechanism. Our main result states that

- The success of learning depends on  $(V, N, d, L, m)$  in a *multiplicative* manner: learning becomes easier as  $(N, d, m)$  increase — reflecting benefits from more data, higher-dimensional (hence more orthogonal) embeddings, and larger MLP width — whereas learning becomes harder as  $(V, L)$  increase, i.e., the task is more difficult with a larger vocabulary or longer sequences. This multiplicative relation is visualized in Figure 1, where we examine how the parameter size  $m \times d$  depend on the vocabulary size  $V$  at different sequence lengths  $L$ .
- Consequently, while optimal capacity and sample complexity can be achieved jointly for short sequences, successful learning on long sequences requires either larger embedding dimension (thus sacrificing capacity) or larger sample sizes (worse statistical complexity).

The multiplicative rate above formalizes the “tradeoff” intuition that smaller embedding dimension  $d$  — which increases superposition and thereby improves storage capacity — simultaneously yields a harder learning problem, as reflected in the required sample size. We complement this with a statistical lower bound showing that the trade-off is inherent for any estimator that accesses only gradient information from the initialized transformer. Finally, although our theory is derived for a specific three-step training algorithm, we empirically observe qualitatively similar multiplicative scaling when the transformer is optimized by gradient descent to low empirical risk.

## 1.2 RELATED WORK

**Learning dynamics of transformers.** A growing line of theory analyzes how transformers acquire specific behaviors from gradient-based training. Much of this literature imposes population-level assumptions and orthogonal/one-hot embeddings to make gradients tractable, often on discrete synthetic tasks (Li et al., 2023; Bietti et al., 2023; Tian et al., 2023a; Nichani et al., 2024; Chen et al., 2024; Ghosal et al., 2024; Chen et al., 2025; Wang et al., 2025). Several works study few-step training regimes as a lens on the “early phase” of feature learning (Bietti et al., 2023; Wang et al., 2025). Beyond discrete settings, related analyses investigate attention learning for continuous inputs and sparse-signal retrieval (Oymak et al., 2023; Marion et al., 2025). A complementary thread focuses on the emergence of in-context learning and induction mechanisms: single- and two-layer attention trained on linear-regression or Markov data provably implements gradient-descent-like updates and generalized induction heads (Von Oswald et al., 2023; Zhang et al., 2024; Chen et al., 2024; Nichani et al., 2024). These results typically rely on simplified settings and do not address storage capacity. In contrast, our work analyzes finite-sample training with random (non-orthogonal) embeddings in an attention+MLP architecture with a particular focus on factual recall.

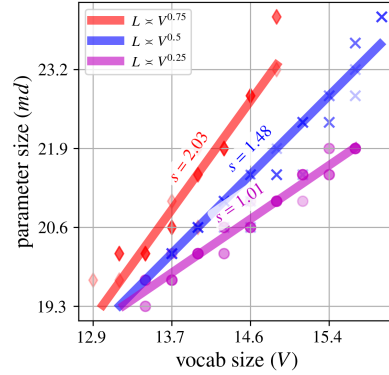


Figure 1: Empirical scaling of parameter count required for GD-trained one-layer transformer to learn factual recall. While the trained model achieves optimal capacity  $V \asymp md$  for small  $L$ , increasing the sequence length  $L$  alters the scaling, suggesting a multiplicative rate.

**Associative memories and storage capacity.** Classical associative memories (Hopfield-type models) study recall of vector patterns and established foundational capacity results (Hopfield, 1982; Amit et al., 1985; McEliece et al., 1988; Krotov & Hopfield, 2016; Demircigil et al., 2017; Ramsauer et al., 2020; Schlag et al., 2021). Recent works adapt associative-memory viewpoints to transformers, modeling inner weights as superpositions of outer products and deriving scaling laws and optimization behaviors (Bietti et al., 2023; Cabannes et al., 2024a;b). In factual recall specifically, random (non-orthogonal) embeddings enable near-parameter-count storage, whereas strictly orthogonal embeddings are not capacity-optimal (Nichani et al., 2025). Various empirical works have studied the mechanisms and scaling behaviors of LLMs in factual association tasks (Petroni et al., 2019; Jiang et al., 2020; Geva et al., 2020; Allen-Zhu & Li, 2024). We provide a theoretical analysis of such mechanisms and quantify how vocabulary size, sequence length, embedding dimension, and MLP width jointly govern learning efficiency. **Our work operates in a setting similar to (Nichani et al., 2025) but allows finite samples and explicitly considers gradient descent dynamics. Our result is similar to the finite-sample results in (Oymak et al., 2023), where the required sample size grows with the dimensionality and sparsity level of informative tokens, while we allow non-orthogonal embeddings and show optimal capacity as in (Nichani et al., 2025) under certain conditions.**

## 2 PROBLEM SETTING

Our goal is to understand the capacity of transformers trained on finite data with non-orthogonal embeddings, in a setting where the relevant information is hidden in a potentially large sequence of non-informative noisy tokens. The attention operation should then identify the relevant token, while the subsequent linear or MLP block can then recall the correct label via an associative memory mechanism. This is similar to the factual recall task studied by Nichani et al. (2025), with simplifications that make the analysis more tractable, as detailed below.

**Notation.**  $\sigma$  denotes the softmax function.  $\mathbb{1}_V := (1, \dots, 1)^\top \in \mathbb{R}^V$  is the  $V$ -dimensional all-ones vector;  $e_i$  is the one-hot vector with a 1 in the  $i$ -th position (dimension understood from context). We use  $\gtrsim$  (resp.  $\lesssim$ ) to mean “ $\geq$ ” (resp. “ $\leq$ ”) up to polylogarithmic factors in  $V$ :  $f_V \gtrsim g_V \iff f_V \geq \text{poly}(\log V)g_V$  and  $f_V \lesssim g_V \iff f_V \leq \text{poly}(\log V)g_V$ , for some fixed polynomial. Lastly,  $\|\cdot\|_2$  denotes the Euclidean norm for vectors and the operator (spectral) norm for matrices.

**Problem setup.** Let the input/output tokens take values from a finite alphabet  $[V] := \{1, \dots, V\}$ . For notational convenience, we represent the alphabet by the one-hot vocabulary  $\mathcal{V} = \{e_1, \dots, e_V\}$ . Each example in the data consists of a length- $L$  input sequence  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_L] \in \mathcal{V}^L$  and a label  $\mathbf{p} \in \mathcal{V}$  generated as follows:

- *Input* tokens are sampled independently and uniformly:  $[\mathbf{x}_1, \dots, \mathbf{x}_L] \sim \text{Unif}(\mathcal{V}^L)$ .
- *Informative position* is a random index  $\ell \sim \text{Unif}([L])$  independent of  $\mathbf{X}$ .
- *Ground-truth function* is a permutation matrix  $\Pi_* \in \{0, 1\}^{V \times V}$ . Labels are generated as the permuted informative token,  $\mathbf{p} = \Pi_* \mathbf{x}_\ell$ , while the remaining tokens are non-informative.

The goal is to identify the correct token position  $\ell$  and learn the target function (permutation)  $\Pi_*$ .

**Transformer architecture.** We consider a basic transformer block which first maps input tokens into a  $d$ -dimensional embedding space where  $d < V$ . The embedding layer is parameterized by  $(\mathbf{Z}_{\text{in}}, \mathbf{Z}_{\text{out}}, \mathbf{z}_{\text{trig}}, \mathbf{z}_{\text{EOS}}) \in \mathbb{R}^{d \times V} \times \mathbb{R}^{d \times V} \times \mathbb{R}^d \times \mathbb{R}^d$ , where

- The input tokens are embedded by the columns of the matrix  $\mathbf{Z}_{\text{in}} \in \mathbb{R}^{d \times V}$ .
- Output tokens are associated with unembedding vectors, which are collected in  $\mathbf{Z}_{\text{out}} \in \mathbb{R}^{d \times V}$ .
- $\mathbf{z}_{\text{trig}}$  is a trigger vector that marks the informative token.
- $\mathbf{z}_{\text{EOS}}$  is the special embedding vector that marks the end-of-sequence.

Given the embedding parameters, we define the self-attention head, parameterized by the key-query matrix  $\mathbf{W}_{\text{KQ}} \in \mathbb{R}^{d \times d}$ , which operates on the embedded sequence of inputs  $\mathbf{Z}_{\text{in}} \mathbf{X} \in \mathbb{R}^{d \times L}$ :

$$\text{attn}(\mathbf{X}; \mathbf{W}_{\text{KQ}}) := \mathbf{Z}_{\text{in}} \mathbf{X} \sigma \left( (\mathbf{z}_{\text{trig}} e_\ell^\top + \mathbf{Z}_{\text{in}} \mathbf{X})^\top \mathbf{W}_{\text{KQ}} \mathbf{z}_{\text{EOS}} \right). \quad (1)$$

The trigger embedding  $\mathbf{z}_{\text{trig}}$  is used to “mark” the informative token with a special direction, mimicking the behavior of previous transformer layers that may learn to flag particular tokens by adding to its residual stream<sup>1</sup> (note that the number of trainable parameters inside softmax can be reduced to  $d$  by collapsing  $\mathbf{W}_{\text{KQ}}\mathbf{z}_{\text{EOS}}$  into a vector). We consider two different learning models: an *Attention-only* model and a width- $m$ , two-layer neural network model *Attention-MLP*, defined as:

$$\hat{p}(\mathbf{X}; \mathbf{V}, \mathbf{W}_{\text{KQ}}) = \begin{cases} \sigma(\mathbf{Z}_{\text{out}}^\top \mathbf{V} \text{attn}(\mathbf{X}; \mathbf{W}_{\text{KQ}})), & \text{Attention only} \\ \sigma(\mathbf{Z}_{\text{out}}^\top \mathbf{V} \phi(\mathbf{W}_{\text{in}} \text{attn}(\mathbf{X}; \mathbf{W}_{\text{KQ}}))), & \text{Attention-MLP} \end{cases} \quad (2)$$

where  $\mathbf{V} \in \mathbb{R}^{d \times d}$  for the *Attention-only* and  $\mathbf{V} \in \mathbb{R}^{d \times m}$ ,  $\mathbf{W}_{\text{in}} \in \mathbb{R}^{m \times d}$  for the *Attention-MLP* model. Note that compared with *Attention-only* model, the *Attention-MLP* model contains an additional set of trainable parameters and nonlinear activation function  $\phi$  before the value matrix. Constructions of the two models for a related factual recall task can be found in (Nichani et al., 2025, Figure 3).

For the *Attention-MLP*, we keep  $\mathbf{W}_{\text{in}}$  fixed at its random initialization. The trainable parameters for both of our models are  $(\mathbf{V}, \mathbf{W}_{\text{KQ}})$ . We use cross-entropy loss to train our model:

$$\mathcal{L}((\mathbf{V}, \mathbf{W}_{\text{KQ}}), (\mathbf{X}, \mathbf{p})) = - \sum_{i=1}^V p_i \log \hat{p}_i.$$

**Training algorithm.** Following Oymak et al. (2023), we consider a 3-step gradient-based algorithm with dataset  $\{(\mathbf{X}_i, \mathbf{p}_i)\}_{i=1}^N$  with a sample size of  $N$ . We initialize our parameters as  $\mathbf{V}^{(0)} = \mathbf{0}$ ,  $\mathbf{W}_{\text{KQ}}^{(0)} = \mathbf{0}$  and use the learning rates  $\eta, \gamma > 0$ :

$$\mathbf{V}^{(1)} = \mathbf{V}^{(0)} - \eta \cdot \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{V}} \mathcal{L}((\mathbf{V}^{(0)}, \mathbf{W}_{\text{KQ}}^{(0)}); (\mathbf{X}_i, \mathbf{p}_i)) \quad (3)$$

$$\mathbf{W}_{\text{KQ}}^{(1)} = \mathbf{W}_{\text{KQ}}^{(0)} - \gamma \cdot \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{W}_{\text{KQ}}} \mathcal{L}((\mathbf{V}^{(1)}, \mathbf{W}_{\text{KQ}}^{(0)}); (\mathbf{X}_i, \mathbf{p}_i)) \quad (4)$$

$$\mathbf{V}^{(2)} = \mathbf{V}^{(1)} - \gamma \cdot \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{V}} \mathcal{L}((\mathbf{V}^{(1)}, \mathbf{W}_{\text{KQ}}^{(1)}); (\mathbf{X}_i, \mathbf{p}_i)). \quad (5)$$

**Network prediction and storage.** Given our model and training method, we use argmax decoding at inference and define the test accuracy as

$$\text{Accuracy} := \mathbb{P}_{(\mathbf{X}, \mathbf{p})}[\mathbf{p} = \mathbf{e}_{\text{pred}(\mathbf{X})}], \quad \text{where} \quad \text{pred}(\mathbf{X}) := \arg \max_{j \in [V]} \hat{p}_j(\mathbf{X}; \mathbf{V}^{(2)}, \mathbf{W}_{\text{KQ}}^{(1)}),$$

where  $\hat{p}(\mathbf{X}; \mathbf{V}^{(2)}, \mathbf{W}_{\text{KQ}}^{(1)})$  is the network output defined in (2). In what follows, we characterize conditions under which the model stores the informative tokens asymptotically, i.e.,  $\text{Accuracy} \rightarrow 1$  as  $V \rightarrow \infty$ , in terms of the relevant parameters  $(V, N, d, L, m)$ .

### 3 MAIN RESULTS

We first present our general theorem on learnability via gradient descent, and then specialize into different regimes to derive more interpretable scaling behaviors in Section 4. We provide a proof sketch in Section C.1, and defer the full proof to Appendix C.

#### 3.1 TECHNICAL ASSUMPTIONS

We first state generic assumptions that apply to both the *Attention-only* and *Attention-MLP* models.

##### Assumption 1.

- **Parameter range:** Let  $L = V^c$  for  $c \in (0, 1)$ ,  $\Omega(V \log V) \leq N = o(VL)$ , and  $V \geq \Omega(1)$ .
- **Learning rate:** We use a sufficiently small learning rate  $\eta = o(1)$  for the initial step (3), and sufficiently large learning rate  $\gamma = \omega(1)$  for the remaining steps (4)-(5) that satisfy Assumption 4.
- **Embeddings:** Let  $\mathbf{Z}_{\text{in}}, \mathbf{Z}_{\text{out}} \in \mathbb{R}^{d \times V}$  be independent Gaussian matrices, and let  $\mathbf{z}_{\text{trig}}, \mathbf{z}_{\text{EOS}} \in \mathbb{R}^d$  be independent Gaussian vectors, all with i.i.d. entries distributed as  $\mathcal{N}(0, 1/d)$ .

<sup>1</sup>The “trigger” terminology is borrowed from (Bietti et al., 2023), where a special previous token “triggers” a retrieval operation in the context of induction heads. Our setup resembles learning only the “induction head” layer assuming the first “previous token head” layer is already in place.

We assume  $c \in (0, 1)$  since in many practical pretraining setups, the context length is smaller than the vocabulary size, and the condition  $L \ll V$  simplifies several terms in the proofs. The lower bound  $N \gtrsim V \log V$  is required so that each element from the alphabet of size  $V$  is seen at least once with high probability. The learning rates follow prior analyses (Oymak et al., 2023; Nichani et al., 2024): a small  $\eta$  ensures that the network’s predictions remain close to uniform after the first step, whereas a large  $\gamma$  is needed to push the attention scores and predictions toward one-hot vectors.

In addition to the above assumptions, we require the transformer model to have sufficient capacity to reach perfect test accuracy. Such conditions are characterized by Nichani et al. (2025). For the *Attention-only* model, we have the following condition (see Nichani et al., 2025, Theorem 3).

**Assumption 2** (Attention-only). *For the Attention-only model, we require  $d \gtrsim \sqrt{V}$ .*

With a nonlinear MLP layer, a smaller embedding dimension can suffice if the width is large enough. Hence for *Attention-MLP* we require the following condition.

**Assumption 3** (Attention-MLP). *For the Attention-MLP model, we assume that*

- **Polynomial activation:**  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  satisfies  $\phi(0), \phi'(0), \phi''(0) \neq 0$ .
- **MLP width:**  $md \gtrsim V$  and  $d \gtrsim V^{\frac{1}{k_*+1}}$ , where  $k_*$  denotes the smallest nonzero Hermite mode of  $\phi$ , i.e.,  $k_* := \min\{k > 0 : \mathbb{E}_{Z \sim \mathcal{N}(0,1)}[\phi(Z)h_k(Z)] \neq 0\}$  where  $h_k$  is the  $k^{\text{th}}$  Hermite polynomial.
- **Initialization:**  $\mathbf{W}_{\text{in}} \in \mathbb{R}^{m \times d}$  are fixed with entries i.i.d. distributed as  $\mathcal{N}(0, 1)$ .

The nonlinear MLP layer allows us to compensate for the embedding dimension and go beyond the  $d \gtrsim \sqrt{V}$  lower bound required by the *Attention-only* model (Assumption 2). Note that  $md \gtrsim V$  is a necessary condition for capacity as shown in (Nichani et al., 2025). The additional requirements imposed on the polynomial activation function appear to be artifacts of our three-step GD analysis, and we anticipate that they could be relaxed when considering a longer training horizon.

### 3.2 LEARNABILITY STATEMENT

Now we are ready to present our main theorem on the complexity of learning the factual recall task. Specifically, transformer learns the desired mechanism when the signal term dominates the noise and bias terms as stated below.

**Theorem 1.** *Let Assumptions 1 and 3 hold for Attention-MLP, and 1 and 2 hold for Attention-only. The Attention-MLP model achieves Accuracy =  $1 - o_V(1)$  with probability  $1 - o_V(1)$  whenever*

$$\underbrace{\frac{1}{VL^2}}_{\text{Signal}} \gtrsim \underbrace{\frac{1}{N\sqrt{Ld}(d \wedge L)}}_{\text{Gradient noise}} + \underbrace{\frac{1}{N\sqrt{Vd}(d \wedge L)}}_{\text{Mean bias}} + \underbrace{\frac{1}{Nd\sqrt{m}}}_{\text{MLP noise}}. \quad (6)$$

*For the Attention-only model, the same holds with the last MLP noise term removed.*

Theorem 1 characterizes learnability as a function of  $(V, N, d, L, m)$  and identifies the following terms that impact the gradient signal-to-noise ratio:

1. *Signal* measures the alignment between the key–query weights  $\mathbf{W}_{\text{KQ}}^{(1)}$  and the trigger  $\mathbf{z}_{\text{trig}}$ .
2. *Gradient noise* is due to the concentration error in the update of  $\mathbf{W}_{\text{KQ}}^{(1)}$ .
3. *Mean bias* arises from the nonzero mean of token vectors  $\{\mathbf{X}_i\}_{i=1}^N$ .
4. *MLP noise* reflects the randomness in the MLP weight matrix  $\mathbf{W}_{\text{in}}$  in *Attention-MLP*.

We make the following observations.

- **Multiplicative scaling.** Note that the parameters  $(V, N, d, L, m)$  interact in a multiplicative fashion. For example, the noise and bias terms in (6) all decay with  $(N \times d)$ , suggesting that increasing the embedding dimensions  $d$  can lower the statistical complexity of learning the correct recall mechanism. While the full 5-parameter trade-off can be opaque, in Section 4 we focus on specific regimes that lead to simplification of the scaling relationship and validate the rate empirically.



- **Optimal storage & sample complexity.** Recall that the capacity-optimal construction for the factual recall task requires  $md \gtrsim V$  parameters (or  $d^2 \gtrsim V$  for *Attention-only*); and as discussed earlier, a sample size  $N \asymp V \log V$  is necessary to observe all distinct tokens. (6) implies that in the small- $L$  regime, the optimized transformer achieve optimal capacity and sample complexity simultaneously. For longer sequences, however, these two conditions may not be achieved at the same time, i.e., one must increase either the network width or sample size beyond optimality to learn the task — this confirms the empirical observation in Figure 1.

### 3.3 STATISTICAL LOWER BOUND

Theorem 1 provides an upper bound (i.e., sufficient condition) on the model and sample size for learning factual recall under a 3-gradient-step optimization procedure. We complement this sufficient condition with a lower bound indicating that the multiplicative dependence on the problem parameters is partly statistical; that is, the scaling behavior will be observed in any model satisfying the broader conditions stated below. Our lower bound applies to statistical methods that can query the dataset through the attention outputs at initialization,  $\mathbf{h}_i := \text{attn}(\mathbf{X}_i, \mathbf{W}_{\text{KQ}}^{(0)})$ . In particular, we consider queries of the form  $\{\mathbf{h}_i, \mathbf{h}_i \mathbf{h}_i^\top\}_{i=1}^N$  as the gradient with respect to the key–query matrix  $\mathbf{W}_{\text{KQ}}$  depends on these quantities (see (10)). The statement is given below:

**Theorem 2 (Informal).** *Any method that relies on the noisy version of the queries  $\{\mathbf{h}_i, \mathbf{h}_i \mathbf{h}_i^\top\}_{i=1}^N$  fails, i.e., Accuracy  $\not\rightarrow 1$  with finite probability, if  $N \lesssim V \min\{1, L/d^2\}$ .*

The complete statement of Theorem 2 is deferred to Theorem 4 in Appendix D. We observe that the lower bound does not exactly match our upper bound in Theorem 1, as *Signal*  $\lesssim$  *Gradient Noise* in (6) is stronger than the stated lower bound. This being said, Theorem 2 also confirms the multiplicative scaling, hence suggesting the trade-off between capacity and sample efficiency is present in a boarder class of learning algorithms. A stronger computational lower bound for transformers and gradient-based optimization is an interesting problem we leave for future work.

## 4 IMPLICATIONS AND EMPIRICAL VERIFICATIONS

In this section, we leverage our main theorem to obtain more concrete scalings between parameters, and present empirical evidence on the derived multiplicative rate.

### 4.1 ATTENTION-ONLY MODEL

We start with the *Attention-only model* which gives a simpler phase diagram.

**Corollary 1.** *For the Attention-only model, the bottleneck term in (6) is the Mean bias term, and Theorem (1) is equivalent to requiring  $d \gtrsim \max\{\sqrt{V}, V^{\frac{1}{3}} L^{\frac{4}{3}} / N^{\frac{2}{3}}\}$ .*

We make the following observations:

- The condition in Corollary 1 is the maximum of two terms, where  $d \gtrsim \sqrt{V}$  is due to the capacity requirement in Assumption 2, whereas the second term ensures *Signal*  $\gtrsim$  *Mean bias* and implies a multiplicative scaling between the sample size  $N$  and embedding dimension  $d$  (i.e., increasing one of the parameters can compensate for the other).
- Note that the *Mean bias* term arises from a nonzero token mean, which can potentially be alleviated by centering the tokens, as is effectively done by normalization layer. Exploring the effect of applying normalization in this model is an interesting direction for future work.

**Empirical Findings.** We run the three-step gradient descent algorithm on an *Attention-only* model over varying  $V$  and  $d$ , and report the accuracies in the heatmaps (Figure 2). The plots are in log-log scale; therefore, the slopes give the exponent  $s$  in  $d \asymp V^s$ . As shown in the top row of Figures 2a-2b, the slope for relatively small  $L$  (where  $L \asymp \sqrt{V}$ ) matches the optimal capacity condition  $d \asymp \sqrt{V}$ . By contrast, when the context window is larger ( $L \asymp V$ ), the requirement becomes  $d \asymp V$ , which is also reflected in the experimental results, as observed in the bottom panel of Figure 2a.

In Figure 2b we run experiments with increasing sample size to observe the multiplicative trade-off. As seen in the bottom figure of Figure 2b, increasing the sample size from  $V \log V$  to  $V^{1.5}$  reduces the dimension exponent from 1 to 0.7 (the theoretical value is  $s = 0.66$ ). Finally, the learnability thresholds for  $L \asymp V$  in Figures 2a and 2b are plotted together in Figure 2c, to illustrate that increasing the sample size can compensate for the number of parameters in the network.

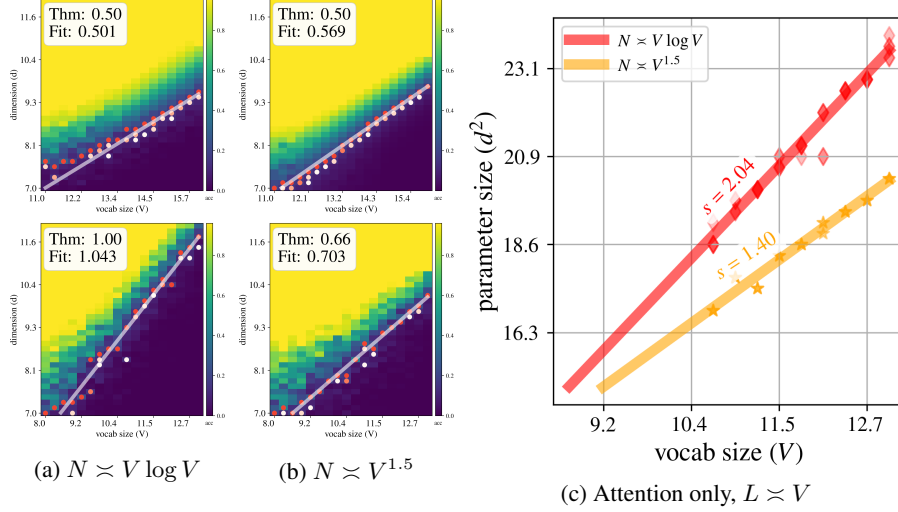


Figure 2: Empirical scaling of embedding dimension (left) and parameter count (right) via three-step GD for the *Attention-only* model. In (a) and (b), top-left and top-right use  $L \asymp \sqrt{V}$ ; bottom-left and bottom-right use  $L \asymp V$ . In the right panel, the  $L \asymp V$  case is shown under two sample-size regimes,  $N \asymp V \log V$  and  $N \asymp V^{1.5}$ . *Line fitting*: We identify in the heatmaps the smallest embedding dimension that achieves accuracies  $\{0.1, 0.125, 0.15\}$  and perform a least squares fit. The slopes of the fitted lines and their theoretical counterparts are reported on the heatmaps. Differences in transparency in (c) are due to overlapping points.

## 4.2 ATTENTION-MLP MODEL

For the attention-MLP model, the nonlinear MLP layer introduces additional phases as stated below.

**Corollary 2.** *For the Attention-MLP model, Theorem 1 translates to  $md \gtrsim V$  and*

$$\text{Signal} \gtrsim \begin{cases} \text{MLP noise,} & m = o(d^2 L) \text{ and } m = o(dV) \\ \text{Gradient noise 2,} & V \gtrsim dL \text{ and } m \gtrsim d^2 L \\ \text{Mean Bias,} & V = o(dL) \text{ and } m \gtrsim dV, \end{cases}$$

where

- Signal  $\gtrsim$  MLP noise is equivalent to  $Nd \gtrsim VL^2/\sqrt{m}$ .
- Signal  $\gtrsim$  Gradient noise 2 is equivalent to  $d\sqrt{N} \gtrsim VL^{\frac{1}{4}}$
- Signal  $\gtrsim$  Mean Bias is equivalent to  $dN^{\frac{2}{3}} \gtrsim L^{\frac{4}{3}}V^{\frac{1}{3}}$ .

The phase diagram for the Attention-MLP model is richer than *Attention-only*, as we can trade off  $m$  and  $d$  and hence use a smaller embedding dimension; this results in potentially different dominant terms in the gradient. In particular, since large  $L$  and  $d$  entails larger magnitude of *Mean Bias* (as in the *Attention-only* setting), we know that by increasing the MLP width  $m$  and thereby reducing the required embedding dimension  $d$ , we may suppress this bias term.

**Empirical Findings.** We run the 3-step gradient descent algorithm on an Attention-MLP network over varying  $V$  and  $d$  and plot the accuracies in Figures 3 and 4. We take the nonlinearity to be the mixture of two Hermite polynomials  $\phi = 0.7h_2 + 0.3h_3$ , satisfying the conditions in Assumption 3. We run experiments with width  $m \asymp d^2$  and  $m \asymp d^3$ . Due to the prohibitive cost of increasing the width further, we restrict ourselves to the *MLP noise*-dominated region.

In Figure 1, we plot the scaling of the number of parameters ( $md$ ) as a function of vocabulary size  $V$  for different sequence-length regimes in  $L$ . We observe that  $L \asymp V^{0.25}$  requires  $md \asymp V$ , which

is the optimal capacity, as predicted by our theory. As  $L$  increases, we need more parameters to achieve the same capacity, as observed in the  $L \asymp V^{0.5}$  and  $L \asymp V^{0.75}$  cases in Figure 1, where the slopes agree with our theoretical predictions as well (see Figures 3a and 3b).

We further test the effect of sample size in Figure 3, where we use  $L \asymp V^{0.5}$  and  $m \asymp d^2$ . We plot both heat maps in Figures 3a and 3b, and the fitted lines for  $L \asymp V^{0.5}$  together in Figure 3c. Note that we state the plot in terms of parameter count, which scales as  $md \asymp d^3$ , so the slopes from the heat map are scaled accordingly. We observe that increasing  $N$  from  $N \asymp V \log V$  to  $N \asymp V^{1.5}$  reduces the network size to the optimal level, aligning with our theoretical prediction. The heatmap versions of these experiments are shown in Figures 3a and 3b.

Lastly, we probe the width scaling by keeping the sample size  $N \asymp V \log V$  and  $L \asymp V^{0.5}$  fixed in Figure 4. Here, we observe that we can reduce the embedding-dimension requirement by increasing  $m$  (Figures 4a and 4b), though it increases the total parameter count overall, as seen in Figure 4c, since width must grow proportionally more than  $d$  to achieve the same accuracy. This is also consistent with our theoretical prediction.

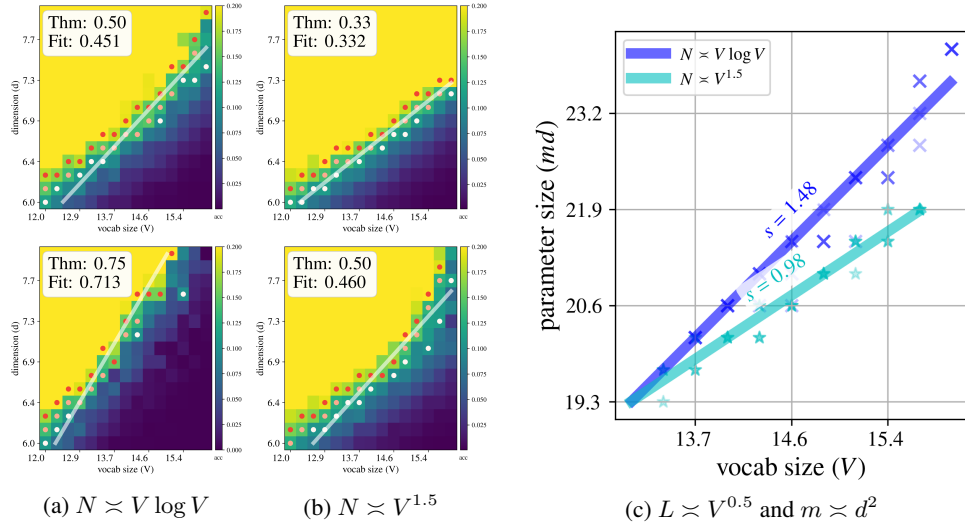


Figure 3: Empirical scaling of embedding dimension (left) and parameter count (right) for the *Attention-MLP* model under  $N \asymp V \log V$  and  $N \asymp V^{1.5}$ . In (a) and (b), top-row uses  $L \asymp V^{0.5}$ ; bottom-row uses  $L \asymp V^{0.75}$ . The right panel also shows  $L \asymp V^{0.5}$  under both sample-size regimes.

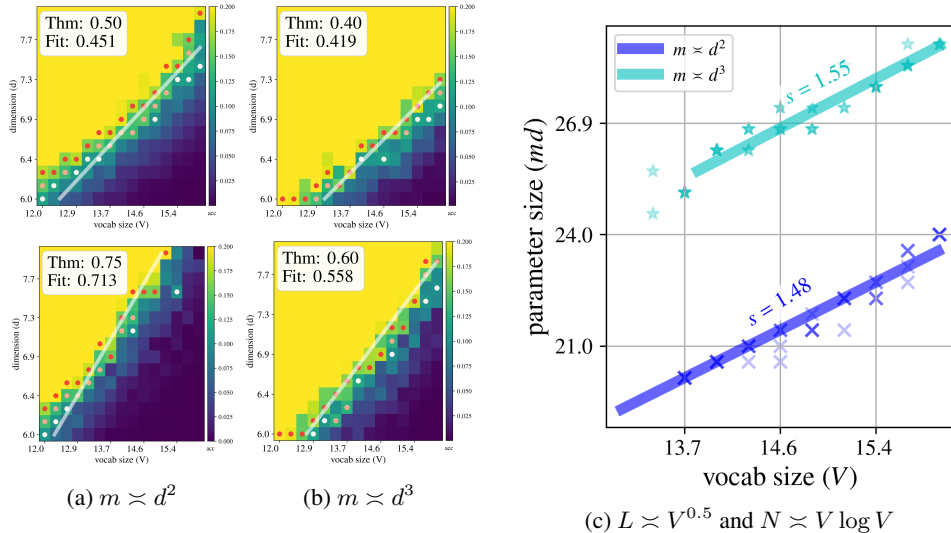


Figure 4: Empirical scaling of embedding dimension (left) and parameter count (right) for the *Attention-MLP* model under two width regimes,  $m \asymp d^2$  and  $m \asymp d^3$ . In (a) and (b), top-row uses  $L \asymp V^{0.5}$ ; bottom-row uses  $L \asymp V^{0.75}$ . The right panel also shows  $L \asymp V^{0.5}$  under both width regimes.



### 4.3 BEYOND EARLY PHASE OF TRAINING

While our theoretical analysis handles a particular 3-gradient-step training procedure, we empirically observe qualitatively similar multiplicative scalings when the transformer model is optimized beyond the “early phase”. Specifically, we train our *Attention-only* model for multiple steps using (i) full-batch gradient descent and (ii) Adam (Kingma & Ba, 2015) with mini-batch gradients. Throughout this section, we use a sample size  $N \asymp V \log V$ .

*Full-batch gradient descent:* We use learning rate  $\eta = 0.5$  and continue training until the test accuracy does not improve by more than 0.01 for 10 consecutive checks. In Figures 5a and 5b, we provide heatmaps for  $L \asymp V$  and  $L \asymp \sqrt{V}$ . We observe that when  $L \asymp \sqrt{V}$ , the slope indicates the network is at the optimal capacity condition; this is also reflected by the slope 1.06 in Figure 5c. By contrast, for large  $L$  the slope significantly shifts and becomes suboptimal, confirming the multiplicative relation established in Section 3.

*Adam with mini-batch gradients:* We use layer normalization in both the attention and output layers and choose learning rate  $\eta = 0.005$ . We specifically use batch size  $\lfloor N/40 \rfloor$  and run the algorithm for 3 epochs. In Figure 6, we provide heatmaps for  $L \asymp V$  and  $L = 8$  at the end of epochs 2 and 3. We observe that for  $L \asymp V$  in early training, the slope is suboptimal, while training the network for one more epoch improves the capacity condition to a near optimal level. In contrast, for  $L = 8$ , we observe that the network does not exhibit a suboptimal phase at the end of epoch 2, which is in line with our theoretical findings. A rigorous analysis of the full gradient descent dynamics is left for future work.

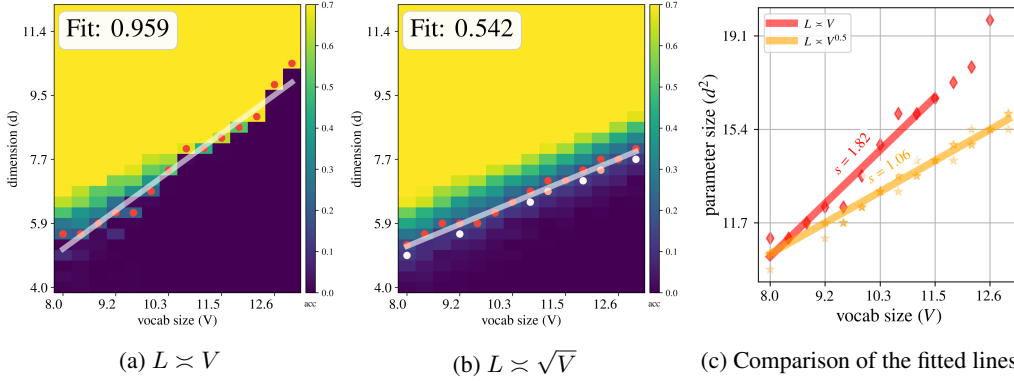


Figure 5: Empirical scaling of embedding dimension (a,b) and parameter count (c) for the *Attention-only* model trained by multiple-step GD.

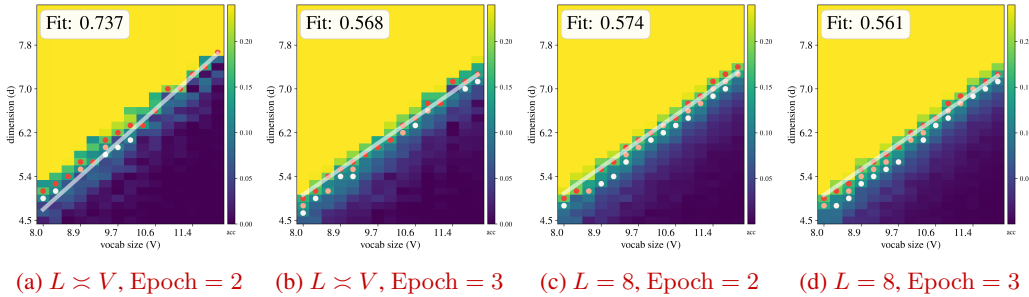


Figure 6: Empirical scaling of embedding dimension for the *Attention-only* model trained by Adam.

## 5 PROOF OVERVIEW: POPULATION ANALYSIS

In this section, we outline the ideas for the proof of Theorem 1. For presentation, we consider the *Attention-only* model with population dynamics and orthogonal embeddings. Since we do not use positional encoding in the model, without loss of generality, we fix the correct position to  $\ell = 1$ .

In the proof, we study the attention scores in (1) and characterize the conditions under which they align with the trigger vector. Once attention can distinguish informative tokens, the remaining part

reduces to learning a linearly separable problem, which is well understood. The pre-softmax scores evaluated on a fresh sequence  $\mathbf{X}_{\text{in}}$ , with the key-query matrix given by the first gradient-descent iterate  $\mathbf{W}_{\text{KQ}}^{(1)}$ , is given as

$$\text{scores} := (\mathbf{z}_{\text{trig}} \mathbf{e}_1^\top + \mathbf{Z}_{\text{in}} \mathbf{X}_{\text{in}})^\top \mathbf{W}_{\text{KQ}}^{(1)} \mathbf{z}_{\text{EOS}}. \quad (7)$$

For the proof overview part, we analyze the following simplified form of the scores (for the full expression, see (10)):

$$\text{scores} \approx \underbrace{\gamma \mathbf{X}_{\text{in}}^\top \mathbf{Z}_{\text{in}}^\top \left( \frac{1}{NL} \sum_{i=1}^N \mathbf{Z}_{\text{in}} \mathbf{X}_i \mathbf{X}_i^\top \mathbf{Z}_{\text{in}}^\top (\mathbf{V}^{(1)})^\top \mathbf{Z}_{\text{out}} (\mathbf{p}_i - \frac{1}{V} \mathbb{1}_V) \right)}_{\text{Non-informative}} \quad (8)$$

$$+ \underbrace{\gamma \|\mathbf{z}_{\text{trig}}\|_2^2 \mathbf{e}_1 \left( \frac{1}{NL} \sum_{i=1}^N \mathbf{x}_{i,1}^\top \mathbf{Z}_{\text{in}}^\top (\mathbf{V}^{(1)})^\top \mathbf{Z}_{\text{out}} (\mathbf{p}_i - \frac{1}{V} \mathbb{1}_V) \right)}_{\text{Informative}}. \quad (9)$$

Here  $\mathbf{V}^{(1)}$  denotes the first iteration of the value matrix given in (3). The informative term in (9) captures the alignment between the trigger vector in the fresh input and the one in the learned weights  $\mathbf{W}_{\text{KQ}}^{(1)}$ , and therefore contains the position information of the informative token. By contrast, the non-informative term in (8) reflects correlations between tokens and does not carry any information about the token’s position. The proof characterizes conditions under which the informative term in (9) dominates, which is sufficient for attention to identify the correct position.

To study population dynamics with orthogonal embeddings, we set  $\mathbf{Z}_{\text{in}} = \mathbf{Z}_{\text{out}} = \mathbf{I}_V$  and take  $N \rightarrow \infty$  while other parameters remain fixed. Under population dynamics, we first observe that  $\mathbf{V}^{(1)} = O(\eta) \mathbf{\Pi}_*$ , where  $\eta$  is the learning rate of the first step in (3). Then, we can write

$$\begin{aligned} \text{Non-informative} &= \mathbf{X}_{\text{in}}^\top \frac{O(\eta\gamma)}{NL} \sum_{i=1}^N \mathbf{X}_i \mathbf{X}_i^\top \mathbf{\Pi}_*^\top (\mathbf{p}_i - \frac{1}{V} \mathbb{1}_V) \\ &\xrightarrow{N \rightarrow \infty} \frac{O(\eta\gamma)}{L} \mathbf{X}_{\text{in}}^\top \mathbb{E}[\mathbf{X} \mathbf{X}^\top (\mathbf{x}_1 - \frac{1}{V} \mathbb{1}_V)] = (1 - \frac{1}{V}) \frac{O(\eta\gamma)}{VL} \mathbb{1}_L, \end{aligned}$$

where the last equality follows since  $\mathbf{X}_{\text{in}}$  has one-hot columns. On the other hand, we have

$$\text{Informative} = O(\eta\gamma) \mathbf{e}_1 \left( \frac{1}{NL} \sum_{i=1}^N \mathbf{x}_{i,1}^\top \mathbf{\Pi}_*^\top (\mathbf{p}_i - \frac{1}{V} \mathbb{1}_V) \right) = (1 - \frac{1}{V}) \frac{O(\eta\gamma)}{L} \mathbf{e}_1,$$

where we used  $\mathbf{p}_i = \mathbf{\Pi}_* \mathbf{x}_{i,1}$ , in the last step. By choosing learning rates that guarantee  $\eta\gamma \rightarrow \infty$ , we can show that the attention probabilities align with  $\mathbf{e}_1$  and eventually select the correct position. The reader may refer to Appendix C.1 for an extended proof overview of the empirical dynamics with non-orthogonal embeddings, where we detail how each term in Theorem 1 arises from the terms in (8)-(9).

## 6 CONCLUSION

In this paper, we derived precise asymptotic rates for learning with gradient descent on transformers trained on a simple recall task with random embeddings and finite samples. Our analysis and experiments reveal a rich picture of multiplicative scalings between various problem parameters, showing that parameter count is not the only important factor controlling capacity when learning with finite samples on large noisy sequences. Our results suggest that finer control of the data distribution may be necessary for learning efficiently at optimal capacity, for instance by ensuring sequences are less noisy and more informative, hoping that the discovered mechanisms are robust to harder settings. This is reminiscent of the procedures used for long context extension in LLMs, where most of training happens on shorter sequences, but the final models are extended to work with very long sequences, and empirically do well on retrieval tasks such as “needle-in-a-haystack” (e.g., [Gemini Team, 2024](#)), which resembles our theoretical setup. Analyzing similar scalings in more structured data distributions and architectures is thus an interesting avenue for future work.

## REFERENCES

- Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.3, knowledge capacity scaling laws. *arXiv preprint arXiv:2404.05405*, 2024.
- Daniel J Amit, Hanoch Gutfreund, and Haim Sompolinsky. Storing infinite numbers of patterns in a spin-glass model of neural networks. *Physical Review Letters*, 55(14):1530, 1985.
- Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-dimensional asymptotics of feature learning: How one gradient step improves the representation. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a transformer: A memory viewpoint. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- Vivien Cabannes, Elvis Dohmatob, and Alberto Bietti. Scaling laws for associative memories. In *International Conference on Learning Representations (ICLR)*, 2024a.
- Vivien Cabannes, Berfin Simsek, and Alberto Bietti. Learning associative memories with gradient descent. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2024b.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. Quantifying memorization across neural language models. In *The Eleventh International Conference on Learning Representations*, 2022.
- Lei Chen, Joan Bruna, and Alberto Bietti. Distributional associations vs in-context reasoning: A study of feed-forward and attention layers. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025.
- Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Unveiling induction heads: Provable training dynamics and feature learning in transformers. *Advances in Neural Information Processing Systems*, 37:66479–66567, 2024.
- Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations with gradient descent. In *Conference on Learning Theory*, pp. 5413–5452. PMLR, 2022.
- Yatin Dandi, Florent Krzakala, Bruno Loureiro, Luca Pesce, and Ludovic Stephan. Learning two-layer neural networks, one (giant) step at a time. *arXiv preprint arXiv:2305.18270*, 2023.
- Mete Demircigil, Judith Heusel, Matthias Löwe, Sven Upgang, and Franck Vermet. On a model of associative memory with huge storage capacity. *Journal of Statistical Physics*, 168:288–299, 2017.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of superposition. *arXiv preprint arXiv:2209.10652*, 2022.
- Gemini Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. *arXiv preprint arXiv:2012.14913*, 2020.
- Gaurav Ghosal, Tatsunori Hashimoto, and Aditi Raghunathan. Understanding finetuning for factual knowledge extraction. *arXiv preprint arXiv:2406.14785*, 2024.
- John J Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.
- Yixiao Huang, Hanlin Zhu, Tianyu Guo, Jiantao Jiao, Somayeh Sojoudi, Michael I Jordan, Stuart Russell, and Song Mei. Generalization or hallucination? understanding out-of-context reasoning in transformers. *arXiv preprint arXiv:2506.10887*, 2025.

- Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438, 2020.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR*, 2015.
- Dmitry Krotov and John J Hopfield. Dense associative memory for pattern recognition. *Advances in neural information processing systems*, 29, 2016.
- Yuchen Li, Yuanzhi Li, and Andrej Risteski. How do transformers learn topic structure: Towards a mechanistic understanding. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2023.
- Pierre Marion, Raphael Berthier, Gerard Biau, and Claire Boyer. Attention layers provably solve single-location regression. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025.
- Robert J McEliece, Edward C Posner, Eugene R Rodemich, and Santosh S Venkatesh. The capacity of the hopfield associative memory. *IEEE transactions on Information Theory*, 33(4):461–482, 1988.
- John X Morris, Chawin Sitawarin, Chuan Guo, Narine Kokhlikyan, G Edward Suh, Alexander M Rush, Kamalika Chaudhuri, and Saeed Mahloujifar. How much do language models memorize? *arXiv preprint arXiv:2505.24832*, 2025.
- Eshaan Nichani, Alex Damian, and Jason D Lee. How transformers learn causal structure with gradient descent. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2024.
- Eshaan Nichani, Jason D Lee, and Alberto Bietti. Understanding factual recall in transformers via associative memories. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025.
- Samet Oymak, Ankit Singh Rawat, Mahdi Soltanolkotabi, and Christos Thrampoulidis. On the role of attention in prompt-tuning. In *International Conference on Machine Learning*, 2023.
- Yang Peng, Yuchen Xin, and Zhihua Zhang. Matrix rosenthal and concentration inequalities for markov chains with applications in statistical learning. *arXiv preprint arXiv:2508.04327*, 2025.
- Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. Language models as knowledge bases? *arXiv preprint arXiv:1909.01066*, 2019.
- Hubert Ramsauer, Bernhard Schäfl, Johannes Lehner, Philipp Seidl, Michael Widrich, Thomas Adler, Lukas Gruber, Markus Holzleitner, Milena Pavlović, Geir Kjetil Sandve, et al. Hopfield networks is all you need. *arXiv preprint arXiv:2008.02217*, 2020.
- Adam Roberts, Colin Raffel, and Noam Shazeer. How much knowledge can you pack into the parameters of a language model? *arXiv preprint arXiv:2002.08910*, 2020.
- Jonathan Scarlett and Volkan Cevher. An introductory guide to fano’s inequality with applications in statistical estimation. *arXiv preprint arXiv:1901.00555*, 2019.
- Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. Linear transformers are secretly fast weight programmers. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2021.
- Yuangdong Tian, Yiping Wang, Beidi Chen, and Simon S Du. Scan and snap: Understanding training dynamics and token composition in 1-layer transformer. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023a.
- Yuangdong Tian, Yiping Wang, Zhenyu Zhang, Beidi Chen, and Simon Du. Joma: Demystifying multilayer transformers via joint dynamics of mlp and attention. *arXiv preprint arXiv:2310.00535*, 2023b.

- Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvinov, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pp. 35151–35174. PMLR, 2023.
- Zixuan Wang, Eshaan Nichani, Alberto Bietti, Alex Damian, Daniel Hsu, Jason D Lee, and Denny Wu. Learning compositional functions with transformers from easy-to-hard data. In *Conference on Learning Theory (COLT)*, 2025.
- Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. *Journal of Machine Learning Research*, 25(49):1–55, 2024.
- Nicolas Zucchet, Jörg Bornschein, Stephanie Chan, Andrew Lampinen, Razvan Pascanu, and Soham De. How do language models learn facts? dynamics, curricula and hallucinations. *arXiv preprint arXiv:2503.21676*, 2025.



## A APPENDIX

## CONTENTS

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Our Contributions . . . . .	2
1.2	Related Work . . . . .	2
<b>2</b>	<b>Problem Setting</b>	<b>3</b>
<b>3</b>	<b>Main Results</b>	<b>4</b>
3.1	Technical Assumptions . . . . .	4
3.2	Learnability Statement . . . . .	5
3.3	Statistical Lower Bound . . . . .	6
<b>4</b>	<b>Implications and Empirical Verifications</b>	<b>6</b>
4.1	Attention-only Model . . . . .	6
4.2	Attention-MLP Model . . . . .	7
4.3	Beyond Early Phase of Training . . . . .	9
<b>5</b>	<b>Proof Overview: Population analysis</b>	<b>9</b>
<b>6</b>	<b>Conclusion</b>	<b>10</b>
<b>A</b>	<b>Appendix</b>	<b>14</b>
<b>B</b>	<b>Preliminaries</b>	<b>16</b>
<b>C</b>	<b>Proof of Theorem 1</b>	<b>17</b>
C.1	Proof sketch for Theorem 1 . . . . .	17
C.2	Attention scores and their asymptotic scaling . . . . .	18
C.3	Proof of Theorem 3 . . . . .	20
C.3.1	Concentration bound for $\mathbf{s}_1$ . . . . .	20
C.3.2	Concentration bound for $\mathbf{s}_2$ . . . . .	35
C.3.3	Concentration bound for $\mathbf{s}_3$ . . . . .	42
<b>D</b>	<b>Lower Bound</b>	<b>47</b>
<b>E</b>	<b>Auxiliary Statements</b>	<b>49</b>
E.1	A nice event characterization . . . . .	49
E.2	Gaussian matrices and related statements . . . . .	56
E.3	Multinomial distribution and related statements . . . . .	58
<b>F</b>	<b>Miscellaneous</b>	<b>68</b>

756	F.1 Rosenthal-Burkholder inequality and corollaries . . . . .	69
757		
758		
759		
760		
761		
762		
763		
764		
765		
766		
767		
768		
769		
770		
771		
772		
773		
774		
775		
776		
777		
778		
779		
780		
781		
782		
783		
784		
785		
786		
787		
788		
789		
790		
791		
792		
793		
794		
795		
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		

**LLM Usage.** Large language models are used to polish the abstract and find relevant references for the related work section.

## B PRELIMINARIES

**Proof organization.** We combine the proof for both models: Using  $\phi(x) = x$  and  $m = \infty$  is valid for applying the arguments for *Attention-only* model in the proof. As the network succeeds in storing all informative tokens only when attention selects the correct position, we focus on how attention learns the correct index and under what conditions. This is the bottleneck in our analysis under Assumptions 1, 2, and 3. Accordingly, we study the pre-softmax scores in (1). Theorem 3 characterizes the scaling of these terms and yields (6). Because the proof involves lengthy expressions, we provide a proof sketch in Section C.1 and refer readers to the corresponding parts of the formal proof.

**Additional Notation.** For a vector  $\mathbf{x} \in \mathbb{R}^V$ , we use  $\text{diag}(\mathbf{x}) \in \mathbb{R}^{V \times V}$  denotes the diagonal matrix which has the same diagonal entries with  $\mathbf{x}$ , while for a matrix  $\mathbf{A}$ ,  $\text{diag}(\mathbf{A}) \in \mathbb{R}^V$  denotes the column vector whose elements coincide with the diagonal entries of  $\mathbf{A}$ . For a random variable  $\mathbf{w}$ ,  $\mathbb{E}_{\mathbf{w}}[\cdot]$  denotes taking expectation with respect to  $\mathbf{w}$  and keeping the remaining independent terms fixed. Similarly, we use  $\mathbb{E}[\cdot|\mathbf{w}]$  for conditional expectation, conditioned on  $\mathbf{w}$ . We use  $\mathbb{1}_{\text{Event}}$  as an indicator function, which takes values  $\{0, 1\}$  depending on the event holds or not. We use  $C$  to denote any constant in the upper-bound, which might depend on  $\phi$ .

Since we do not use positional encoding in the model, without loss of generality we can fix the informative index  $\ell = 1$ . We define the sequence of non-informative tokens as  $\mathbf{N}_i := [\mathbf{x}_{i,2}, \dots, \mathbf{x}_{i,L}]^\top$ . We will denote the rows of  $\mathbf{W}_{\text{in}}$  with  $\{\mathbf{w}_k\}_{k=1}^m$ . For compact representation the attention with the trigger we define

$$\mathbf{Z}_{\text{in}} =: [\mathbf{Z}_{\text{in}} \quad \mathbf{z}_{\text{trig}}] \quad \text{and} \quad \mathbf{X}_i = \begin{bmatrix} \mathbf{x}_{i,1}^\top & 1 \\ \mathbf{N}_i & 0 \end{bmatrix} \in \mathbb{R}^{L \times (V+1)}$$

With this notation, we can write the iterates in three-step GD. Let

$$\boldsymbol{\alpha}^{(0)} := \sigma\left((\mathbf{z}_{\text{trig}} \mathbf{e}_\ell^\top + \mathbf{Z}_{\text{in}} \mathbf{X})^\top \mathbf{W}_{\text{KQ}}^{(0)} \mathbf{z}_{\text{EOS}}\right).$$

We have

$$\begin{aligned} \mathbf{V}^{(1)} &= \mathbf{Z}_{\text{out}} \left( \frac{\eta}{N} \sum_{i=1}^N (\mathbf{p}_i - \hat{\mathbf{p}}_i^{(0)}) \phi((\boldsymbol{\alpha}_i^{(0)})^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{W}_{\text{in}}^\top) \right) \\ \mathbf{W}_{\text{KQ}}^{(1)} &= \mathbf{Z}_{\text{in}} \frac{\gamma}{N} \sum_{i=1}^N \mathbf{X}_i^\top (\text{diag}(\boldsymbol{\alpha}_i^{(0)}) - \boldsymbol{\alpha}_i^{(0)} (\boldsymbol{\alpha}_i^{(0)})^\top) \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{W}_{\text{in}}^\top \\ &\quad \times \text{diag}\left(\phi'(\mathbf{W}_{\text{in}} \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \boldsymbol{\alpha}_i^{(0)})\right) (\mathbf{V}^{(1)})^\top \mathbf{Z}_{\text{out}} (\mathbf{p}_i - \hat{\mathbf{p}}_i^{(1)}) \mathbf{z}_{\text{EOS}}^\top \end{aligned} \quad (10)$$

For notational convenience, we define the noise due to finite width as

$$\begin{aligned} \text{FW}(\mathbf{W}_{\text{in}}; \mathbf{Z}_{\text{in}}, \mathbf{X}_i, \mathbf{X}_j) &:= \frac{1}{m} \left( \mathbf{W}_{\text{in}}^\top \text{diag}\left(\phi'\left(\frac{1}{L} \mathbf{W}_{\text{in}} \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L\right)\right) \phi\left(\frac{1}{L} \mathbf{W}_{\text{in}} \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L\right) \right. \\ &\quad \left. - \mathbb{E}_{\mathbf{W}_{\text{in}}} \left[ \mathbf{W}_{\text{in}}^\top \text{diag}\left(\phi'\left(\frac{1}{L} \mathbf{W}_{\text{in}} \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L\right)\right) \phi\left(\frac{1}{L} \mathbf{W}_{\text{in}} \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L\right) \right] \right). \end{aligned}$$

and the terms arising in the expected value term as

$$\begin{aligned} \alpha_{ij} &:= \mathbb{E}_{\mathbf{w}} \left[ \phi'\left(\frac{1}{L} \mathbf{w}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L\right) \phi'\left(\frac{1}{L} \mathbf{w}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L\right) \right], \\ \beta_{ij} &:= \mathbb{E}_{\mathbf{w}} \left[ \phi''\left(\frac{1}{L} \mathbf{w}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L\right) \phi\left(\frac{1}{L} \mathbf{w}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L\right) \right]. \end{aligned}$$

## C PROOF OF THEOREM 1

### C.1 PROOF SKETCH FOR THEOREM 1

For convenience, we fix  $\Pi_* = \mathbf{I}_V$  and, accordingly,  $\mathbf{p}_i = \mathbf{x}_{i,1}$ , and consider the *Attention-only* model unless stated otherwise; however the derivations for *Attention-only* model holds also for *Attention-MLP*. We study the pre-attention scores given in (7), with the explicit formula in (8)–(9). We derive the terms in (6) in three parts:

- In the first part, we analyze the informative component (9) and derive the scaling of the *Signal* term in (6).
- In the second part, we analyze the non-informative component (8) and derive the scaling of *Gradient noise* and *Mean bias* in (6), corresponding to its mean and bias components.
- In the third part, we consider the *Attention-MLP* model and derive the scaling of *MLP noise* in (8).

Before proceeding, we note that both the informative and non-informative terms in (8)–(9) depend on the first iterate of the output layer,  $\mathbf{V}^{(1)}$ , which can be decomposed into mean, bias and gradient noise components as

$$\mathbf{V}^{(1)} \approx \mathbf{Z}_{\text{out}} \left( \frac{1}{NL} \sum_{i=1}^N (\mathbf{x}_{i,1} - \frac{1}{V} \mathbb{1}_V) (\mathbf{X}_i \mathbb{1}_L)^\top \right) \mathbf{Z}_{\text{in}}^\top \quad (11)$$

$$\approx \mathbf{Z}_{\text{out}} \left( \underbrace{\frac{1}{VL} (\mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top)}_{\text{Mean}} + \underbrace{\frac{1}{VN} \sum_{i=1}^N (\mathbf{x}_{i,1} - \frac{1}{V} \mathbb{1}_V) \mathbb{1}_V^\top}_{\text{Bias}} + \underbrace{\frac{1}{\sqrt{LVN}} \Xi}_{\text{Gradient noise}} \right) \mathbf{Z}_{\text{in}}^\top \quad (12)$$

where the gradient noise component is given by

$$\Xi := \sqrt{\frac{V}{LN}} \left( \sum_{i=1}^N (\mathbf{x}_{i,1} - \frac{1}{V} \mathbb{1}_V) (\mathbf{X}_i \mathbb{1}_L - \frac{L}{V} \mathbb{1}_V)^\top - \frac{1}{V} (\mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top) \right).$$

Here, the bias term arises from aggregating tokens at initialization: the aggregate-token averages  $\frac{1}{L} \mathbf{X}_i \mathbb{1}_L$  concentrate around their mean  $\frac{1}{V} \mathbb{1}_V$  as  $L$  grows, so this effect appears as the bias term. The gradient-noise term captures finite-sample fluctuations of tokens around this mean. We explicitly factor out the typical size  $1/\sqrt{VLN}$  in (12) so that the remaining matrix  $\Xi$  stays of constant size on average, i.e.,  $\mathbb{E}[\|\Xi\|_2^2] = O(1)$ . We are now ready to consider the cases listed above.

**Informative term.** With the decomposition in (12), the informative term in (9) can be written as the sum of two contributions:

$$\text{Informative} \approx \frac{1}{VL^2} \frac{1}{N} \sum_{i=1}^N \mathbf{x}_{i,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top) \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_{i,1} - \frac{1}{V} \mathbb{1}_V) \underbrace{= O(1)}$$

The first term is due to the mean component; the second term is due to the gradient-noise  $\Xi$  component in (12). The bias-related terms are ignored, as they do not contribute. By standard concentration arguments for Gaussian matrices, the first term remains  $O(1)$ , whereas the second term concentrates within  $\pm(\log V)/d$ , yielding the *Signal* component (6) (the noise component is weaker than the remaining terms in (6)). See the “*Concentration bound for score<sub>12</sub>*” part for the formal proof.

**Non-informative term.** In this part, we consider large  $L$  regime where

$$\frac{1}{L} \mathbf{Z}_{\text{in}} \mathbf{X}_i \mathbf{X}_i^\top \mathbf{Z}_{\text{in}}^\top \approx \frac{1}{d} \mathbf{I}_d. \quad (13)$$

We consider an arbitrary row of  $\mathbf{X}_{\text{in}}$ , which we denote with  $\mathbf{x}_{\text{in}}$ . With this approximation, we can write the non-informative term as

$$\text{Non-informative} \approx \frac{1}{d\sqrt{LVN}} \underbrace{\mathbf{x}_{\text{in}}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \Xi \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_{i,1} - \frac{1}{V} \mathbb{1}_V)}_{\in \sqrt{\frac{V}{N}} \left[ -\frac{\log V}{d}, \frac{\log V}{d} \right]}$$

$$+ \frac{1}{Vd} \underbrace{\mathbf{x}_{\text{in}}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V}_{\in \sqrt{\frac{V}{d}} [-\log V, \log V]} \underbrace{\left\| \mathbf{Z}_{\text{out}} \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_{i,1} - \frac{1}{V} \mathbb{1}_V) \right\|_2^2}_{\approx \frac{1}{N}},$$

where we ignore terms depending on the mean component in (11), as they do not contribute. Here, the first term is due to the gradient-noise component  $\Xi$ ; by standard concentration arguments, this yields the scaling of the *Gradient noise* term. The second term arises from the bias term in (12); using standard concentration, this gives the scaling of the *Mean bias* term in (6). See the “*Concentration bound for score<sub>11</sub>*” part for the formal proof.

**MLP-noise in Attention-MLP.** We denote the rows of  $\mathbf{W}_{\text{in}}$  by  $\{\mathbf{w}_k\}_{k=1}^m$ , where  $\mathbf{w}_k \sim \mathcal{N}(0, \mathbf{I}_d)$ . We work in the large- $L$  regime for illustration, adopting the approximation in (13); however, the result extends to general  $L$ . Under this approximation, MLP-noise can be written as

$$\text{MLP-noise} \approx \mathbf{x}_{\text{in}}^\top \mathbf{Z}_{\text{in}}^\top \frac{1}{N^2 d} \sum_{i,j=1}^N \left( \frac{1}{m} \sum_{k=1}^m \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i \mathbb{1}_L \right) \phi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j \mathbb{1}_L \right) \right. \\ \left. \times \left( \mathbf{x}_{i,1} - \frac{1}{V} \mathbb{1}_V \right) \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} \left( \mathbf{x}_{j,1} - \frac{1}{V} \mathbb{1}_V \right) \right).$$

For large  $L$ , we have  $\left\| \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i \mathbb{1}_L \right\|_2 \approx L^{-1/2} \rightarrow 0$ , hence

$$\phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i \mathbb{1}_L \right) \phi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j \mathbb{1}_L \right) \rightarrow \underbrace{\phi(0)\phi'(0)}_{\text{nonzero constant}},$$

where Assumption 4 ensures  $\phi(0)\phi'(0) \neq 0$ . Replacing the  $\phi$ -dependent factors by this constant yields

$$\text{MLP noise} \approx \frac{\phi(0)\phi'(0)}{d} \frac{1}{m} \sum_{k=1}^m \underbrace{\mathbf{x}_{\text{input}}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k}_{\in \left[ -\frac{\log V}{\sqrt{m}}, \frac{\log V}{\sqrt{m}} \right]} \underbrace{\left\| \mathbf{Z}_{\text{out}} \frac{1}{N} \sum_{i=1}^N \left( \mathbf{x}_{i,1} - \frac{1}{V} \mathbb{1}_V \right) \right\|_2^2}_{\approx \frac{1}{N}},$$

which, by standard concentration arguments, gives the scaling of the *MLP noise* term in (6). See the “*Concentration bound for s<sub>3</sub>*” part for the formal proof.

## C.2 ATTENTION SCORES AND THEIR ASYMPTOTIC SCALING

**Assumption 4** (Technical conditions). *We work under the following conditions:*

- **Permutation.** Without loss of generality, assume  $\mathbf{\Pi} = \mathbf{I}_V$ .
- **Learning rates.** Take  $\eta = o_V(1)$ , chosen sufficiently small so that any  $o_\eta(1)$  terms are negligible; in particular, we may write  $\hat{\mathbf{p}}_1 = \frac{1}{V} \mathbb{1}_V + o_\eta(1)$ .
- **Activation.** We consider a polynomial activation  $\phi$  with a degree of  $p_\star$  satisfying:
  - $\phi(0), \phi'(0), \phi''(0) \neq 0$
  - The smallest non-zero Hermite component of  $\phi$  has index  $q_\star$ , i.e.  $q^\star := \min\{k > 0 \mid \mathbb{E}[\phi(Z)H_{e_k}] \neq 0\}$ , for  $Z \sim \mathcal{N}(0, 1)$ .

By using the technical condition above and ignoring the vanishing terms due to learning rate, we decompose the attention scores in to three terms  $\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3 \in \mathbb{R}^L$ :

$$\mathbf{X} \mathbf{Z}_{\text{in}}^\top \mathbf{W}_{\text{KQ}}^{(1)} \\ = \frac{\eta\gamma}{N^2 L^2} \mathbf{X} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \sum_{i,j=1}^N \alpha_{ij} \mathbf{X}_i^\top \left( \mathbf{I}_L - \frac{1}{L} \mathbb{1}_L \mathbb{1}_L^\top \right) \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L$$



$$\begin{aligned}
& \times (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \\
& + \frac{\eta\gamma}{N^2 L^2} \mathbf{X} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \sum_{i,j=1}^N \beta_{ij} \mathbf{X}_i^\top (\mathbf{I}_L - \frac{1}{L} \mathbb{1}_L \mathbb{1}_L^\top) \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \\
& \times (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \\
& + \frac{\eta\gamma}{N^2 L} \mathbf{X} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \sum_{i,j=1}^N \mathbf{X}_i^\top (\mathbf{I}_L - \frac{1}{L} \mathbb{1}_L \mathbb{1}_L^\top) \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \text{FW}(\mathbf{W}_{\text{in}}; \mathbf{Z}_{\text{in}}, \mathbf{X}_i, \mathbf{X}_j) \\
& \times (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \\
& =: \eta\gamma(\mathbf{s}_1 + \mathbf{s}_2 + \mathbf{s}_3).
\end{aligned}$$

The following theorem characterizes the scaling of each term:

**Theorem 3.** *With probability at least  $1 - o_V(1)$ , we have the following:*

$$\begin{aligned}
\mathbf{e}_l^\top \mathbf{s}_1 & \asymp \mathbb{1}_{l=1} \left( \frac{1}{VL^2} \pm \frac{1}{\sqrt{NV}L^{3/2}d} \right) \pm \frac{1}{N\sqrt{L}d(d \wedge L^2)^{1/2}(d \wedge L)^{1/2}} \\
\mathbf{e}_l^\top \mathbf{s}_2 & \asymp \pm \left( \frac{1}{N\sqrt{L}d(L \wedge d)} + \frac{1}{NLd(L \wedge d)^{1/2}} \right) \\
\mathbf{e}_l^\top \mathbf{s}_3 & \asymp \frac{\pm 1}{Nd\sqrt{m}}.
\end{aligned}$$

Moreover, for notational convenience, we define

$$\begin{aligned}
\mathbf{A}_{1,ir} & := \mathbf{Z}_{\text{in}} \left( \frac{1}{LN} \sum_{j=1}^N \alpha_{ij} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\
& \times \left( \frac{1}{LN} \sum_{j=1}^N \alpha_{rj} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \mathbf{Z}_{\text{in}}^\top \quad (14)
\end{aligned}$$

$$\begin{aligned}
\mathbf{A}_{2,ir} & := \mathbf{Z}_{\text{in}} \left( \frac{1}{LN} \sum_{j=1}^N \alpha_{ij} (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\
& \times \left( \frac{1}{LN} \sum_{j=1}^N \alpha_{rj} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right) \mathbf{Z}_{\text{in}}^\top \quad (15)
\end{aligned}$$

$$\mathbf{A}_{3,ir} := \frac{1}{L^2 V^2} \left( \frac{1}{N} \sum_{j=1}^N \alpha_{ij} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right)^\top \left( \frac{1}{N} \sum_{j=1}^N \alpha_{rj} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right) \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top$$

and

$$\mathbf{S}_1 := \left( \frac{1}{LN} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \left( \frac{1}{LN} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \quad (16)$$

$$\begin{aligned}
\mathbf{S}_2 & := \left( \frac{1}{LN} \sum_{j=1}^N (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\
& \times \left( \frac{1}{LN} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right) \quad (17)
\end{aligned}$$

$$\mathbf{S}_3 := \frac{1}{L^2 V^2} \left( \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right)^\top \left( \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right) \mathbb{1}_V \mathbb{1}_V^\top. \quad (18)$$

We first make an observation that we will frequently rely on in the following:

**Proposition 1.** For any  $p \in \mathbb{N}$ , we have

$$\mathbb{E}[\|\mathbf{A}_{1,ir}\|_2^2] \vee \mathbb{E}[\|\mathbf{A}_{2,ir}\|_2^2] \vee \mathbb{E}[\|\mathbf{A}_{3,ir}\|_2^2] \leq \text{poly}_{p,p_*}(d, V, L).$$

where  $\text{poly}_{p,p_*}(N, d, V, L)$  denotes a polynomial function of  $(d, V, L)$  whose degree depends on  $(p, p_*)$ .

*Proof.* By Proposition 12, we observe that  $\alpha_{ij} \leq \text{poly}_{p_*}(d, V, L)$ . Therefore, we have

$$\|\mathbf{A}_{1,ir}\|_2 \vee \|\mathbf{A}_{2,ir}\|_2 \vee \|\mathbf{A}_{3,ir}\|_2 \leq \text{poly}_{p_*}(d, V, L) \|\mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top\|_2.$$

from which the result follows.  $\square$

### C.3 PROOF OF THEOREM 3

We observe that

$$\begin{aligned} \mathbf{X} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top &= \left( \mathbf{Z}_{\text{in}} \mathbf{X}^\top + \delta \mathbf{Z}_{\text{in}} \mathbf{e}_{V+1} \mathbf{e}_1^\top \right)^\top \left( \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top + \delta \mathbf{Z}_{\text{in}} \mathbf{e}_{V+1} \mathbf{e}_1^\top \right) \\ &= \mathbf{X} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top + \delta \mathbf{e}_1 \mathbf{z}_{\text{trig}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top + \delta \mathbf{X} \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\text{trig}} \mathbf{e}_1^\top + \delta^2 \|\mathbf{z}_{\text{trig}}\|_2^2 \mathbf{e}_1 \mathbf{e}_1^\top. \end{aligned} \quad (19)$$

In the following, we will consider  $\mathbf{x}_l = \mathbf{e}_\nu$ , for  $\nu \in [V]$ . We will write

$$\delta \mathbf{e}_1^\top \mathbf{e}_l \mathbf{z}_{\text{trig}} + \mathbf{Z}_{\text{in}}^\top \mathbf{X}^\top \mathbf{e}_l = \mathbf{z}_\nu + \mathbb{1}_{k=1} \delta \mathbf{z}_{\text{trig}} =: \mathbf{z}_{\nu, \delta} \quad (20)$$

and

$$\delta \mathbf{e}_1 \mathbf{z}_{\text{trig}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}^\top \mathbf{e}_l + \delta^2 \|\mathbf{z}_{\text{trig}}\|_2^2 \mathbf{e}_1 \mathbf{e}_1^\top \mathbf{e}_l = \delta \mathbf{z}_{\nu, \delta}^\top \mathbf{z}_{\text{trig}} \mathbf{e}_1 =: \delta s_{\nu, \delta} \mathbf{e}_1. \quad (21)$$

In the following, we will consider the event.

$$\text{Event} := (E.1) \cap (E.2).$$

#### C.3.1 CONCENTRATION BOUND FOR $\mathbf{s}_1$

By (19)-(20)-(21), we can write that

$$\begin{aligned} \mathbf{e}_l^\top \mathbf{s}_1 &= \frac{1}{N^2 L^2} \sum_{i,j=1}^N \alpha_{ij} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \left( \mathbf{I}_L - \frac{1}{L} \mathbb{1}_L \mathbb{1}_L^\top \right) \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \\ &\quad + \frac{\delta s_{\nu, \delta}}{N^2 L^2} \sum_{i,j=1}^N \alpha_{ij} (\mathbf{e}_1 - \frac{1}{L} \mathbb{1}_L)^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \\ &=: \text{score}_{11} + \text{score}_{12}. \end{aligned}$$

We will analyze  $\text{score}_{11}$  and  $\text{score}_{12}$  separately. We define

$$\mathbf{V}_i := \mathbf{V}_{i,1} + \mathbf{V}_{i,2} + \mathbf{V}_{i,3},$$

where

$$\begin{aligned} \mathbf{V}_{i,1} &:= \left( \frac{1}{NL} \sum_{j=1}^N \alpha_{ij} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\ \mathbf{V}_{i,2} &:= \left( \frac{1}{NL} \sum_{j=1}^N \alpha_{ij} (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\ \mathbf{V}_{i,3} &:= \frac{1}{V} \mathbb{1}_V \left( \frac{1}{NL} \sum_{j=1}^N \alpha_{ij} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right). \end{aligned}$$

**Concentration bound for  $\text{score}_{11}$ :** We start with  $\text{score}_{11}$ . We define

$$\begin{aligned} C_i &:= \frac{1}{L}(\mathbf{x}_i - \frac{1}{V}\mathbb{1}_V)\mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top (\mathbf{I}_L - \frac{1}{L}\mathbb{1}_L \mathbb{1}_L^\top) \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \\ &= \underbrace{\frac{1}{L}(\mathbf{x}_i - \frac{1}{V}\mathbb{1}_V)\mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}}}_{:=C_{i,1}} - \underbrace{\frac{1}{L^2}(\mathbf{x}_i - \frac{1}{V}\mathbb{1}_V)\mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}}}_{:=C_{i,2}} \end{aligned}$$

By Chebyshev's inequality, with probability  $1 - o_V(1)$ ,

$$\begin{aligned} \text{score}_{11} &= \frac{1}{N} \sum_{i=1}^N \text{tr}(\mathbf{V}_i \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} C_i) = \text{tr}\left(\mathbf{Z}_{\text{out}} \frac{1}{N} \sum_{i=1}^N C_i \mathbf{V}_i \mathbf{Z}_{\text{out}}^\top\right) \\ &= \underbrace{\text{tr}\left(\frac{1}{N} \sum_{i=1}^N C_i \mathbf{V}_i\right)}_{\text{score}_{111}} \pm \underbrace{\frac{1}{\sqrt{d}} \left\| \frac{\log V}{N} \sum_{i=1}^N C_i \mathbf{V}_i \right\|_F}_{\text{score}_{112}}. \end{aligned}$$

We start with bounding  $\text{score}_{112}$  term. We have

$$\text{score}_{112} \leq \frac{1}{\sqrt{d}} \left\| \frac{1}{N} \sum_{i=1}^N C_i \mathbf{V}_{i,1} \right\|_F + \frac{1}{\sqrt{d}} \left\| \frac{1}{N} \sum_{i=1}^N C_i \mathbf{V}_{i,2} \right\|_F + \frac{1}{\sqrt{d}} \left\| \frac{1}{N} \sum_{i=1}^N C_i \mathbf{V}_{i,3} \right\|_F$$

We have for  $u \in \{1, 2\}$

$$\left\| \frac{1}{N} \sum_{i=1}^N C_i \mathbf{V}_{i,u} \right\|_F^2 \leq \frac{2}{N^2} \left( \sum_{i,r=1}^N \text{tr}(C_{i,1} \mathbf{V}_{i,u} \mathbf{V}_{r,u}^\top C_{r,1}^\top) + \text{tr}(C_{i,2} \mathbf{V}_{i,u} \mathbf{V}_{r,u}^\top C_{r,2}^\top) \right).$$

We define

$$t_1 := \frac{\phi'(0)^4}{dL^2} \left( \frac{1}{N} + \left(1 - \frac{1}{V}\right) \frac{1}{V} \right), \quad t_2 := \frac{\phi'(0)^4}{d} \left(1 - \frac{1}{V}\right)^2 \frac{L-1}{L^2 N},$$

For  $i \neq r$ , by using the definition in  $\mathbf{A}_{1,ir}$  and  $\mathbf{A}_{2,ir}$  in (14)-(15), we have

$$\begin{aligned} \text{tr}(C_{i,1} \mathbf{V}_{i,u} \mathbf{V}_{r,u}^\top C_{r,1}^\top) &= \frac{1}{L^2} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{A}_{u,ir} \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\ &= \frac{t_u}{L^2} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\ &\quad + \frac{1}{L^2} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top (\mathbf{A}_{u,ir} - t_u \mathbf{I}_d) \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\ &\leq \frac{t_u}{L^2} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\ &\quad + \frac{1}{L^2} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \|\mathbf{A}_{u,ir} - t_u \mathbf{I}_d\|_2 \|\mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2 \|\mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2 \end{aligned}$$

By Proposition 9, we have

$$\frac{t_u}{L^2} \mathbb{E} \left[ (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \mid \mathbf{Z}_{\text{in}} \right] \leq \frac{C t_u}{L^2} \frac{1}{V d}. \quad (22)$$

Moreover, by using Event and Propositions 1 and 9, we have

$$\begin{aligned} &\frac{1}{L^2} \mathbb{E} \left[ (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \|\mathbf{A}_{u,ir} - t_u \mathbf{I}_d\|_2 \|\mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2 \|\mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2 \mid \mathbf{Z}_{\text{in}} \right] \\ &\leq \frac{C}{V d (L \wedge d)} \begin{cases} \phi'(0)^2 \left( \frac{1}{N d L^3} + \frac{1}{V d L^2} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}} \right) + \phi'(0)^4 \left( \frac{\log V}{L^2 V^{3/2} \sqrt{d}} + \frac{\log^2 V}{L^2 N \sqrt{V d}} \right), & u = 1 \\ \frac{\sqrt{V}}{d \sqrt{N L}} \left( \frac{1}{N L^{\frac{3}{2}}} + \frac{1}{V \sqrt{L}} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}} \right) + \phi'(0)^4 \left( \frac{\log V}{N L \sqrt{V d}} + \frac{\log^3 V}{N \sqrt{L V d}} \right), & u = 2 \end{cases} \\ &\leq \frac{C}{N^{3/2} \sqrt{V} d^2 L^2} \frac{1}{L \wedge d} + \frac{C}{V^{3/2} \sqrt{N} L d^2} \frac{1}{L \wedge d} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}} + \frac{C \log^3 V}{N V^{3/2} L^{1/2} d^{3/2}} \frac{1}{(L \wedge d)^{3/2}}. \end{aligned} \quad (23)$$

On the other hand,

$$\begin{aligned}
& \text{tr}(\mathbf{C}_{i,2} \mathbf{V}_{i,u} \mathbf{V}_{r,u}^\top \mathbf{C}_{r,2}^\top) \\
&= \frac{1}{L^4} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{A}_{u,ir} \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\
&= \frac{t_u}{L^4} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\
&+ \frac{1}{L^4} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top (\mathbf{A}_{u,ir} - t_u \mathbf{I}_d) \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\
&\leq \frac{t_u}{L^4} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\
&+ \frac{1}{L^4} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \|\mathbf{A}_{u,ir} - t_u \mathbf{I}_d\|_2 \|\mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2 \|\mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2.
\end{aligned}$$

By using Event,

$$\frac{t_u}{L^4} \mathbb{E} \left[ (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} | \mathbf{Z}_{\text{in}} \right] \leq \frac{C t_u}{V L} \frac{\log^2 V}{V \wedge L^2 \wedge L \sqrt{d}} \frac{1}{L \wedge d}.$$

Moreover, by using Event,

$$\begin{aligned}
& \frac{1}{L^4} \mathbb{E} \left[ (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \|\mathbf{A}_{u,ir} - t_u \mathbf{I}_d\|_2 \|\mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2 \|\mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2 | \mathbf{Z}_{\text{in}} \right] \\
&\leq \frac{C}{V L^2 (L \wedge d)} \left\{ \phi'(0)^2 \left( \frac{1}{N d L^3} + \frac{1}{V d L^2} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}} \right) + \phi'(0)^4 \left( \frac{\log V}{L^2 V^{3/2} \sqrt{d}} + \frac{\log^2 V}{L^2 N \sqrt{V d}} \right), \quad u=1 \right. \\
&\quad \left. \frac{\sqrt{V}}{d \sqrt{N L}} \left( \frac{1}{N L^{\frac{3}{2}}} + \frac{1}{V \sqrt{L}} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}} \right) + \phi'(0)^4 \left( \frac{\log V}{N L \sqrt{V d}} + \frac{\log^3 V}{N \sqrt{L V d}} \right), \quad u=2 \right\} \\
&\leq \frac{C}{N^{3/2} \sqrt{V} d L^4} \frac{1}{L \wedge d} + \frac{C}{V^{3/2} \sqrt{N} L^3 d} \frac{1}{L \wedge d} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}} + \frac{C \log^3 V}{N V^{3/2} L^{5/2} \sqrt{d}} \frac{1}{(L \wedge d)^{3/2}}. \quad (24)
\end{aligned}$$

On the other hand, for  $i = r$ , we have

$$\begin{aligned}
& \text{tr}(\mathbf{C}_{i,1} \mathbf{V}_{i,u} \mathbf{V}_{i,u}^\top \mathbf{C}_{i,1}^\top) + \text{tr}(\mathbf{C}_{i,2} \mathbf{V}_{i,u} \mathbf{V}_{i,u}^\top \mathbf{C}_{i,2}^\top) \\
&= \frac{1}{L^2} (1 - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{A}_{u,ii} \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\
&+ \frac{(1 - \frac{1}{V})}{L^4} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{A}_{u,ii} \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\
&\leq \frac{t_u}{L^2} (1 - \frac{1}{V}) \|\mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2^2 + \frac{t_u}{L^4} \|\mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2^2.
\end{aligned}$$

By using Event and Proposition 9, we have

$$\begin{aligned}
& \frac{t_u}{L^2} \mathbb{E} \left[ \|\mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2^2 | \mathbf{Z}_{\text{in}} \right] + \frac{t_u}{L^4} \mathbb{E} \left[ \|\mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2^2 | \mathbf{Z}_{\text{in}} \right] \\
&\leq \frac{C t_u}{L^2} \left( \frac{L}{d} + \frac{L^2}{d^2} \right) + \frac{C t_u}{L^2} \frac{1}{L \wedge d}. \quad (25)
\end{aligned}$$

Therefore, we have by (22)-(23)-(24)-(25) and using  $N \ll V L$  and  $L \ll V$ , we have

$$\begin{aligned}
& \mathbb{E} \left[ \left\| \frac{1}{N} \sum_{i=1}^N \mathbf{C}_i \mathbf{V}_{i,1} \right\|_F^2 | \mathbf{Z}_{\text{in}} \right] + \mathbb{E} \left[ \left\| \frac{1}{N} \sum_{i=1}^N \mathbf{C}_i \mathbf{V}_{i,2} \right\|_F^2 | \mathbf{Z}_{\text{in}} \right] \\
&\leq \frac{C}{N^2 d L (d \wedge L^2) (d \wedge L)} + \frac{C}{N^{3/2} \sqrt{V} d L^2 (d \wedge L^2) (L \wedge d)} \\
&+ \frac{C}{V^{3/2} \sqrt{N} L d (d \wedge L^2) (L \wedge d)} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}} + \frac{C \log^3 V}{N V^{3/2} \sqrt{L} d (d \wedge L^2) (L \wedge d)^{3/2}} \\
&\leq \frac{C}{N^2 d L (d \wedge L^2) (d \wedge L)} + \frac{C \log^3 V}{N V^{3/2} \sqrt{L} d (d \wedge L^2) (L \wedge d)^{3/2}}. \quad (26)
\end{aligned}$$

On the other hand, we have

$$\left\| \frac{1}{N} \sum_{i=1}^N \mathbf{C}_i \mathbf{V}_{i,3} \right\|_F^2 \leq \frac{2}{N^2} \left( \sum_{i,r=1}^N \text{tr}(\mathbf{C}_{i,1} \mathbf{V}_{i,3} \mathbf{V}_{r,3}^\top \mathbf{C}_{r,1}^\top) + \text{tr}(\mathbf{C}_{i,2} \mathbf{V}_{i,3} \mathbf{V}_{r,3}^\top \mathbf{C}_{r,2}^\top) \right).$$

We define  $t_3 := \frac{\phi'(0)^4}{NV^2L^2}$  and

$$\tilde{\Delta}_{3,ir} := \frac{1}{L^2V^2} \left( \frac{1}{N} \sum_{j=1}^N \alpha_{ij} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right)^\top \left( \frac{1}{N} \sum_{j=1}^N \alpha_{rj} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right) - \frac{\phi'(0)^4}{NV^2L^2}.$$

We have for  $i \neq r$ ,

$$\begin{aligned} & \text{tr}(\mathbf{C}_{i,1} \mathbf{V}_{i,3} \mathbf{V}_{r,3}^\top \mathbf{C}_{r,1}^\top) + \text{tr}(\mathbf{C}_{i,2} \mathbf{V}_{i,3} \mathbf{V}_{r,3}^\top \mathbf{C}_{r,2}^\top) \\ &= \frac{1}{L^2} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{A}_{3,ir} \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\ &+ \frac{1}{L^4} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{A}_{3,ir} \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\ &\leq \frac{t_3}{L^2} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\ &+ \frac{t_3}{L^4} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\ &+ \frac{\tilde{\Delta}_{3,ir}}{L^2} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\ &+ \frac{\tilde{\Delta}_{3,ir}}{L^4} (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \end{aligned}$$

For the first term, by Proposition 9,

$$\frac{t_3}{L^2} \mathbb{E} \left[ (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} | \mathbf{Z}_{\text{in}} \right] \leq \frac{C\phi'(0)^4 \log^2 V}{NV^2L^4 d^2}.$$

For the second term, by using Event, we have

$$\begin{aligned} & \frac{t_3}{L^4} \mathbb{E} \left[ (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} | \mathbf{Z}_{\text{in}} \right] \\ & \leq \frac{\phi'(0)^4}{NV^3L^4} \frac{1}{d} \left( L \vee \frac{V}{d} \right) \end{aligned}$$

For the third item and fourth items, by using Event and Proposition 9,

$$\begin{aligned} & \frac{1}{L^2} \mathbb{E} \left[ (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) |\tilde{\Delta}_{3,ir}| \|\mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2 \|\mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2 | \mathbf{Z}_{\text{in}} \right] \\ & + \frac{1}{L^4} \mathbb{E} \left[ (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_r} - \frac{1}{V}) |\tilde{\Delta}_{3,ir}| \|\mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2 \|\mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2 | \mathbf{Z}_{\text{in}} \right] \\ & \leq \frac{C}{V^3L^4} \left( \frac{L}{d} + \frac{V}{d^2} + \frac{L^2}{d^2} \right) \left( \frac{\phi'(0)^4 \log^2 V}{N\sqrt{V}} + \frac{\phi'(0)^2}{N} \left( \frac{1}{NL} + \frac{1}{\sqrt{N} V \wedge L^2 \wedge L\sqrt{d}} \right) \right) \\ & + \frac{C}{V^3L^4} \left( \frac{L}{d} + \frac{V}{d^2} + \frac{L^2}{d^2} \right) \left( \frac{1}{NL} + \frac{1}{\sqrt{N} V \wedge L^2 \wedge L\sqrt{d}} \right)^2 \\ & \leq \frac{C}{NV^3L^2d(L \wedge d)} \frac{\phi'(0)^4 \log^2 V}{\sqrt{V} \wedge L^4 \wedge L^2d} + \frac{C}{NV^2L^4d^2} \frac{\phi'(0)^4 \log^2 V}{\sqrt{V} \wedge L^4 \wedge L^2d}. \end{aligned} \tag{27}$$

For  $i = r$ , by using Event, we have

$$\begin{aligned} & \text{tr}(\mathbf{C}_{i,1} \mathbf{V}_{i,3} \mathbf{V}_{i,3}^\top \mathbf{C}_{i,1}^\top) + \text{tr}(\mathbf{C}_{i,2} \mathbf{V}_{i,3} \mathbf{V}_{i,3}^\top \mathbf{C}_{i,2}^\top) \\ &= \frac{1}{L^2} (1 - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{A}_{3,ir} \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\ &+ \frac{1}{L^4} (1 - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{A}_{3,ir} \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \end{aligned}$$



$$\leq \frac{2t_3}{L^2} |\mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta}|^2 + \frac{2t_3}{L^4} |\mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta}|^2.$$

Then, by using Event and Proposition 9,

$$\mathbb{E} \left[ \text{tr}(\mathbf{C}_{i,1} \mathbf{V}_{i,3} \mathbf{V}_{i,3}^\top \mathbf{C}_{i,1}^\top) + \text{tr}(\mathbf{C}_{i,2} \mathbf{V}_{i,3} \mathbf{V}_{i,3}^\top \mathbf{C}_{i,2}^\top) | \mathbf{Z}_{\text{in}} \right] \leq \frac{C \phi'(0)^4 \log^2 V}{N V d^2 L^3} \left(1 + \frac{L}{d}\right). \quad (28)$$

Therefore, by using (27)-(28) and using  $L \ll V$  and  $N \ll VL$ , we have

$$\mathbb{E} \left[ \left\| \frac{1}{N} \sum_{i=1}^N \mathbf{C}_i \mathbf{V}_{i,3} \right\|_F^2 | \mathbf{Z}_{\text{in}} \right] \ll \frac{1}{N^2 d L (d \wedge L^2) (d \wedge L)}. \quad (29)$$

Therefore, by (26)-(29), we have

$$\text{score}_{112} \leq \frac{C \log V}{N \sqrt{L} d (d \wedge L^2)^{1/2} (d \wedge L)^{1/2}} + \frac{C \log^{5/2} V}{\sqrt{N} (V d)^{3/4} L^{1/4} (d \wedge L^2)^{1/2} (L \wedge d)^{3/4}}.$$

For  $\text{score}_{111}$ , we write

$$\text{score}_{111} = \underbrace{\text{tr} \left( \frac{1}{N} \sum_{i=1}^N \mathbf{C}_i \mathbf{V}_{i,1} \right)}_{:= \text{score}_{1111}} + \underbrace{\text{tr} \left( \frac{1}{N} \sum_{i=1}^N \mathbf{C}_i \mathbf{V}_{i,2} \right)}_{:= \text{score}_{1112}} + \underbrace{\text{tr} \left( \frac{1}{N} \sum_{i=1}^N \mathbf{C}_i \mathbf{V}_{i,3} \right)}_{:= \text{score}_{1113}}.$$

We have

$$\begin{aligned} \text{score}_{1111} &= \frac{1}{N^2 L^2} \sum_{i,j=1}^N \alpha_{ij} (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \left( \mathbf{I}_L - \frac{1}{L} \mathbb{1}_L \mathbb{1}_L^\top \right) \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L \\ &= \frac{\phi'(0)^2}{N^2 L^2} \sum_{j=1}^N \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \sum_{\substack{i=1 \\ i \neq j}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V} \mathbb{1}_V) \mathbf{X}_i^\top \left( \mathbf{I}_L - \frac{1}{L} \mathbb{1}_L \mathbb{1}_L^\top \right) \mathbf{X}_i \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L \\ &\quad + \frac{(1 - \frac{1}{V})}{N^2 L^2} \sum_{j=1}^N \alpha_{jj} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \left( \mathbf{I}_L - \frac{1}{L} \mathbb{1}_L \mathbb{1}_L^\top \right) \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L \\ &\quad + \frac{1}{N^2 L^2} \sum_{j=1}^N \sum_{\substack{i=1 \\ i \neq j}}^N (\alpha_{ij} - \phi'(0)^2) (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \left( \mathbf{I}_L - \frac{1}{L} \mathbb{1}_L \mathbb{1}_L^\top \right) \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L \\ &=: \text{score}_{11111} + \text{score}_{11112} + \text{score}_{11113}. \end{aligned}$$

We start with the last term. By using Cauchy-Schwartz inequality,

$$\begin{aligned} |\text{score}_{11113}| &\leq \left( \frac{1}{N^2 L^2} \sum_{j=1}^N \sum_{\substack{i=1 \\ i \neq j}}^N |\alpha_{ij} - \phi'(0)^2| |\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}| \right) \\ &\quad \times \sup_{i \neq j \in [N]} |\mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \left( \mathbf{I}_L - \frac{1}{L} \mathbb{1}_L \mathbb{1}_L^\top \right) \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L| \\ &\leq \left( \frac{1}{N^2 L^2} \sum_{j=1}^N \sum_{\substack{i=1 \\ i \neq j}}^N |\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}| \right) \sup_{i \neq j \in [N]} |\alpha_{ij} - \phi'(0)^2| \\ &\quad \times \sup_{i \neq j \in [N]} |\mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \left( \mathbf{I}_L - \frac{1}{L} \mathbb{1}_L \mathbb{1}_L^\top \right) \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L| \end{aligned}$$

$$\leq \frac{C \log V}{V \sqrt{Ld(L \wedge d)}} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}}. \quad (30)$$

where we used Event in (30). Next, we consider  $\text{score}_{11112}$ :

$$\begin{aligned} & |\text{score}_{11112}| \\ &= \frac{(1 - \frac{1}{V})}{N^2 L^2} \sum_{j=1}^N \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \alpha_{jj} \mathbf{X}_j^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L - \mathbb{E} [\alpha_{jj} \mathbf{X}_j^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L | \mathbf{Z}_{\text{in}}] \right) \\ &+ \frac{(1 - \frac{1}{V})}{N L^2} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} [\alpha_{11} \mathbf{X}_1^\top \mathbf{X}_1 \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_1^\top \mathbb{1}_L | \mathbf{Z}_{\text{in}}] \\ &- \frac{(1 - \frac{1}{V})}{N^2 L V} \sum_{j=1}^N \alpha_{jj} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \\ &- \frac{(1 - \frac{1}{V})}{N^2 L^3} \sum_{j=1}^N \alpha_{jj} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L \end{aligned}$$

By using Event, we have

- For the first summand,

$$\begin{aligned} & \mathbb{E} \left[ \left( \frac{(1 - \frac{1}{V})}{N^2 L^2} \sum_{j=1}^N \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \alpha_{jj} \mathbf{X}_j^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L - \mathbb{E} [\alpha_{jj} \mathbf{X}_j^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L | \mathbf{Z}_{\text{in}}] \right) \right)^2 | \mathbf{Z}_{\text{in}} \right] \\ & \leq \frac{(1 - \frac{1}{V})^2}{N^3 L^4} \mathbb{E} \left[ \alpha_{jj}^2 \left( \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_1^\top \mathbf{X}_1 \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_1^\top \mathbb{1}_L \right)^2 | \mathbf{Z}_{\text{in}} \right] \\ & \leq \frac{C \phi'(0)^4}{N^3 L^3} \mathbb{E} \left[ \left\| \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_1^\top \mathbf{X}_1 \mathbf{Z}_{\text{in}}^\top \right\|_2^2 | \mathbf{Z}_{\text{in}} \right] + o_V(1) \\ & \leq \frac{C \phi'(0)^4}{N^3 L d (L \wedge d)} \end{aligned}$$

By Chebyshev's inequality with probability  $1 - o_V(1)$ , we have

$$\begin{aligned} & \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \frac{(1 - \frac{1}{V})}{N^2 L^2} \sum_{j=1}^N \left( \alpha_{jj} \mathbf{X}_j^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L - \mathbb{E} [\alpha_{jj} \mathbf{X}_j^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L | \mathbf{Z}_{\text{in}}] \right) \\ & \leq \frac{C \phi'(0)^2 \log V}{N^{\frac{3}{2}} \sqrt{Ld} \sqrt{L \wedge d}}. \end{aligned}$$

- For the second summand,

$$\begin{aligned} & \frac{(1 - \frac{1}{V})}{N L^2} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} [\alpha_{11} \mathbf{X}_1^\top \mathbf{X}_1 \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_1^\top \mathbb{1}_L | \mathbf{Z}_{\text{in}}] \\ &= \frac{(1 - \frac{1}{V}) \phi'(0)^2}{N L^2} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} [\mathbf{X}_1^\top \mathbf{X}_1 \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_1^\top \mathbf{X}_1 | \mathbf{Z}_{\text{in}}] \mathbb{1}_V \\ &+ \frac{(1 - \frac{1}{V})}{N L^2} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} [(\alpha_{11} - \phi'(0)^2) \mathbf{X}_1^\top \mathbf{X}_1 \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_1^\top \mathbf{X}_1 | \mathbf{Z}_{\text{in}}] \mathbb{1}_V \\ &= \frac{(1 - \frac{1}{V}) \phi'(0)^2}{N L} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} [\mathbf{x}_1 \mathbf{x}_1^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_1 \mathbf{x}_1^\top | \mathbf{Z}_{\text{in}}] \mathbb{1}_V \\ &+ \frac{(1 - \frac{1}{V}) \phi'(0)^2}{N V^2} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \\ &+ \frac{(1 - \frac{1}{V})}{N L V} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} [(\alpha_{11} - \phi'(0)^2) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_1^\top \mathbb{1}_L | \mathbf{Z}_{\text{in}}] \\ &+ \frac{(1 - \frac{1}{V})}{N L^2} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} [(\alpha_{11} - \phi'(0)^2) (\mathbf{X}_1^\top \mathbf{X}_1 - \frac{L}{V} \mathbf{I}_V) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_1^\top \mathbb{1}_L | \mathbf{Z}_{\text{in}}] \end{aligned}$$

$$\leq C \log V \left( \frac{1}{NL\sqrt{V}d} + \frac{1}{N\sqrt{V}d^{3/2}} + \frac{1}{NL^{3/2}d} \right).$$

• For the third summand,

$$\begin{aligned} & \frac{1}{N^2LV} \sum_{j=1}^N \alpha_{jj} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \\ &= \frac{\phi'(0)^2}{NV^2} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbb{1}_V + \frac{\phi'(0)^2}{N^2LV} \sum_{j=1}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_j^\top \mathbf{X}_j - \frac{L}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \\ &+ \frac{1}{N^2LV} \sum_{j=1}^N (\alpha_{jj} - \phi'(0)^2) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V. \end{aligned}$$

The first term:

$$\left| \frac{\phi'(0)^2}{NV^2} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbb{1}_V \right| \leq \frac{C \log V}{N\sqrt{V}d^{\frac{3}{2}}}.$$

The second term: By using

$$\begin{aligned} & \mathbb{E} \left[ \left( \sum_{j=1}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_j^\top \mathbf{X}_j - \frac{L}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right)^2 \middle| \mathbf{Z}_{\text{in}} \right] \\ &= \sum_{j=1}^N \mathbb{E} \left[ \left( \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_j^\top \mathbf{X}_j - \frac{L}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right)^2 \middle| \mathbf{Z}_{\text{in}} \right] = \frac{LVN}{d^2}. \end{aligned}$$

Therefore, by Chebyshev's inequality, we have

$$\left| \frac{\phi'(0)^2}{N^2LV} \sum_{j=1}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_j^\top \mathbf{X}_j - \frac{L}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right| \leq \frac{\phi'(0)^2}{N^{3/2}\sqrt{V}Ld}.$$

Finally,

$$\begin{aligned} & \left| \frac{1}{N^2LV} \sum_{j=1}^N (\alpha_{jj} - \phi'(0)^2) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right| \\ & \leq \frac{1}{NV} \|\mathbf{z}_{\nu,\delta}\|_2 \left\| \frac{1}{NL} \sum_{j=1}^N (\alpha_{jj} - \phi'(0)^2) \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \right\|_2 \|\mathbf{Z}_{\text{in}} \mathbb{1}_V\|_2 \leq \frac{C}{N\sqrt{V}Ld}. \end{aligned}$$

Therefore,

$$\left| \frac{1}{N^2LV} \sum_{j=1}^N \alpha_{jj} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right| \leq C \log V \left( \frac{1}{N\sqrt{V}d^{\frac{3}{2}}} + \frac{1}{N\sqrt{V}Ld} \right).$$

• For the last summand,

$$\begin{aligned} & \frac{(1 - \frac{1}{V})}{N^2L^3} \sum_{j=1}^N \alpha_{jj} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right) \mathbb{1}_L \\ &= \frac{(1 - \frac{1}{V})}{N^2L^3} \sum_{j=1}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \\ & \times \left( \alpha_{jj} \mathbf{X}_j^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right) \mathbb{1}_L - \mathbb{E} \left[ \alpha_{jj} \mathbf{X}_j^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right) \mathbb{1}_L \middle| \mathbf{Z}_{\text{in}} \right] \right) \\ &+ \frac{(1 - \frac{1}{V})\phi'(0)^2}{N^2L^3} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} \left[ \mathbf{X}_1^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_1 \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right) \mathbb{1}_L \middle| \mathbf{Z}_{\text{in}} \right] \end{aligned}$$

$$+ \frac{(1 - \frac{1}{V})}{NL^3} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} [(\alpha_{11} - \phi'(0)^2) \mathbf{X}_1^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_1 \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L | \mathbf{Z}_{\text{in}}]. \quad (31)$$

We have

$$\begin{aligned} & \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} \left[ \left( \alpha_{jj} \mathbf{X}_1^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_1 \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L \right)^2 | \mathbf{Z}_{\text{in}} \right] \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \\ & \leq C \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} \left[ \mathbf{X}_1^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_1 \right] \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} + o_V(1) \\ & \leq C \left( \frac{L}{d} + \frac{L^2}{Vd} \right). \end{aligned}$$

Moreover, by using Proposition 6

$$\begin{aligned} & \mathbb{E} \left[ \mathbf{X}_1^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_1 \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L | \mathbf{Z}_{\text{in}} \right] \\ & = \mathbb{E} \left[ \mathbf{X}_1^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_1 \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_1^\top \mathbb{1}_L | \mathbf{Z}_{\text{in}} \right] - \frac{L}{V} \mathbb{E} \left[ \mathbf{X}_1^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_1 \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V | \mathbf{Z}_{\text{in}} \right] \\ & = L \mathbb{E} \left[ \mathbf{x}_1 \mathbf{x}_1^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_1 | \mathbf{Z}_{\text{in}} \right] + \left( \frac{L(L-1)}{V^2} \text{tr}(\mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}}) - \frac{2L(L-1)}{V^3} \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right) \mathbb{1}_V \\ & + \frac{L(L-2)}{V^2} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \end{aligned}$$

Lastly,

$$\begin{aligned} & \left| \frac{1}{NL^3} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} [(\alpha_{11} - \phi'(0)^2) \mathbf{X}_1^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_1 \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L | \mathbf{Z}_{\text{in}}] \right| \\ & \leq \frac{1}{NL^3} \mathbb{E} \left[ |\alpha_{11} - \phi'(0)^2| |\mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_1^\top \mathbb{1}_L| |\mathbb{1}_L^\top \mathbf{X}_1 \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L| | \mathbf{Z}_{\text{in}} \right] \\ & \leq \frac{C \log V}{NL^{5/2} \sqrt{d}}. \end{aligned}$$

Therefore, by Chebyshev's inequality, with probability  $1 - o_V(1)$ , we have

$$(31) \leq C \log V \left( \frac{1}{NL^2 \sqrt{L \wedge d}} + \frac{1}{NL \sqrt{Vd}} \right)$$

Therefore, we have

$$|\text{score}_{11112}| \leq C \log V \left( \frac{1}{N \sqrt{V} (L \wedge d) \sqrt{d}} + \frac{1}{NL^2 \sqrt{L \wedge d}} \right). \quad (32)$$

Finally, we consider  $\text{score}_{11111}$ :

$\text{score}_{11111}$

$$\begin{aligned} & = (1 - \frac{1}{L}) \frac{\phi'(0)^2}{NL} \sum_{j=1}^N \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{NL} \sum_{\substack{i=1 \\ i \neq j}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L \\ & - \frac{\phi'(0)^2}{NL^2} \sum_{j=1}^N \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{NL} \sum_{\substack{i=1 \\ i \neq j}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) (\mathbf{x}_i \mathbb{1}_L^\top \mathbf{X}_i + \mathbf{X}_i^\top \mathbb{1}_L \mathbf{x}_i^\top) \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L \\ & + \frac{\phi'(0)^2}{NL} \sum_{j=1}^N \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{NL} \sum_{\substack{i=1 \\ i \neq j}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{N}_i^\top \mathbf{N}_i \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \\ & + \frac{\phi'(0)^2}{NL} \sum_{j=1}^N \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{NL} \sum_{\substack{i=1 \\ i \neq j}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{N}_i^\top \mathbf{N}_i \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \end{aligned}$$

$$\begin{aligned}
& - \frac{\phi'(0)^2}{NL^2} \sum_{j=1}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{NL} \sum_{\substack{i=1 \\ i \neq j}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{N}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{N}_i \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \\
& =: \text{score}_{V1} + \text{score}_{V2} + \text{score}_{V3} + \text{score}_{V4} + \text{score}_{V5}.
\end{aligned}$$

For the first summand, we write

$$\begin{aligned}
\text{score}_{V1} & := (1 - \frac{1}{L}) \frac{\phi'(0)^2}{NL} \sum_{j=1}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \underbrace{\left( \frac{1}{NL} \sum_{\substack{i=1 \\ i \neq j}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)}_{:= \text{score}_{V11}} \\
& + (1 - \frac{1}{L}) \frac{\phi'(0)^2}{NL} \sum_{j=1}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \underbrace{\left( \frac{1}{NL} \sum_{\substack{i=1 \\ i \neq j}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1}}_{:= \text{score}_{V12}}.
\end{aligned}$$

We have by Event

$$\begin{aligned}
|\text{score}_{V11}| & \leq \left( \frac{\phi'(0)^2}{N^2 L^2} \sum_{j=1}^N \sum_{\substack{i=1 \\ i \neq j}}^N |\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}| \right) \sup_{i \neq j} \left| \mathbf{x}_i \mathbb{1}_{\mathbf{x}_i \neq \mathbf{e}_\nu} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right| \\
& \leq \left( \frac{\phi'(0)^2}{N^2 L^2} \sum_{j=1}^N \sum_{\substack{i=1 \\ i \neq j}}^N |\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}| \right) \sup_{i \neq j} \left| \mathbb{1}_{\mathbf{x}_i = \mathbf{e}_\nu} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right| \\
& \leq \frac{C \log V}{V L^2 \sqrt{d}}.
\end{aligned}$$

Moreover, let

$$\text{score}_{V12} =: (1 - \frac{1}{L}) \frac{\phi'(0)^2}{NL} \sum_{j=1}^N \text{score}_{V12,j}.$$

We have  $\mathbb{E}[\text{score}_{V12,j}] = 0$  and  $\mathbb{E}[\text{score}_{V12,j} \text{score}_{V12,j'}] = 0$  for  $j \neq j'$ , and

$$\begin{aligned}
& \mathbb{E}[\text{score}_{V12,j}^2] \\
& \leq \frac{CL}{d} \mathbb{E} \left[ \mathbb{1}_{\mathbf{x}_j = \mathbf{e}_k} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{NL} \sum_{\substack{i=1 \\ i \neq j}}^N \mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{NL} \sum_{\substack{i=1 \\ i \neq j}}^N \mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \right] \\
& + \frac{CL}{d} \mathbb{E} \left[ \mathbb{1}_{\mathbf{x}_j \neq \mathbf{e}_k} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{NL} \sum_{\substack{i=1 \\ i \neq j}}^N \mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{NL} \sum_{\substack{i=1 \\ i \neq j}}^N \mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \right] \\
& + \frac{CL}{dV^2} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{NL} \sum_{\substack{i=1 \\ i \neq j}}^N \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{NL} \sum_{\substack{i=1 \\ i \neq j}}^N \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\
& \leq \frac{C}{V L d^2}.
\end{aligned}$$

Therefore, by Chebyshev's inequality with probability  $1 - o_V(1)$ , we have

$$|\text{score}_{V12,j}| \leq \frac{C \log V}{\sqrt{N V L^3/2} d}.$$

Therefore,

$$|\text{score}_{V1}| \leq \frac{C \log V}{V L^2 \sqrt{d}} + \frac{C \log V}{\sqrt{N V L^3/2} d}.$$

Moreover, for the second term, we write

$$\begin{aligned}
& |\text{score}_{V2}| \\
& \leq \left( \frac{\phi'(0)^2}{N^2 L^3} \sum_{j=1}^N \sum_{\substack{i=1 \\ i \neq j}}^N |\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}| \right) \\
& \quad \times \sup_{i \neq j} |\mathbb{1}_{\mathbf{x}_i \neq \mathbf{e}_\nu} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} (\mathbf{x}_i \mathbb{1}_L^\top \mathbf{X}_i + \mathbf{X}_i^\top \mathbb{1}_L \mathbf{x}_i^\top) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L| \\
& \leq \left( \frac{\phi'(0)^2}{N^2 L^3} \sum_{j=1}^N \sum_{\substack{i=1 \\ i \neq j}}^N |\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}| \right) \\
& \quad \times \sup_{i \neq j} |\mathbb{1}_{\mathbf{x}_i \neq \mathbf{e}_\nu} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} (\mathbf{x}_i \mathbb{1}_L^\top \mathbf{X}_i + \mathbf{X}_i^\top \mathbb{1}_L \mathbf{x}_i^\top) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L| \\
& \leq \frac{C \log V}{V L^2 d}.
\end{aligned}$$

For the third term, we write

$$\begin{aligned}
\text{score}_{V3} &= \frac{\phi'(0)^2}{N L^2} \sum_{i=1}^N \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right) \\
&= \frac{\phi'(0)^2}{N L^2} \sum_{i=1}^N \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right) \\
&\quad + \frac{\phi'(0)^2 (L-1)}{N L^2 V} \sum_{i=1}^N \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right) \\
&=: \text{score}_{V31} + \text{score}_{V32}.
\end{aligned}$$

For  $\text{score}_{V32}$ , we note that

$$\begin{aligned}
& \mathbb{E} \left[ \left( \frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right) \left( \frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right)^\top \middle| \mathbf{x}_i \right] \\
& \leq C \left( \frac{1}{NV} + \frac{1}{V^2} \right) (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top + \frac{C}{NV^3} \mathbf{I}_V.
\end{aligned}$$

Therefore, we have

$$\mathbb{E} \left[ \left( \frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right) \left( \frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right)^\top \right] \preceq \frac{C}{V^3} \mathbf{I}_V.$$

We have

$$\mathbb{E}[\text{score}_{V32}^2] \leq \frac{C}{N^2 L^2 V^4 d} \sum_{i=1}^N \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \leq \frac{C}{N L^2 V^2 d^3}.$$

Therefore,

$$|\text{score}_{V32}| \leq \frac{C \log V}{\sqrt{NV} L d^{\frac{3}{2}}}$$

Moreover, we have

$$\mathbb{E}[\text{score}_{V32}^2]$$

$$\begin{aligned}
&\leq \frac{C}{N^2 V^2 L^4 d} \sum_{i=1}^N \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \right] \\
&\leq \frac{C}{N^2 V^2 L^4 d} \frac{L-1}{V} \sum_{i=1}^N \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \right] - \frac{C}{N^2 V^2 L^4 d} \frac{L-1}{V^2} \sum_{i=1}^N \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \right] \\
&\leq \frac{C}{N V^2 L^3 d^2}.
\end{aligned}$$

Therefore, by Chebyshev's inequality, we have

$$|\text{score}_{V32}| \leq \frac{C \log V}{\sqrt{N V L^{\frac{3}{2}} d}}$$

For the fourth term, we have  $\mathbb{E}[\text{score}_{V4}] = 0$  and

$$\begin{aligned}
\mathbb{E}[\text{score}_{V4}^2] &\leq \frac{C}{V N^4 L^4} \sum_{j=1}^N \sum_{\substack{i=1 \\ i \neq j}}^N \mathbb{E} \left[ \left( \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top \right) \mathbb{1}_{L-1} \right)^2 \right] \\
&\leq \frac{C}{V N^4 L^3 d} \sum_{j=1}^N \sum_{\substack{i=1 \\ i \neq j}}^N \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \right] \\
&\leq \frac{C}{V N^2 L d^2 (L \wedge d)}
\end{aligned}$$

Therefore, by Chebyshev's inequality with probability  $1 - o_V(1)$ , we have

$$|\text{score}_{V4}| \leq \frac{C \log V}{N \sqrt{V L d} \sqrt{L \wedge d}}$$

For the last term, we have  $\mathbb{E}[\text{score}_{V5}] = 0$  and

$$\begin{aligned}
\mathbb{E}[\text{score}_{V5}^2] &\leq \frac{C}{N^4 L^6 V} \sum_{j=1}^N \sum_{\substack{i=1 \\ i \neq j}}^N \mathbb{E} \left[ \left( \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top \right) \mathbb{1}_{L-1} \right)^2 \right] \\
&\leq \frac{C}{N^4 L^4 d V} \sum_{j=1}^N \sum_{\substack{i=1 \\ i \neq j}}^N \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbb{1}_{L-1} \mathbb{1}_{L-1}^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \right] \\
&\leq \frac{C}{N^2 L^4 d V} \left( \frac{L}{d} + \frac{L^2}{V d} \right)
\end{aligned}$$

Therefore, by Chebyshev's inequality with probability  $1 - o_V(1)$ , we have

$$|\text{score}_{V5}| \leq \frac{C \log V}{N L d \sqrt{V} (L \wedge \sqrt{V})}.$$

Overall, we have

$$|\text{score}_{11111}| \leq C \log V \left( \frac{1}{V L^2 \sqrt{d}} + \frac{1}{\sqrt{N V L^3/2} d} + \frac{1}{N \sqrt{V L d} \sqrt{L \wedge d}} + \frac{1}{N L d \sqrt{V} (L \wedge \sqrt{V})} \right). \quad (33)$$

Therefore, by (30)-(32)-(33) and using  $N \ll V L$  and  $L \ll V$ , we have

$$\begin{aligned}
|\text{score}_{11111}| &\leq C \log V \left( \frac{1}{V L^2 \sqrt{d}} + \frac{1}{\sqrt{N V L^3/2} d} + \frac{1}{N \sqrt{V L d} \sqrt{L \wedge d}} + \frac{1}{N L d \sqrt{V} (L \wedge \sqrt{V})} \right) \\
&\quad + C \log V \left( \frac{1}{N \sqrt{V} (L \wedge d) \sqrt{d}} + \frac{1}{N L^2 \sqrt{L \wedge d}} \right) + \frac{C \log V}{V \sqrt{L d} (L \wedge d)} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}}
\end{aligned}$$

$$\leq C \log V \left( \frac{1}{N\sqrt{V}(L \wedge d)\sqrt{d}} + \frac{1}{\sqrt{NV}L^{3/2}d} + \frac{1}{VL^{3/2}d^2} + \frac{1}{V^2\sqrt{L}d^{3/2}} \right) + \frac{C \log V}{VL^2(L \wedge d)^{1/2}} \quad (34)$$

Finally,

$$\begin{aligned} \text{score}_{1112} &= \frac{1}{N^2 L^2 V} \sum_{i=1}^N \sum_{j=1}^N \alpha_{ij} (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top (\mathbf{I}_L - \frac{1}{L} \mathbb{1}_L \mathbb{1}_L^\top) \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \\ &= \frac{1}{N^2 L^2 V} \sum_{i=1}^N \sum_{j=1}^N \alpha_{ij} (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_i^\top \mathbf{X}_i - \frac{L}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \\ &\quad + \frac{1}{N^2 L V^2} \sum_{i=1}^N \sum_{j=1}^N \alpha_{ij} (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \\ &\quad - \frac{1}{N^2 L^3 V} \sum_{i=1}^N \sum_{j=1}^N \alpha_{ij} (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \end{aligned}$$

By using Event and Proposition 9, we have

- For the third summand,

$$\frac{1}{L^2} \mathbb{E} \left[ \left( \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right)^2 | \mathbf{Z}_{\text{in}} \right] \leq \frac{C \log^4 V}{d} \left( L + \frac{V}{d} \right).$$

- For the first summand,

$$\begin{aligned} &\mathbb{E} \left[ \left( \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_i^\top \mathbf{X}_i - \frac{L}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right)^2 | \mathbf{Z}_{\text{in}} \right] \\ &= \mathbb{E} \left[ \left( \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right)^2 | \mathbf{Z}_{\text{in}} \right] - \frac{L^2}{V^2} \left( \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right)^2 = \frac{CVL \log^2 V}{d^2} \end{aligned}$$

Therefore

$$\begin{aligned} &\mathbb{E} \left[ \left( \frac{1}{N^2 L^2 V} \sum_{i=1}^N \sum_{j=1}^N \alpha_{ij} (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_i^\top \mathbf{X}_i - \frac{L}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right)^2 | \mathbf{Z}_{\text{in}} \right] \\ &\quad + \mathbb{E} \left[ \left( \frac{1}{N^2 L^3 V} \sum_{i=1}^N \sum_{j=1}^N \alpha_{ij} (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right)^2 | \mathbf{Z}_{\text{in}} \right] \\ &\leq \frac{C \log^2 V}{NL^3 V d^2}. \end{aligned}$$

- For the second summand, by Event,

$$\begin{aligned} &\frac{1}{N^2 L V^2} \sum_{i=1}^N \sum_{j=1}^N \alpha_{ij} (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \\ &\leq \frac{1}{N^2 L V^2} \left( \sum_{i=1}^N \sum_{j=1}^N |\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}| \right) \sup_{i,j} |\alpha_{ij} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V| \leq \frac{C \log V}{VL\sqrt{V}d^{3/2}}. \end{aligned}$$

Therefore, by Chebyshev's inequality with probability  $1 - o_V(1)$ , we have

$$|\text{score}_{1112}| \leq \frac{C \log^3 V}{NL^3 V d^2} + \frac{C \log^2 V}{VL\sqrt{V}d^{3/2}}. \quad (35)$$

By (34)-(35), overall we have

$$\text{score}_{111} \leq C \log V \left( \frac{1}{VL^2 \sqrt{d}} + \frac{1}{\sqrt{NV}L^{3/2}d} + \frac{1}{N\sqrt{V}Ld\sqrt{L \wedge d}} + \frac{1}{NLd\sqrt{V}(L \wedge \sqrt{V})} \right)$$



$$\begin{aligned}
& + C \log V \left( \frac{1}{N\sqrt{V}(L \wedge d)\sqrt{d}} + \frac{1}{N\sqrt{V}d^{\frac{3}{2}}} + \frac{1}{NL^2\sqrt{L \wedge d}} \right) + \frac{C \log V}{V\sqrt{Ld}(L \wedge d)} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \\
& \leq C \log V \left( \frac{1}{N\sqrt{V}(L \wedge d)\sqrt{d}} + \frac{1}{\sqrt{NV}L^{3/2}d} + \frac{1}{VL^{3/2}d(L \wedge d)} + \frac{1}{V^2\sqrt{Ld}(L \wedge d)} \right) \\
& + \frac{C \log V}{VL^2(L \wedge d)^{1/2}}.
\end{aligned}$$

**Concentration bound for  $\text{score}_{12}$ :** We recall that

$$\text{score}_{12} = \frac{\delta s_{\nu, \delta}}{N^2 L^2} \sum_{i,j=1}^N \alpha_{ij} \left( \mathbf{e}_1 - \frac{1}{L} \mathbb{1}_L \right)^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)$$

In this part, we will focus on the term

$$\begin{aligned}
& \frac{1}{N^2 L^2} \sum_{i,j=1}^N \alpha_{ij} \left( \mathbf{e}_1 - \frac{1}{L} \mathbb{1}_L \right)^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \\
& = \frac{1}{N^2 L^2} \sum_{i,j=1}^N \text{tr}(\mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} \alpha_{ij} (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{x}_i^\top) \\
& - \frac{1}{N^2 L^3} \sum_{i,j=1}^N \alpha_{ij} \text{tr}(\mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbb{1}_L^\top \mathbf{X}_i) \\
& = \text{score}_{121} + \text{score}_{122}
\end{aligned}$$

For the first term, we write

$$\begin{aligned}
\text{score}_{121} & = \frac{\phi'(0)^2}{N^2 L^2} \sum_{i,j=1}^N \text{tr}(\mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{x}_i^\top) \\
& + \frac{1}{N^2 L^2} \sum_{i,j=1}^N \text{tr}(\mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\alpha_{ij} - \phi'(0)^2) (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{x}_i^\top) \\
& = \text{score}_{1211} + \text{score}_{1212}
\end{aligned}$$

We start with the second term. By Chebyshev's inequality, we have

$$\begin{aligned}
\text{score}_{1212} & = \frac{1}{N^2 L^2} \sum_{i,j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \\
& \pm \frac{1}{N^2 L^2 \sqrt{d}} \left\| \sum_{i,j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right\|_F
\end{aligned}$$

By Event

$$\begin{aligned}
& \sum_{i,j=1}^N (\alpha_{ij} - \phi'(0)^2)^2 (\mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L)^2 \\
& \leq \frac{1}{(V \wedge L^2 \wedge L\sqrt{d})^2} \sum_{i \neq j=1}^N (\mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L)^2 + \frac{1}{L^2} \sum_{i=1}^N (\mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L)^2 \\
& \leq \left( \frac{N^2}{(V \wedge L^2 \wedge L\sqrt{d})^2} + \frac{N}{L^2} \right) \left( 1 + \frac{L}{d} \right)
\end{aligned}$$

Therefore,

$$\frac{1}{N^2 L^2 \sqrt{d}} \left\| \sum_{i,j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right\|_F$$

$$\begin{aligned}
&\leq \frac{1}{N^2 L^{3/2} \sqrt{d} (L \wedge d)^{1/2}} \left( \frac{N}{V \wedge L^2 \wedge L \sqrt{d}} + \frac{\sqrt{N}}{L} \right) \left\| \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2 \\
&\leq \frac{1}{N V L^{3/2} \sqrt{d} (L \wedge d)^{1/2}} \left( \frac{N}{V \wedge L^2 \wedge L \sqrt{d}} + \frac{\sqrt{N}}{L} \right) = \frac{o_V(1)}{V L^2}.
\end{aligned}$$

Moreover,

$$\begin{aligned}
&\left\| \frac{1}{N^2 L^2} \sum_{i,j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}) \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right\| \\
&\leq \frac{1}{N^2 L^2} \left( \sum_{i,j=1}^N |\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}| \right) \sup_{i,j \in [N]} |\alpha_{ij} - \phi'(0)^2| \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \leq \frac{(1 - \frac{1}{V})}{N L^2} \frac{1}{L} \left( 1 + \sqrt{\frac{L}{d}} \right).
\end{aligned}$$

Therefore,  $|\text{score}_{1212}| \ll \frac{1}{V L^2}$ . Next, we consider  $\text{score}_{122}$ :

$$\begin{aligned}
&\text{score}_{122} \\
&= \frac{\phi'(0)^2}{N^2 L^3} \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}) \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \\
&\quad \pm \frac{\phi'(0)^2}{\sqrt{d}} \left\| \frac{1}{N^2 L^3} \sum_{i,j=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right\|_F \\
&\quad + \frac{1}{N^2 L^3} \sum_{i,j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}) \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \\
&\quad \pm \frac{1}{\sqrt{d}} \left\| \frac{1}{N^2 L^3} \sum_{i,j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right\|_F \\
&=: \text{score}_{1221} + \text{score}_{1222} + \text{score}_{1223} + \text{score}_{1224}.
\end{aligned}$$

For  $\text{score}_{1224}$ , by Event,

$$\begin{aligned}
&\frac{1}{L^2} \sum_{i,j=1}^N (\alpha_{ij} - \phi'(0)^2)^2 (\mathbb{1}_L^\top \mathbf{X}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L)^2 \\
&\leq \frac{1}{(V \wedge L^2 \wedge L \sqrt{d})^2} \frac{1}{L^2} \sum_{i \neq j=1}^N (\mathbb{1}_L^\top \mathbf{X}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L)^2 + \frac{1}{L^2} \frac{1}{L^2} \sum_{i=1}^N (\mathbb{1}_L^\top \mathbf{X}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L)^2 \\
&\leq \frac{N^2}{(V \wedge L^2 \wedge L \sqrt{d})^3} + \frac{N}{L^2}.
\end{aligned}$$

Therefore,

$$|\text{score}_{1224}| \leq \frac{1}{N V L^2 \sqrt{d}} \left( \frac{N}{(V \wedge L^2 \wedge L \sqrt{d})^{3/2}} + \frac{\sqrt{N}}{L} \right) \leq \frac{o_V(1)}{V L^2}.$$

For  $\text{score}_{1223}$ , by Event,

$$\begin{aligned}
|\text{score}_{1223}| &\leq \frac{1}{N^2 L^3} \left( \sum_{i \neq j=1}^N |\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}| \right) \sup_{i \neq j} |(\alpha_{ij} - \phi'(0)^2) \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L| \\
&\quad + \frac{1}{N L^3} \sup_i |(\alpha_{ii} - \phi'(0)^2) \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L| \leq \frac{o_V(1)}{V L^2}.
\end{aligned}$$

For the first two terms, we define

$$\mathbf{V}_0 := \frac{1}{N L} \sum_{j=1}^N \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top, \quad \mathbf{V}_{0,1} := \frac{1}{N L} \sum_{i=1}^N \mathbf{x}_i (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top$$

$$\mathbf{V}_{0,2} := \frac{1}{NL} \sum_{i=1}^N (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top, \quad \mathbf{V}_{0,3} := \frac{1}{V} \mathbb{1}_V \frac{1}{NL} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top.$$

We have by Event,

$$|\text{score}_{1222}| \leq \frac{\phi'(0)^2}{\sqrt{d}} \left\| \frac{1}{L} \mathbf{Z}_{\text{in}} \mathbf{V}_0 \mathbf{V}_0^\top \mathbf{Z}_{\text{in}}^\top \right\|_F \leq \frac{\phi'(0)^2}{NV L^2 \sqrt{d}} \left\| \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \right\|_F \leq \frac{2\phi'(0)^2}{LdLN} \leq \frac{o_V(1)}{VL^2}.$$

Lastly,

$$|\text{score}_{1221}| = \frac{\phi'(0)^2}{L} \text{tr}(\mathbf{Z}_{\text{in}} \mathbf{V}_0 \mathbf{V}_0^\top \mathbf{Z}_{\text{in}}^\top) \leq \frac{\phi'(0)^2}{NV L^2 \sqrt{d}} \text{tr}(\mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top) \leq \frac{2\phi'(0)^2}{NL^2 \sqrt{d}} \leq \frac{o_V(1)}{VL^2}.$$

Therefore,  $|\text{score}_{122}| \ll \frac{1}{VL^2}$ . Lastly, we consider  $\text{score}_{121}$ . By Chebyshev's inequality, we have

$$\begin{aligned} & |\text{score}_{121}| \\ &= \phi'(0)^2 \text{tr}(\mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,1}) + \phi'(0)^2 \text{tr}(\mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,2}) + \phi'(0)^2 (L-1) \text{tr}(\mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,3}) \\ &\pm \frac{1}{\sqrt{d}} \left\| \mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,1} \right\|_F \pm \frac{1}{\sqrt{d}} \left\| \mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,2} \right\|_F \pm \frac{L-1}{\sqrt{d}} \left\| \mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,3} \right\|_F. \end{aligned}$$

We have the following:

- The first summand:

$$\text{tr}(\mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,1}) = \text{tr}(\mathbf{Z}_{\text{in}} \mathbf{V}_{0,1} \mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top) \asymp \frac{1}{VL^2}$$

- The second summand: By Chebyshev's inequality, we have

$$\begin{aligned} \text{tr}(\mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,2}) &= \frac{1}{N^2 L^2} \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}) \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \\ &= \pm \frac{1}{\sqrt{d}} \left\| \mathbf{V}_{0,2} \mathbf{V}_{0,1}^\top \right\|_F. \end{aligned}$$

We have by Event,

$$\left\| \mathbf{V}_{0,2} \mathbf{V}_{0,1}^\top \right\|_F \leq \frac{C \log V}{VL} \left\| \mathbf{V}_{0,2} \right\|_F \leq \frac{C \log V}{VL \sqrt{LN}}.$$

- The third summand: We have

$$\begin{aligned} & (L-1) \text{tr}(\mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,2}) \\ &= \frac{L-1}{N^2 L^2 V} \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}) \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \\ &= \frac{L-1}{N^2 L^2 V} \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}) \pm \frac{L-1}{N^2 L^2 V \sqrt{d}} \left\| \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}) \mathbf{x}_i \right\|_2. \end{aligned}$$

We have by Event

$$\left| \frac{L-1}{N^2 L^2 V} \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}) \right| \leq \frac{C \log V}{NLV^{3/2}}.$$

and

$$\frac{L-1}{N^2 L^2 V \sqrt{d}} \left\| \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}) \mathbf{x}_i \right\|_2 = \frac{L-1}{LV \sqrt{d}} \left\| \mathbf{V}_{0,1} \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \leq \frac{1}{\sqrt{NV^2 L \sqrt{d}}}$$

- The fourth summand: We have by Event

$$\frac{1}{\sqrt{d}} \left\| \mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,1} \right\|_F = \frac{1}{\sqrt{d}} \left\| \mathbf{Z}_{\text{in}} \mathbf{V}_{0,1} \mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \right\|_F \leq \frac{C}{VL^2d}$$

- The fifth summand: We have by Event

$$\begin{aligned} \left\| \mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,2} \right\|_F^2 &\leq \text{tr} \left( \mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,2} \mathbf{V}_{0,2}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,1} \right) \\ &\leq \frac{1}{NLd} \text{tr}(\mathbf{Z}_{\text{in}} \mathbf{V}_{0,1} \mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top) \leq \frac{V}{NLd} \frac{1}{V^2L^2} = \frac{1}{NVL^3d}. \end{aligned}$$

Therefore,

$$\frac{1}{\sqrt{d}} \left\| \mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,2} \right\|_F \leq \frac{1}{\sqrt{NV}L^{3/2}d}$$

- The sixth summand:

$$(L-1)^2 \left\| \mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,3} \right\|_F^2 \leq \frac{1}{V^2N} \mathbb{1}_V \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,1} \mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \leq \frac{1}{V^2NL^2d}$$

Therefore,

$$\frac{L-1}{\sqrt{d}} \left\| \mathbf{V}_{0,1}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{V}_{0,3} \right\|_F \leq \frac{C}{VL\sqrt{Nd}}.$$

Then,

$$\text{score}_{121} = \frac{1 \pm o_V(1)}{VL^2} \pm \frac{1}{\sqrt{NV}L^{3/2}d}.$$

Therefore, we have

$$\text{score}_{12} = \delta s_{\nu,\delta} \left( \frac{1 \pm o_V(1)}{VL^2} \pm \frac{1}{\sqrt{NV}L^{3/2}d} \right).$$

### C.3.2 CONCENTRATION BOUND FOR $\mathbf{s}_2$

We have

$$\begin{aligned} &\mathbf{e}_l^\top \mathbf{s}_2 \\ &= \frac{1}{N^2L^2} \sum_{i,j=1}^N \beta_{ij} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \\ &\quad - \frac{1}{N^2L^3} \sum_{i,j=1}^N \beta_{ij} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \\ &\quad + \frac{\delta s_{\nu,\delta}}{N^2L^2} \sum_{i,j=1}^N \beta_{ij} (\mathbf{e}_1 - \frac{1}{L} \mathbb{1}_L)^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \\ &=: \text{score}_{21} + \text{score}_{22} + \text{score}_{23}. \end{aligned}$$

**Concentration for  $\text{score}_{21}$ :** We will write  $\text{score}_{21}$  as follows:

$$\begin{aligned} \text{score}_{21} &= \frac{1}{N^2L^2} \sum_{i,j=1}^N (\beta_{ij} - \phi'(0)^2) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \\ &\quad \times (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \\ &\quad + \frac{\phi'(0)^2}{N^2L^2} \sum_{i,j=1}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{x}_i + \frac{L-1}{V} \mathbb{1}_V) \end{aligned}$$

$$\begin{aligned}
& \times (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \\
& + \frac{\phi'(0)^2}{N^2 L^2} \sum_{i,j=1}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{x}_i + \frac{L-1}{V} \mathbb{1}_V) \\
& \times (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \\
& + \frac{\phi'(0)^2}{N^2 L^2} \sum_{i,j=1}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \\
& \times (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \\
& + \frac{\phi'(0)^2}{N^2 L^2} \sum_{i,j=1}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i \mathbf{N}_i^\top - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \\
& \times (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \\
& + \frac{\phi'(0)^2}{N^2 L V} \sum_{i,j=1}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \\
& =: \text{score}_{211} + \text{score}_{212} + \text{score}_{213} + \text{score}_{214} + \text{score}_{215} + \text{score}_{216}
\end{aligned}$$

By Chebyshev's inequality, we have

$$\begin{aligned}
\text{score}_{211} &= \frac{1}{N^2 L^2} \sum_{i,j=1}^N (\beta_{ij} - \phi'(0)^2) (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \\
&\pm \frac{1}{N^2 L^2 \sqrt{d}} \left\| \sum_{i,j=1}^N (\beta_{ij} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \right\|_F \\
&=: \text{score}_{2111} + \text{score}_{2112}.
\end{aligned}$$

By using Events

$$\begin{aligned}
|\text{score}_{2111}| &\leq \frac{1}{N^2 L^2} \left( \sum_{i=1}^N \left( \sum_{j=1}^N (\beta_{ij} - \phi'(0)^2) (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}) \right)^2 \right)^{\frac{1}{2}} \left( \sum_{i=1}^N (\mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L)^2 \right)^{\frac{1}{2}} \\
&\leq \frac{C \log^2 V}{N^{3/2} L^2} \left( \frac{\sqrt{N}}{L} + \frac{N^{3/2}}{V} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}} \right) \left( \frac{L}{\sqrt{d} (L \wedge d)^{\frac{1}{2}}} \right) \\
&+ \frac{C \log^2 V}{N^{3/2} L V} \left( \frac{\sqrt{N}}{L} + \frac{N^{3/2}}{V} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}} \right) \left( \frac{\sqrt{L} V}{d^{3/2}} \right) \\
&\leq \frac{C \log^2 V}{V L^2 \sqrt{d} (L \wedge d)^{\frac{1}{2}}} + \frac{C \log^2 V}{N L^{3/2} d^{3/2}} + \frac{C \log^2 V}{V \sqrt{L} d^{3/2}} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}}.
\end{aligned}$$

Moreover, by Chebyshev's inequality

$$\begin{aligned}
|\text{score}_{2112}| &\leq \frac{1}{N^{3/2} L^2 \sqrt{d}} \left\| \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2 \\
&\times \left( \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N (\beta_{ij} - \phi'(0)^2)^2 |\mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L|^2 \right)^{\frac{1}{2}} \\
&\leq \frac{C \log^2 V}{N^{3/2} L^2 \sqrt{d}} \frac{N}{V} \left( \frac{\sqrt{N}}{V \wedge L^2 \wedge L \sqrt{d}} + \frac{1}{L} \right) \left( \frac{\sqrt{L}}{\sqrt{d}} + \frac{L^{3/2}}{d^{3/2}} \right) \\
&= \frac{C \log^2 V}{\sqrt{N} L V d (L \wedge d)} \left( \frac{\sqrt{N}}{V \wedge L^2 \wedge L \sqrt{d}} + \frac{1}{L} \right)
\end{aligned}$$

$$\leq \frac{C \log^2 V}{V \sqrt{L} d (L \wedge d)} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}} + \frac{C \log^2 V}{\sqrt{N} V L^{3/2} d (L \wedge d)}.$$

Therefore,

$$|\text{score}_{211}| \leq \frac{1}{V L^2 \sqrt{d}} + \frac{C \log^2 V}{N L^{3/2} d^{3/2}} + C \log^2 V \left( \frac{1}{V \sqrt{L} d (L \wedge d)} + \frac{1}{V \sqrt{L} d^{3/2}} \right) \frac{1}{V \wedge L^2 \wedge L \sqrt{d}}.$$

Moreover, by Chebyshev's inequality, we have

$$\begin{aligned} \text{score}_{212} &= \frac{(1 - \frac{1}{V})}{N^2 L^2} \left( \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right)^\top \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i + \frac{L-1}{V} \mathbb{1}_V \right) \\ &\pm \frac{1}{N^2 L^2 \sqrt{d}} \left\| \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \left\| \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i + \frac{L-1}{V} \mathbb{1}_V \right) \right\|_2. \end{aligned}$$

Let  $n_i := |\{j \leq N | \mathbf{x}_j = \mathbf{e}_i\}|$ . We have

$$\begin{aligned} &\frac{1}{N^2 L^2} \left\| \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \left\| \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i + \frac{L-1}{V} \mathbb{1}_V \right) \right\|_2 \\ &\frac{1}{N^2 L^2} \left\| \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \left\| \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i + \frac{L-1}{V} \mathbb{1}_V \right) \right\|_2 \\ &+ \frac{L-1}{L} \frac{1}{N^2 L V} \left\| \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \left\| \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right\|_2 \\ &\leq \frac{C}{N L^2} \left( \frac{1}{N} \sum_{i=1}^V n_i^2 \left| \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{e}_i \mathbf{e}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{e}_i \right|^2 \right)^{\frac{1}{2}} \\ &+ \frac{C}{\sqrt{N} L V} \left\| \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right\|_2 \end{aligned}$$

By Event, we have

$$\left| \frac{1}{N} \sum_{i=1}^V n_i^2 \right| \leq \frac{C N}{V} \quad \text{and} \quad \sup_{i \leq N} \left| \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{e}_i \mathbf{e}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{e}_i \right| \leq \frac{C \log V}{\sqrt{d}}.$$

Moreover,

$$\begin{aligned} &\mathbb{E} \left[ \left\| \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right\|_2^2 | \mathbf{Z}_{\text{in}} \right] \\ &\leq \frac{1}{N} \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} | \mathbf{Z}_{\text{in}} \right] \\ &+ \frac{1}{N^2} \mathbb{E} \left[ \sum_{i \neq j=1}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_j \mathbf{x}_j^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} | \mathbf{Z}_{\text{in}} \right] \\ &\leq \left( \frac{1}{N} + \frac{N-1}{N V} \right) \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} | \mathbf{Z}_{\text{in}} \right] \end{aligned}$$

$$\leq \left(\frac{1}{N} + \frac{N-1}{NV}\right) \frac{CV \log V}{d^2},$$

where we used Events in the last step. Therefore, by Chebyshev's inequality, we have

$$|\text{score}_{212}| \leq \frac{C \log V}{\sqrt{NV} L^2 \sqrt{d}} \left(1 + \frac{L}{\sqrt{Vd}}\right).$$

Moreover, by using Chebyshev's inequality

$$\begin{aligned} \text{score}_{213} &= \frac{1}{N^2 L^2} \left( \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right)^\top \\ &\times \left( \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i + \frac{L-1}{V} \mathbb{1}_V \right) \right) \\ &\pm \frac{1}{N^2 L^2 \sqrt{d}} \left\| \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \\ &\times \left\| \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i + \frac{L-1}{V} \mathbb{1}_V \right) \right\|_2 \end{aligned}$$

We have

$$\begin{aligned} &\mathbb{E} \left[ \left( \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i + \frac{L-1}{V} \mathbb{1}_V \right) \right)^2 \middle| \mathbf{Z}_{\text{in}} \right] \\ &= \sum_{i=1}^N \mathbb{E} \left[ \left( \sum_{j=1}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \right)^2 \middle| \mathbf{Z}_{\text{in}} \right] \\ &\leq 2(1 - \frac{1}{V})^2 \sum_{i=1}^N \mathbb{E} \left[ \left( \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i + \frac{L-1}{V} \mathbb{1}_V \right) \right)^2 \middle| \mathbf{Z}_{\text{in}} \right] \\ &+ \frac{2(1 - \frac{1}{V})}{V} \sum_{i=1}^N \sum_{j \neq i}^N \mathbb{E} \left[ \left( \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i + \frac{L-1}{V} \mathbb{1}_V \right) \right)^2 \middle| \mathbf{Z}_{\text{in}} \right] \quad (36) \end{aligned}$$

We have

$$\mathbb{E} \left[ \left( \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i + \frac{L-1}{V} \mathbb{1}_V \right) \right)^2 \middle| \mathbf{Z}_{\text{in}} \right] \leq \frac{CL}{d^2} \left(1 + \frac{L^2}{V}\right).$$

Therefore, we have (36)  $\leq \frac{CN^2 L}{Vd^2}$ . Also,

$$\mathbb{E} \left[ \left\| \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i + \frac{L-1}{V} \mathbb{1}_V \right) \right\|_2^2 \middle| \mathbf{Z}_{\text{in}} \right] \leq \frac{CL}{d^2} \left(1 + \frac{L^2}{V}\right).$$

Therefore, by Events and Chebyshev's inequality, we have

$$|\text{score}_{213}| \leq C \log V \left( \frac{1}{N\sqrt{V}L^{3/2}d} + \frac{1}{NV Ld} + \frac{1}{NL^{3/2}d^{3/2}} + \frac{1}{N\sqrt{V}\sqrt{L}d^{3/2}} \right).$$

Moreover, by Chebyshev's inequality

$$\begin{aligned} \text{score}_{214} &= \\ &\frac{1}{N^2 L^2} \left( \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right)^\top \end{aligned}$$

$$\begin{aligned}
& \times \left( \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right) \\
& \pm \frac{1}{N^2 L^2 \sqrt{d}} \left\| \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \\
& \times \left\| \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\|_2.
\end{aligned}$$

We have

$$\begin{aligned}
& \mathbb{E} \left[ \left( \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right)^2 \middle| \mathbf{Z}_{\text{in}} \right] \\
& = \sum_{i=1}^N \mathbb{E} \left[ \left( \sum_{j=1}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right)^2 \middle| \mathbf{Z}_{\text{in}} \right] \\
& \leq 2(1 - \frac{1}{V})^2 \sum_{i=1}^N \mathbb{E} \left[ \left( \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right)^2 \middle| \mathbf{Z}_{\text{in}} \right] \\
& + \frac{2(1-\frac{1}{V})}{V} \sum_{i=1}^N \sum_{j \neq i}^N \mathbb{E} \left[ \left( \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right)^2 \middle| \mathbf{Z}_{\text{in}} \right] \quad (37)
\end{aligned}$$

We have

$$\begin{aligned}
& \mathbb{E} \left[ \left( \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right)^2 \middle| \mathbf{Z}_{\text{in}} \right] \\
& \leq \frac{CL}{d} \mathbf{z}_{\nu, \delta}^\top \mathbb{E} \left[ \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \middle| \mathbf{Z}_{\text{in}} \right] \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \leq \frac{CL}{d^2}.
\end{aligned}$$

Therefore, (37)  $\leq \frac{CN^2L}{Vd^2}$ . Also,

$$\begin{aligned}
& \mathbb{E} \left[ \left\| \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\|_2^2 \middle| \mathbf{Z}_{\text{in}} \right] \\
& = \frac{CL}{d} \sum_{i=1}^N \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} \left[ \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{V} \mathbf{I}_V \right) \middle| \mathbf{Z}_{\text{in}} \right] \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \leq \frac{CNL}{d^2}.
\end{aligned}$$

Therefore, by Events and by Chebyshev's inequality, we have

$$|\text{score}_{214}| \leq C \log V \left( \frac{1}{N\sqrt{V}L^{3/2}d} + \frac{1}{NL^{3/2}d^{3/2}} \right)$$

Moreover, let

$$\begin{aligned}
\gamma_i &:= \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \\
& - \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \middle| \mathbf{Z}_{\text{in}} \right].
\end{aligned}$$

We have

$$\begin{aligned}
& \text{score}_{215} = \\
& \frac{1}{N^2 L^2} \left( \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right)^\top \left( \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \gamma_i \right) \\
& \pm \frac{1}{N^2 L^2 \sqrt{d}} \left\| \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \left\| \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \gamma_i \right\|_2
\end{aligned}$$



$$\begin{aligned}
& + \frac{1}{N^2 L^2} \left( \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right)^\top \\
& \times \left( \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} \left[ \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top \right) \mathbb{1}_{L-1} | \mathbf{Z}_{\text{in}} \right] \right) \\
& \pm \frac{1}{N^2 L^2 \sqrt{d}} \left\| \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \\
& \times \left\| \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} \left[ \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top \right) \mathbb{1}_{L-1} | \mathbf{Z}_{\text{in}} \right] \right\|_2.
\end{aligned}$$

By Proposition 9

$$\begin{aligned}
& \mathbb{E} \left[ \left( \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \gamma_i \right)^2 \right] = \sum_{i=1}^N \mathbb{E} \left[ \left( \sum_{j=1}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \gamma_i \right)^2 \right] \\
& \leq 2(1 - \frac{1}{V})^2 \sum_{i=1}^N \mathbb{E}[\gamma_i^2] + \frac{2(1 - \frac{1}{V})}{V} \sum_{i=1}^N \sum_{j \neq i}^N \mathbb{E}[\gamma_i^2] \leq C \log^2 V \frac{N^2}{V} \left( \frac{L}{d} + \frac{L^2}{d^2} \right).
\end{aligned}$$

Then,

$$\mathbb{E} \left[ \left\| \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \gamma_i \right\|_2^2 \right] \leq \sum_{i=1}^N \mathbb{E}[\gamma_i^2] \leq C N \log^2 V \left( \frac{L}{d} + \frac{L^2}{d^2} \right).$$

Moreover, by Proposition 9 and Event, we have

$$\begin{aligned}
& \left\| \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} \left[ \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top \right) \mathbb{1}_{L-1} | \mathbf{Z}_{\text{in}} \right] \right\|_2 \\
& \leq C \log V \frac{L \sqrt{N}}{\sqrt{V} d}.
\end{aligned}$$

Therefore, by Chebyshev's inequality, we have

$$|\text{score}_{215}| \leq C \log^2 V \left( \frac{1}{N L \sqrt{V} d (L \wedge d)^{1/2}} + \frac{1}{N L d (L \wedge d)^{1/2}} + \frac{1}{N L \sqrt{V} d} \right)$$

Lastly, by Chebyshev's inequality, we have

$$\begin{aligned}
\text{score}_{216} &= \frac{\phi'(0)^2}{N^2 L V} \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \\
& \pm \frac{\phi'(0)^2}{N^2 L V \sqrt{d}} \left\| \sum_{i,j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \right\|_F \\
&= \frac{\phi'(0)^2}{N^2 L V} \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \\
& \pm \frac{\phi'(0)^2}{V \sqrt{d}} \left\| \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \left\| \frac{1}{N L} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right\|_2
\end{aligned}$$

We have

$$\left\| \frac{1}{N L} \sum_{i=1}^N \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2$$

$$\begin{aligned}
&\leq \left\| \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \frac{1}{NL} \sum_{i=1}^N (\mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2 \\
&+ \frac{1}{V} \left\| \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right\| \left\| \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2 \\
&\leq \frac{CV}{d} \left\| \mathbf{Z}_{\text{in}} \frac{1}{NL} \sum_{i=1}^N (\mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2 + C \log V \frac{\sqrt{V}}{d^{3/2} \sqrt{N}} \\
&\leq \frac{CV}{\sqrt{NL} d^{3/2}} + C \log V \frac{\sqrt{V}}{d^{3/2} \sqrt{N}} \leq \frac{CV}{\sqrt{NL} d^{3/2}}.
\end{aligned}$$

Moreover,

$$\begin{aligned}
&\mathbb{E} \left[ \left( \frac{1}{N^2 L V} \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right)^2 | \mathbf{Z}_{\text{in}} \right] \\
&= \frac{1}{N^4 L^2 V^2} \sum_{j=1}^N \mathbb{E} \left[ \left( \sum_{i=1}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right)^2 | \mathbf{Z}_{\text{in}} \right] \\
&\leq \frac{2}{N^4 L^2 V^2} \sum_{j=1}^N \mathbb{E} \left[ \left( \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right)^2 | \mathbf{Z}_{\text{in}} \right] \\
&+ \frac{2}{N^4 L^2 V^2} \sum_{j=1}^N \mathbb{E} \left[ \left( \sum_{\substack{i=1 \\ i \neq j}}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right)^2 | \mathbf{Z}_{\text{in}} \right] \\
&\leq \frac{2}{N^4 L^2 V^2} \sum_{j=1}^N \mathbb{E} \left[ \left( \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right)^2 | \mathbf{Z}_{\text{in}} \right] \\
&+ \frac{2}{N^2 V^3} \sum_{j=1}^N \mathbb{E} \left[ \left\| \frac{1}{NL} \sum_{\substack{i=1 \\ i \neq j}}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2^2 | \mathbf{Z}_{\text{in}} \right] \\
&\leq \frac{2}{N^4 L^2 V^2} \sum_{j=1}^N \mathbb{E} \left[ \left( \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right)^2 | \mathbf{Z}_{\text{in}} \right] + \frac{C}{N^2 V L d^3},
\end{aligned}$$

where we used C.3.2 in the last step. We have

$$\begin{aligned}
&\frac{1}{N^4 L^2 V^2} \sum_{j=1}^N \mathbb{E} \left[ \left( \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right)^2 | \mathbf{Z}_{\text{in}} \right] \\
&\leq \frac{1}{N^3 L^2 V^2} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \frac{L^2}{V^2} \mathbb{1}_V \mathbb{1}_V^\top + \frac{L}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}} \mathbf{z}_{\nu,\delta} \leq \frac{C \log^2 V}{N^3 L d^3}
\end{aligned}$$

Therefore, by Chebyshev's inequality, we have

$$|\text{score}_{216}| \leq C \log V \left( \frac{1}{N \sqrt{L} d^2} + \frac{1}{N \sqrt{V} L d^{3/2}} \right) \leq \frac{C \log V}{N \sqrt{L} d^2}$$

Overall, by using  $N \ll VL$ ,

$$\begin{aligned}
&|\text{score}_{21}| \\
&\leq C \log^2 V \left( \frac{1}{N \sqrt{L} d (L \wedge d)} + \frac{1}{N L d (L \wedge d)^{1/2}} + \frac{1}{N L \sqrt{V} d} + \frac{1}{\sqrt{N} V L d} + \frac{1}{\sqrt{N} V L^2 \sqrt{d}} \right) \\
&+ C \log^2 V \left( \frac{1}{V L^2 \sqrt{d}} \frac{1}{V^2 \sqrt{L} d^{3/2}} \right).
\end{aligned}$$

### C.3.3 CONCENTRATION BOUND FOR $\mathbf{s}_3$

We have

$$\begin{aligned}
& \mathbf{e}_l^\top \mathbf{s}_3 \\
&= \frac{1}{N^2 L} \sum_{i,j=1}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \\
&\quad \times \left( \frac{1}{m} \sum_{k=1}^m \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \phi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) \right. \\
&\quad \left. - \mathbb{E} \left[ \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \phi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) \right] \right) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \\
&\quad - \frac{1}{N^2 L^2} \sum_{i,j=1}^N \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mathbb{1}_L^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \\
&\quad \times \left( \frac{1}{m} \sum_{k=1}^m \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \phi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) \right. \\
&\quad \left. - \mathbb{E} \left[ \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \phi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) \right] \right) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \\
&\quad + \frac{\delta s_{\nu,\delta}}{N^2 L} \sum_{i,j=1}^N (\mathbf{e}_1 - \frac{1}{L} \mathbb{1}_L)^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \\
&\quad \times \left( \frac{1}{m} \sum_{k=1}^m \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \phi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) \right. \\
&\quad \left. - \mathbb{E} \left[ \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \phi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) \right] \right) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \\
&=: \text{score}_{31} + \text{score}_{32} + \text{score}_{33}.
\end{aligned}$$

**Concentration bound for  $\text{score}_{31}$ :** We start with  $\text{score}_{31}$ . We have

$$\begin{aligned}
\text{score}_{31k} &:= \text{tr} \left( \frac{1}{NL} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right. \\
&\quad \left. \times \frac{1}{N} \sum_{j=1}^N \phi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{Z}_{\text{out}}^\top \mathbf{Z}_{\text{out}} \right) \\
&= \text{tr} \left( \frac{1}{NL} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right. \\
&\quad \left. \times \frac{1}{N} \sum_{j=1}^N \phi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\
&\quad \pm \frac{1}{\sqrt{d}} \left\| \frac{1}{NL} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right\|_2 \\
&\quad \left\| \frac{1}{N} \sum_{j=1}^N \phi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \\
&=: \text{score}_{31k_1} + \text{score}_{31k_2},
\end{aligned}$$

where we used Chebyshev's inequality for the second step. We define

$$\phi(t) =: \phi(0) + t\psi(t) \quad \text{and} \quad \phi'(t) =: \phi(0) + t\psi_1(t) \quad \text{and} \quad \psi(t) =: \psi(0) + t\psi_2(t).$$

and write

$$\begin{aligned} \text{score}_{31k_1} &= \phi(0)\phi'(0) \operatorname{tr} \left( \frac{1}{NL} \sum_{i=1}^N \mathbf{x}_i \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\ &+ \phi(0) \operatorname{tr} \left( \frac{1}{NL^2} \sum_{i=1}^N \mathbf{x}_i \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \psi_1 \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\ &+ \phi(0) \operatorname{tr} \left( \frac{1}{NL^2} \sum_{i=1}^N \mathbf{x}_i \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \psi_1 \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\ &+ \operatorname{tr} \left( \frac{1}{NL} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right. \\ &\quad \left. \times \frac{1}{N} \sum_{j=1}^N \psi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\ &=: \text{score}_{31k_{11}} + \text{score}_{31k_{12}} + \text{score}_{31k_{13}} + \text{score}_{31k_{14}}. \end{aligned}$$

In the following, we bound each term separately.

• We have

$$\begin{aligned} \text{score}_{31k_{11}} &= \frac{1}{L} \sum_{w=1}^V \left( \frac{n_w}{N} - \frac{1}{V} \right) \frac{n_w}{N} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{e}_w \mathbf{e}_w^\top + \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \\ &\quad + \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \frac{1}{NL} \sum_{w=1}^V \left( \frac{n_w}{N} - \frac{1}{V} \right) \sum_{i \in \{i_1, \dots, i_{n_w}\}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \end{aligned}$$

We have

$$\begin{aligned} &\mathbb{E} \left[ \left( \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \frac{1}{NL} \sum_{w=1}^V \left( \frac{n_w}{N} - \frac{1}{V} \right) \sum_{i \in \{i_1, \dots, i_{n_w}\}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \right)^2 \middle| \mathbf{Z}_{\text{in}} \right] \\ &= \mathbb{E} \left[ \left\| \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \frac{1}{NL} \sum_{w=1}^V \left( \frac{n_w}{N} - \frac{1}{V} \right) \sum_{i \in \{i_1, \dots, i_{n_w}\}} \left( \mathbf{N}_i^\top \mathbf{N}_i - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \right\|_2^2 \middle| \mathbf{Z}_{\text{in}} \right] \\ &= \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \frac{1}{N^2 L^2} \sum_{w=1}^V \mathbb{E} \left[ \left( \frac{n_w}{N} - \frac{1}{V} \right)^2 n_w \right] \mathbb{E} \left[ \left( \mathbf{N}_1^\top \mathbf{N}_1 - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_1^\top \mathbf{N}_1 - \frac{L-1}{V} \mathbf{I}_V \right) \middle| \mathbf{Z}_{\text{in}} \right] \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \\ &\leq \frac{C}{N^2 L^2 V} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} \left[ \left( \mathbf{N}_1^\top \mathbf{N}_1 - \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{N}_1^\top \mathbf{N}_1 - \frac{L-1}{V} \mathbf{I}_V \right) \middle| \mathbf{Z}_{\text{in}} \right] \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} = \frac{C}{N^2 V d(L \wedge d)} \end{aligned}$$

Moreover, by using Event, we write

$$\begin{aligned} &\mathbb{E} \left[ \left( \frac{1}{L} \sum_{w=1}^V \left( \frac{n_w}{N} - \frac{1}{V} \right) \frac{n_w}{N} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{e}_w \mathbf{e}_w^\top + \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \right)^2 \middle| \mathbf{Z}_{\text{in}} \right] \\ &= \frac{1}{L^2} \mathbb{E} \left[ \left\| \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \sum_{w=1}^V \left( \frac{n_w}{N} - \frac{1}{V} \right) \frac{n_w}{N} \mathbf{e}_w \mathbf{e}_w^\top + \left( \frac{n_w}{N} - \frac{1}{V} \right)^2 \frac{L-1}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \right\|_2^2 \middle| \mathbf{Z}_{\text{in}} \right] \end{aligned}$$

$$\begin{aligned}
&\leq \frac{V^2}{L^2 d^2} \mathbb{E} \left[ \left\| \sum_{w=1}^V \left( \frac{n_w}{N} - \frac{1}{V} \right) \frac{n_w}{N} \mathbf{e}_w \mathbf{e}_w^\top + \left( \frac{n_w}{N} - \frac{1}{V} \right)^2 \frac{L-1}{V} \mathbf{I}_V \right\|_2^2 | \mathbf{Z}_{\text{in}} \right] \\
&\leq \frac{CV^2}{L^2 d^2} \mathbb{E} \left[ \sup_{w \in [N]} \left| \left( \frac{n_w}{N} - \frac{1}{V} \right) \frac{n_w}{N} \right|^2 \right] + \frac{C}{d^2} \mathbb{E} \left[ \left( \sum_{w=1}^V \left( \frac{n_w}{N} - \frac{1}{V} \right)^2 \right) \right] \leq \frac{C}{d^2 N^2}.
\end{aligned}$$

Therefore,

$$\mathbb{E} [\text{score}_{31k_{11}}^2 | \mathbf{Z}_{\text{in}}] \leq \frac{C}{d^2 N^2}.$$

• Moreover,

$$\text{score}_{31k_{12}}^2 \leq \frac{C}{N} \left\| \frac{1}{NL^2} \sum_{i=1}^N \mathbf{x}_i \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \psi_1 \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right\|_2^2.$$

We have for any  $i \in [N]$ , by Event,

$$\left| \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \psi_1 \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right| \leq C(\log^4 V) \sqrt{L} (\mathbb{1}_{\mathbf{x}_i = \mathbf{e}_\nu} + \frac{1}{\sqrt{d}})$$

Then,

$$\begin{aligned}
&\left\| \frac{1}{NL^2} \sum_{i=1}^N \mathbf{x}_i \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \psi_1 \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right\|_2^2 \\
&\leq \frac{C(\log^8 V)}{L^3} \left\| \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i (\mathbb{1}_{\mathbf{x}_i = \mathbf{e}_\nu} + \frac{1}{\sqrt{d}}) \right\|^2 \leq \frac{C(\log^8 V)}{V d L^3}
\end{aligned}$$

Then,

$$\mathbb{E} [\text{score}_{31k_{12}}^2 | \mathbf{Z}_{\text{in}}] \leq \frac{C \log^8 V}{N V d L^3}.$$

• Moreover,

$$\text{score}_{31k_{13}}^2 \leq \frac{C}{N} \left\| \frac{1}{NL^2} \sum_{i=1}^N \mathbf{x}_i \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \psi_1 \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right\|_2^2.$$

We have for any  $i \in [N]$ , by Events

$$\begin{aligned}
&\left| \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \psi_1 \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right| \\
&\leq C(\log^8 V) \sqrt{L} \|\mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta}\|_2 \leq C(\log^8 V) \sqrt{L} (\mathbf{e}_\nu^\top \mathbf{N}_i^\top \mathbb{1}_{L-1} + \frac{L}{d} + \log^6 V \sqrt{\frac{L}{d}})
\end{aligned}$$

Then, by Events

$$\begin{aligned}
&\left\| \frac{1}{NL^2} \sum_{i=1}^N \mathbf{x}_i \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \psi_1 \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right\|_2^2 \\
&\leq \frac{C(\log V)^2}{N^2 L^3} \left\| \sum_{i=1}^N \mathbf{x}_i \left( \mathbb{1}_{L-1}^\top \mathbf{N}_i \mathbf{e}_\nu + \frac{L}{d} + \log^6 V \sqrt{\frac{L}{d}} \right) \right\|_2^2 \\
&\leq \frac{C \log^{14} V}{V L d (L \wedge d)} + \frac{C(\log V)^2}{N^2 L^3} \left\| \sum_{i=1}^N \mathbf{x}_i \mathbb{1}_{L-1}^\top \mathbf{N}_i \mathbf{e}_\nu \right\|_2^2
\end{aligned}$$

We have

$$\frac{C(\log V)^2}{N^2 L^3} \mathbb{E} \left[ \left\| \sum_{i=1}^N \mathbf{x}_i \mathbb{1}_{L-1}^\top \mathbf{N}_i \mathbf{e}_\nu \right\|_2^2 \right] \leq \frac{C(\log V)^2}{N^2 L^3} \left( \frac{N^2 L^2}{V} + N \frac{L}{V} \right)$$

$$= \frac{C(\log V)^2}{V^3 L} + \frac{C(\log V)^2}{NV L^2}.$$

Then,

$$\mathbb{E}[\text{score}_{31k_{13}}^2 | Z_{\text{in}}] \leq \frac{C \log^{20} V}{NV L d (L \wedge d)}.$$

• Lastly, we have

$$\begin{aligned} |\text{score}_{31k_{14}}| &\leq \left\| \frac{1}{NL} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right\|_2 \\ &\quad \times \left\| \frac{1}{N} \sum_{j=1}^N \psi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2. \end{aligned}$$

By using the derivations in the two previous items, we have

$$\begin{aligned} &\left\| \frac{1}{NL} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right\|_2 \\ &\leq \left\| \frac{1}{NL} \sum_{i=1}^N \mathbf{x}_i \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right\|_2 \\ &\quad + \left\| \frac{1}{NL^2} \sum_{i=1}^N \mathbf{x}_i \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \psi_1 \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right\|_2 \\ &\quad + |\phi'(0)| \left\| \frac{1}{NL} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \right\|_2 \\ &\leq \frac{C \log^7 V}{\sqrt{V L d (L \wedge d)}} + \phi'(0) \left\| \frac{1}{NL} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \right\|_2 \end{aligned}$$

We have

$$\begin{aligned} &\mathbb{E} \left[ \left\| \frac{1}{NL} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \right\|_2^2 \middle| Z_{\text{in}} \right] \\ &\leq \frac{1}{N^2 L^2} \mathbb{E} \left[ \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i = \mathbf{x}_j} - \frac{1}{V}) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{N}_j^\top \mathbf{N}_j \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \middle| Z_{\text{in}} \right] \\ &\leq \frac{(1 - \frac{1}{V})}{NL^2} \frac{L}{V} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \text{Diag}(\mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}}) \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} + \frac{(1 - \frac{1}{V})}{NL^2} \frac{L^2}{V^2} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \\ &= \frac{(1 - \frac{1}{V})}{Nd(L \wedge d)} \end{aligned} \tag{38}$$

Moreover,

$$\begin{aligned} &\left\| \frac{1}{NL} \sum_{j=1}^N \psi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2 \\ &= |\psi(0)| \left\| \frac{1}{NL} \sum_{j=1}^N \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2 \\ &\quad + \left\| \frac{1}{NL^2} \sum_{j=1}^N \mathbf{x}_j \psi_2 \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) (\mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L)^2 \right\|_2. \end{aligned}$$

By Event, for all  $j \in [N]$ ,

$$|\psi_2 \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) (\mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L)^2| \leq C(\log^6 V) L.$$

Therefore,

$$\left\| \frac{1}{NL^2} \sum_{j=1}^N \mathbf{x}_j \psi_2 \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) (\mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L)^2 \right\|_2 \leq \frac{C(\log^6 V)}{\sqrt{VL}}.$$

Moreover,

$$\begin{aligned} & \left\| \frac{1}{NL} \sum_{j=1}^N \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2 \\ & \leq \left\| \frac{1}{NL} \sum_{j=1}^N \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2 \\ & \quad + \frac{1}{V} |\mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V| \left\| \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2. \end{aligned}$$

By using Event, we have

$$- \frac{1}{V} |\mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V| \left\| \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2 \leq \frac{C \log^2 V}{\sqrt{VN}}$$

– Moreover,

$$\begin{aligned} & \left\| \frac{1}{NL} \sum_{j=1}^N \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right\|_2^2 \\ & = \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \left( \frac{1}{NL} \sum_{j=1}^N (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\ & \quad \left( \frac{1}{NL} \sum_{j=1}^N (\mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right)^\top \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \\ & \leq \frac{C \log^2 V}{NL}. \end{aligned} \tag{39}$$

Then, for  $N \ll VL$

$$\mathbb{E}[\text{score}_{31k_2}^2 | \mathbf{Z}_{\text{in}}] \leq \frac{C \log^{16} V}{N^2 L d (L \wedge d)}$$

• On the other hand, we have

$$\begin{aligned} |\text{score}_{31k_2}| & \leq \frac{1}{\sqrt{d}} \left\| \frac{1}{NL} \sum_{i=1}^N (\mathbf{x}_i - \frac{1}{V} \mathbb{1}_V) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{w}_k \phi' \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \right) \right\|_2 \\ & \quad \times \left\| \frac{1}{N} \sum_{j=1}^N \phi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \end{aligned}$$

Note that by Event and (39), we have

$$\begin{aligned} & \left\| \frac{1}{N} \sum_{j=1}^N \phi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 = |\phi(0)| \left\| \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \\ & \quad + \left\| \frac{1}{N} \sum_{j=1}^N \psi \left( \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L \right) \frac{1}{L} \mathbf{w}_k^\top \mathbf{Z}_{\text{in}} \mathbf{X}_j^\top \mathbb{1}_L (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \\ & \leq \frac{1}{\sqrt{N}} + \frac{C \log^2 V}{\sqrt{NL}} \leq \frac{C}{\sqrt{N}} \end{aligned}$$

Therefore by (38), we have

$$\mathbb{E} \left[ |\text{score}_{31k_2}|^2 | \mathbf{Z}_{\text{in}} \right] \leq \frac{C}{N^2 d^2 (L \wedge d)}$$

Therefore, we have

$$\mathbb{E}[\text{score}_{31}|\mathbf{Z}_{\text{in}}] = 0 \text{ and } \text{Variance}(\text{score}_{31}|\mathbf{Z}_{\text{in}}) \leq \frac{C}{N^2 d^2 m}.$$

## D LOWER BOUND

To prove a lower bound, we construct a Bayesian setting with the same likelihood distribution in our setting. In particular, the ground truth permutation is chosen from the set of permutation matrices:

$$\mathcal{H} := \{\mathbf{P} \in \{0, 1\}^{V \times V} \mid \mathbf{P} \text{ is a permutation matrix}\}.$$

We describe our Bayesian setting as a game between `Environment` and `Learner` as follows:

- At the beginning, `Environment` samples  $\mathbf{P}_* \sim \text{Unif}(\mathcal{H})$ , probability vectors without revealing them to the learner.
- `Learner` observes  $L + 1$  channel that generates words from the set  $\mathcal{V} = \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_V\}$  sequentially for  $t = 1, 2, \dots, N$  with distributions:
  - At every round, `Environment` randomly picks a channel  $\ell_t$
  - *Label*: Channel 0 generates  $\mathbf{p}_t \sim_{\text{iid}} \text{Unif}(\mathcal{V})$
  - *Input*: Given  $\ell_t$  and  $\mathbf{p}_t$ , Channel  $\ell_t$  generates  $\mathbf{X}_{\ell_t, t} = \mathbf{P}_* \mathbf{p}_t$
  - *Noise distribution*: Channel  $j \in [L] \setminus \{\ell_t\}$  generate  $\mathbf{X}_{j, t} \sim \text{Unif}(\mathcal{V})$  independent of Channel 0.
- Let  $\mathcal{D} := \{(\mathbf{X}_t, \mathbf{p}_t)\}_{t \leq N}$  be the dataset. We study the Bayes estimator with 0 – 1 loss given the representation of the past:  $S = f(\mathcal{D}, \ell_{1:N})$ :

$$\hat{\mathbf{P}} = \arg \max_{\mathbf{P} \in \mathcal{H}} \mathbb{P}[\mathbf{P} = \mathbf{P}_* | S, \mathbf{Z}_{\text{in}}]. \quad (40)$$

In the following we consider the empirical mean and covariance of embedded words as the given data, i.e.,  $S := \{(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t, \mathbf{p}_t)\}_{t \leq N}$ , where

$$\boldsymbol{\mu}_t := \frac{1}{L} \mathbf{Z}_{\text{in}} \mathbf{X}_t^\top \mathbb{1}_L + \frac{\sigma_\mu}{\sqrt{L}} \mathbf{g}_t \text{ and } \boldsymbol{\Sigma}_t := \frac{1}{L} \mathbf{Z}_{\text{in}} \mathbf{X}_t^\top \mathbf{X}_t \mathbf{Z}_{\text{in}}^\top + \frac{\sigma_\Sigma}{\sqrt{dL}} \mathbf{G}_t.$$

where  $\{(\mathbf{g}_t, \mathbf{G}_t)\}_{t \leq N}$  are i.i.d. measurement noise with distributions  $\mathbf{g}_t \sim \mathcal{N}(0, \frac{1}{d} \mathbf{I}_d)$  and  $\mathbf{G}_{t, ij} = \mathbf{G}_{t, ji}$  with  $\mathbf{G}_{t, ij} \sim \mathcal{N}(0, \frac{(1+\delta_{ij})}{d})$  i.i.d. for  $i < j$ .

**Theorem 4.** *The following lower bound holds:*

$$\mathbb{P}[\hat{\mathbf{P}} \neq \mathbf{P}_* | \mathbf{Z}_{\text{in}}] \geq 1 - o_V(1) - \frac{\Omega(N)}{V} \left( 1 \wedge \left( \frac{1}{\sigma_\mu^2} \frac{d}{L \log V} + \frac{C}{\sigma_\Sigma^2} \frac{d^2}{L \log V} \right) \right)$$

We use an information-theoretic argument to prove Theorem 4. For the proof, let  $H(A)$  and  $H(A|C)$  denote the entropy and conditional entropy of  $A$  given  $C$ ; let  $I(A; B) = H(A) - H(A|B)$  and  $I(A; B|C) = H(A|C) - H(A|B, C)$  denote the mutual information between random variables  $A$  and  $B$  and the conditional mutual given  $C$ , respectively. We let  $D_{\text{KL}}$  denote the Kullback-Leibler (KL) divergence. We start with an auxiliary statement for the proof.

**Lemma 1.** *Let  $A, B, C, D$  be discrete random variables defined on the same probability space. The following statements hold:*

- In general,  $H(A|B, C) \leq H(A|B)$ . The equality is satisfied if and only if  $A \perp C | B$ .
- If  $B \perp D | (A, C)$ , we have  $I(A, B|C, D) \leq I(A, B|C)$ .
- Let  $S = g(A, C)$  be a measurable function of  $(A, C)$ . If  $B \perp A | (S, C, D)$ , then  $I(A; B|C, D) = I(S; B|C, D)$ .
- Given,  $\boldsymbol{\mu}, \boldsymbol{\mu}' \in \mathbb{R}^d$ , positive definite  $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$  and  $\text{supp}(A) \subseteq \mathbb{R}^d$ , we have

$$D_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu} + A, \boldsymbol{\Sigma}) || \mathcal{N}(\boldsymbol{\mu}' + A, \boldsymbol{\Sigma})) \leq \frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}')^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}').$$



*Proof.* We have

$$H(A|B) - H(A|B, C) = \mathbb{E} \left[ \log \frac{\mathbb{P}(A|B, C)}{\mathbb{P}(A|B)} \right] = \mathbb{E} \left[ \log \frac{\mathbb{P}(A, C|B)}{\mathbb{P}(A|B)\mathbb{P}(C|B)} \right] = I(A, C|B).$$

Since the mutual information is non-negative, the first item follows. Moreover, since  $I(A, C|B) = 0$  if and only if  $A \perp\!\!\!\perp C|B$ . For the second item, by using the first item,

$$I(A, B|C, D) = H(B|C, D) - H(B|A, C, D) \leq H(B|C) - H(B|A, C) = I(A, B|C).$$

For the third item, since  $S$  is a function of  $(A, C)$ , we have

$$\begin{aligned} I(A; B|C, D) &= I((A, S); B|C, D) = H(B|C, D) - H(B|A, S, C, D) \\ &= H(B|C, D) - H(B|S, C, D) = I(S; B|C, D). \end{aligned}$$

Let  $f$  denotes the Gaussian pdf with 0 and covariance  $\Sigma$ . For any  $\mathbf{x} \in \mathbb{R}^d$ , since  $t \rightarrow t \log t$  is convex

$$\begin{aligned} \left( \sum_{\mathbf{a} \in \text{supp}(A)} p(\mathbf{a}) f(\mathbf{x} - \boldsymbol{\mu} - \mathbf{a}) \right) \log \frac{\left( \sum_{\mathbf{a} \in \text{supp}(A)} p(\mathbf{a}) f(\mathbf{x} - \boldsymbol{\mu} - \mathbf{a}) \right)}{\left( \sum_{\mathbf{a} \in \text{supp}(A)} p(\mathbf{a}) f(\mathbf{x} - \boldsymbol{\mu}' - \mathbf{a}) \right)} \\ \leq \sum_{\mathbf{a} \in \text{supp}(A)} p(\mathbf{a}) f(\mathbf{x} - \boldsymbol{\mu} - \mathbf{a}) \log \frac{f(\mathbf{x} - \boldsymbol{\mu} - \mathbf{a})}{f(\mathbf{x} - \boldsymbol{\mu}' - \mathbf{a})}. \end{aligned}$$

Therefore, we have

$$\begin{aligned} D_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu} + A, \Sigma) || \mathcal{N}(\boldsymbol{\mu}' + A, \Sigma)) &\leq \sum_{\mathbf{a} \in \text{supp}(A)} p(\mathbf{a}) D_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu} + \mathbf{a}, \Sigma) || \mathcal{N}(\boldsymbol{\mu}' + \mathbf{a}, \Sigma)) \\ &= D_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu}, \Sigma) || \mathcal{N}(\boldsymbol{\mu}', \Sigma)), \end{aligned}$$

where the last inequality follows the invariance of KL divergence in the second line to constant shifts. The final bound follows the known formula for the KL divergence between Gaussian distributions.  $\square$

The proof of Theorem 4 is given in the following:

*Proof of Theorem 4.* Since we assume  $Z_{\text{in}}$  is known by the learner, we will fix it in the following without explicitly conditioning the terms on it. Note that we consider the Bayes decision rule in (40) and use Fano's inequality (Scarlett & Cevher, 2019) to lower bound its error probability:

$$\mathbb{P}[\hat{\mathbf{P}} \neq \mathbf{P}_* | Z_{\text{in}}] \geq 1 - \frac{I(\mathbf{P}_*; S) + \log 2}{\log |\mathcal{H}|}. \quad (41)$$

We have

$$\begin{aligned} I(\mathbf{P}_*; S) &= I(\mathbf{P}_*; \{(\boldsymbol{\mu}_t, \Sigma_t, \mathbf{p}_t)\}_{t \leq N}) = I(\mathbf{P}_*; \{\mathbf{p}_t\}_{t \leq N}) + I(\mathbf{P}_*; \{(\boldsymbol{\mu}_t, \Sigma_t)\}_{t \leq N} | \{\mathbf{p}_t\}_{t \leq N}) \\ &\stackrel{(a)}{=} I(\mathbf{P}_*; \{(\boldsymbol{\mu}_t, \Sigma_t)\}_{t \leq N} | \{\mathbf{p}_t\}_{t \leq N}) \\ &= \sum_{t=1}^N I(\mathbf{P}_*; (\boldsymbol{\mu}_t, \Sigma_t) | \{(\boldsymbol{\mu}_u, \Sigma_u)\}_{u < t}, \{\mathbf{p}_t\}_{t \leq N}) \end{aligned}$$

Given fixed  $Z_{\text{in}}$ , we observe that  $(\boldsymbol{\mu}_t, \Sigma_t) \perp\!\!\!\perp \{(\boldsymbol{\mu}_u, \Sigma_u)\}_{u < t} \mid \mathbf{P}_*, \{\mathbf{p}_t\}_{t \leq N}$  and  $(\boldsymbol{\mu}_t, \Sigma_t) \perp\!\!\!\perp \{\mathbf{p}_u\}_{u \neq t} \mid \mathbf{P}_*$ . Therefore, by Lemma 1,

$$I(\mathbf{P}_*; S) \leq \sum_{t=1}^N I(\mathbf{P}_*; (\boldsymbol{\mu}_t, \Sigma_t) | \{\mathbf{p}_t\}_{t \leq N}) \leq \sum_{t=1}^N I(\mathbf{P}_*; (\boldsymbol{\mu}_t, \Sigma_t) | \mathbf{p}_t).$$

Moreover, we have  $\mathbf{P}_* \perp\!\!\!\perp (\boldsymbol{\mu}_t, \Sigma_t) \mid \mathbf{X}_{\ell_t, t}, \mathbf{p}_t$ , where  $\mathbf{X}_{\ell_t, t}$  is a function of  $(\mathbf{P}_*, \mathbf{p}_t)$ . Therefore, by Lemma 1,

$$I(\mathbf{P}_*; S) \leq \sum_{t=1}^N I(\mathbf{X}_{\ell_t, t}; (\boldsymbol{\mu}_t, \Sigma_t) | \mathbf{p}_t).$$

We have

$$I(\mathbf{X}_{\ell_t, t}; (\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t) | \mathbf{p}_t, \mathbf{Z}_{\text{in}}) = \frac{1}{V} \sum_{k=1}^V D_{\text{KL}}(\mathbb{P}_{(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)}^k || \mathbb{P}_0) \stackrel{(a)}{\leq} \frac{1}{V^2} \sum_{j,k=1}^V D_{\text{KL}}(\mathbb{P}_{(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)}^k || \mathbb{P}_{(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)}^j)$$

where  $\mathbb{P}_{(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)}^k$  denotes the distribution of  $(\mathbf{s}_t, \boldsymbol{\Sigma}_t) | \mathbf{X}_{\ell_t, t} = \mathbf{e}_k$ ,  $\mathbb{P}_0$  denotes  $\mathbb{P}_0 = \frac{1}{V} \sum_{k=1}^V \mathbb{P}_{(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)}^k$ , and (a) follows the convexity of KL divergence in its second argument. For  $k \neq j$ , by the last item of Lemma 1, we have

$$D_{\text{KL}}(\mathbb{P}_{(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)}^k || \mathbb{P}_{(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)}^j) \leq \frac{C}{\sigma_{\boldsymbol{\mu}}^2} \frac{d}{L} \|\mathbf{z}_k - \mathbf{z}_j\|_2^2 + \frac{C}{\sigma_{\boldsymbol{\Sigma}}^2} \frac{d^2}{L} \|\mathbf{z}_k \mathbf{z}_k^\top - \mathbf{z}_j \mathbf{z}_j^\top\|_F^2 \leq \frac{C}{\sigma_{\boldsymbol{\mu}}^2} \frac{d}{L} + \frac{C}{\sigma_{\boldsymbol{\Sigma}}^2} \frac{d^2}{L}.$$

Therefore, we have

$$I(\mathbf{P}_*; S) \leq N \left( \frac{C}{\sigma_{\boldsymbol{\mu}}^2} \frac{d}{L} + \frac{C}{\sigma_{\boldsymbol{\Sigma}}^2} \frac{d^2}{L} \right).$$

Moreover, we can write

$$\begin{aligned} I(\mathbf{P}_*; S) &\leq I(\mathbf{P}_*; \mathcal{D}, \ell_{1:N}) = I(\mathbf{P}_*; \{\mathbf{X}_t\}_{t \leq N} | \{\mathbf{p}_t\}_{t \leq N}, \ell_{1:N}) \\ &\leq \sum_{t=1}^N I(\mathbf{P}_*; \mathbf{X}_{\ell_t, t} | \{\mathbf{p}_t, \ell_t\}_{t \leq N}) \\ &\leq \sum_{t=1}^N I(\mathbf{P}_*; \mathbf{X}_{\ell_t, t} | \mathbf{p}_t, \ell_t) \end{aligned}$$

where the first inequality follows data processing inequality, third and fourth inequalities follow the first and second items in Lemma 1. We have

$$I(\mathbf{P}_*; \mathbf{X}_{\ell_t, t} | \mathbf{p}_t, \ell_t) = \underbrace{H(\mathbf{X}_{\ell_t, t} | \mathbf{p}_t, \ell_t)}_{\log V} - \underbrace{H(\mathbf{X}_{\ell_t, t} | \mathbf{p}_t, \ell_t, \mathbf{P}_*)}_{=0} = \log V.$$

Therefore, we have  $I(\mathbf{P}_*; S) \leq N \log V$ . Finally, we have

$$I(\mathbf{P}_*; S) \leq N \left( \log V \wedge \left( \frac{C}{\sigma_{\boldsymbol{\mu}}^2} \frac{d}{L} + \frac{C}{\sigma_{\boldsymbol{\Sigma}}^2} \frac{d^2}{L} \right) \right).$$

The result follows from (41).  $\square$

## E AUXILIARY STATEMENTS

### E.1 A NICE EVENT CHARACTERIZATION

We characterize a “nice event” under which we use in the proof of Theorem 1 holds.

**Lemma 2.** We assume  $V^3 \gg N \gg V \gg L$  and  $L \asymp V^{\epsilon_1}$  and  $d \asymp V^{\epsilon_2}$  for some  $\epsilon_1, \epsilon_2 \in (0, 1)$ . For the following we define,  $m_{ij} := (1 - 1/V)\delta_{ij} + \frac{L}{V}$ . We define the following events:

(E.1) Let  $\mathbf{z}_\nu := \mathbf{Z}_{\text{in}} \mathbf{e}_\nu$  and  $\mathbf{z}_{\nu, \delta} := (\mathbf{z}_\nu + \mathbb{1}_{l=1} \delta \mathbf{z}_{\text{trig}})$ . We have

$$(E1.1) \quad \frac{1}{V} \|\mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top\|_2 \leq \frac{2}{d} \text{ and } \max_{k \leq V} \|\mathbf{z}_k\|_2 \vee \|\mathbf{z}_{\text{trig}}\|_2 \leq 2$$

$$(E1.2) \quad \frac{1}{\sqrt{V}} \|\mathbf{Z}_{\text{in}} \mathbb{1}_V\|_2 \leq 2 \text{ and } \frac{1}{\sqrt{V}} \|\mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V\|_\infty \leq \frac{\log V}{\sqrt{d}}$$

$$(E1.3) \quad \left| \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right| \leq 2 \log V \sqrt{\frac{V}{d}} \text{ and } \left| \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right| \leq C_K \left( \frac{V}{d} \right)^{\frac{3}{2}} \text{ and } \left| \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \text{diag}(\mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}}) \right| \leq C_K \log V \sqrt{\frac{V}{d}}$$

$$(E1.4) \text{ For all } i \in [N], \|\mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L\| \leq \mathbf{e}_\nu^\top \mathbf{X}_i^\top \mathbb{1}_L + C_K \log V \frac{\|\mathbf{X}_i^\top \mathbb{1}_L\|_2}{\sqrt{d}}$$

$$(E1.5) \text{ For all } i \in [N], \|\mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L\| \leq L + C_K \log V \|\mathbf{X}_i^\top \mathbb{1}_L\|_2 \sqrt{\frac{V}{d}}.$$

$$(E1.6) \text{ For all } i \in [N], |\mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L| \leq \frac{V}{d} (\mathbf{e}_\nu^\top \mathbf{X}_i^\top \mathbb{1}_L + C_K \log V \frac{\|\mathbf{X}_i^\top \mathbb{1}_L\|_2}{\sqrt{d}}).$$

$$(E1.7) \text{ For all } i, j \in [N], |\frac{1}{L} \mathbb{1}_L^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L - m_{ij}| \leq |\frac{1}{L} \mathbb{1}_L^\top \mathbf{X}_j \mathbf{X}_i^\top \mathbb{1}_L - m_{ij}| + C_K \frac{\|\mathbf{X}_i^\top \mathbb{1}_L\|_2 \|\mathbf{X}_j^\top \mathbb{1}_L\|_2 \log V}{L \sqrt{d}},$$

$$(E1.8) \text{ For all } i \in [N], \|\mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2 \leq C_K (\mathbf{e}_\nu^\top \mathbf{N}_i^\top \mathbb{1}_{L-1} + \frac{L}{d} + \log^6 V \frac{\|\mathbf{N}_i^\top \mathbb{1}_{L-1}\|_2}{\sqrt{d}}).$$

(E.2) We have

$$(E2.1) \text{ For all } i, j \in [N], |\frac{1}{L} \mathbb{1}_L^\top \mathbf{X}_j \mathbf{X}_i^\top \mathbb{1}_L - m_{ij}| \leq C_K \frac{\log^2 V}{\sqrt{V} \wedge L},$$

$$(E2.2) \text{ For all } i \in [N], \|\mathbf{X}_i^\top \mathbb{1}_L\|_\infty \leq \log L \text{ and } \|\mathbf{X}_i^\top \mathbb{1}_L\|_0 \geq \frac{L}{2}$$

$$(E2.3) \left| \left\| \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \right\|_2 - \frac{1}{N} - \frac{1}{V} \right| \leq C_K \frac{\log^2 N}{N \sqrt{V}} \text{ and } \left| \left\| \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i - \frac{1}{V} \mathbb{1}_V \right\|_2 - \frac{1}{N} \right| \leq C_K \frac{\log^2 N}{N \sqrt{V}}$$

$$\text{and } \left\| \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i - \frac{1}{V} \mathbb{1}_V \right\|_\infty \leq \frac{(e+1)L}{V}$$

$$(E2.4) \sum_{i,j=1}^N |\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}| \leq \frac{4N^2}{V} \text{ and } \sum_{i,j=1}^N (\mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V}) \leq \frac{4N^2}{V}$$

$$(E2.5) \|\mathbf{S}_1\|_2 \leq \frac{e}{L^2 V^2} \text{ and } |\text{tr}(\mathbf{S}_1) - \frac{1}{L^2} (\frac{1}{N} + (1 - \frac{1}{V}) \frac{1}{V})| \leq \frac{C_K \log^2 V}{L^2 N \sqrt{V}}$$

$$(E2.6) \|\mathbf{S}_2\|_2 \leq \frac{C_K \log^2 V}{N L V} \text{ and } |\text{tr}(\mathbf{S}_2) - (1 - \frac{1}{V})^2 \frac{L-1}{L^2 N}| \leq \frac{K C_K \log^3 V}{N \sqrt{L V}}$$

$$(E2.7) \frac{-C_K \log^2 V}{N \sqrt{V}} \frac{1}{V^2 L^2} \mathbb{1}_V \mathbb{1}_V^\top \preceq \mathbf{S}_3 - \frac{1}{N} \frac{1}{V^2 L^2} \mathbb{1}_V \mathbb{1}_V^\top \preceq \frac{C_K \log^2 V}{N \sqrt{V}} \frac{1}{V^2 L^2} \mathbb{1}_V \mathbb{1}_V^\top$$

For any  $K > 0$ , there exists a universal constant  $C_K > 0$  depending only on  $K$  such that

$$\mathbb{P}[(E.1)] \geq 1 - \frac{1}{V^K} \text{ and } \mathbb{P}[(E.2)] \geq 1 - \frac{1}{V^K}.$$

*Proof.* For (E.1):

- By Proposition 3, we have  $\|\frac{1}{V} \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top - \frac{1}{d} \mathbf{I}_d\|_2 \leq \frac{2 \log V}{\sqrt{V} d}$  and by Proposition 4, we have  $\max_{k \leq V} \|\mathbf{z}_k\|_2 \vee \|\mathbf{z}_{\text{trig}}\|_2 \leq 2$  with probability at least  $1 - CVd \exp(-c \log^2 V)$ .

- By Proposition 4,  $\frac{1}{\sqrt{V}} \|\mathbf{Z}_{\text{in}} \mathbb{1}_V\|_2 \leq 2$  and  $\frac{1}{\sqrt{V}} \|\mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V\|_\infty \leq \frac{2 \log V}{\sqrt{d}}$  with probability at least  $1 - CVd \exp(-c \log^2 V)$ .

- By Propositions 4 and 5, we have  $\frac{1}{\sqrt{V}} |\mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V| \leq \frac{2 \log V}{\sqrt{d}}$  and  $|\mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V| \leq C_K (\frac{V}{d})^{\frac{3}{2}}$  with probability at least  $1 - CVd \exp(-c \log^2 V)$ . Moreover

$$\frac{1}{V} |\mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \text{diag}(\mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}})| = \frac{1}{V} \sum_{\substack{i=1 \\ i \neq \nu}}^V \|\mathbf{z}_i\|_2^2 \langle \mathbf{z}_i, \mathbf{z}_\nu \rangle + \frac{\mathbb{1}_{l=1} \delta}{V} \sum_{\substack{i=1 \\ i \neq \nu}}^V \|\mathbf{z}_i\|_2^2 \langle \mathbf{z}_i, \mathbf{z}_{\text{trig}} \rangle + \underbrace{\frac{1}{V} \mathbf{z}_{\nu,\delta}^\top \mathbf{z}_\nu}_{\in \frac{1}{V} [-C_K, C_K]},$$

where we used previous items to bound the last term. For  $i \neq k$ , by using Lemma 3, we have for  $p \leq \frac{d}{6}$ ,

$$\begin{aligned} \mathbb{E}[\|\mathbf{z}_i\|_2^{4p} |\langle \mathbf{z}_i, \mathbf{z}_k \rangle|^{2p}] &\leq d^{-p} \mathbb{E}[\|\mathbf{z}_i\|_2^{6p}] (2p)^p \\ &\leq d^{-p} 2^p p^p \frac{d(d+2) \cdots (d+6p-2)}{d^{3p}} \leq d^{-p} 2^{4p} p^p \end{aligned}$$

Therefore,

$$\mathbb{E}[\|\mathbf{z}_i\|_2^{4p} |\langle \mathbf{z}_i, \mathbf{z}_\nu \rangle|^{2p}]^{\frac{1}{2p}} \leq 4d^{-1/2} \sqrt{p}.$$

By Proposition 13, we have for  $2 \leq p \leq \frac{d}{6}$ ,

$$\mathbb{E} \left[ \left| \frac{1}{V} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \text{diag}(\mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}}) \right|^{2p} \right]^{\frac{1}{2p}} \leq C d^{-1/2} \left[ \sqrt{\frac{p}{V}} + V^{\frac{1}{p}} \frac{p^{3/2}}{V} \right]$$

By using  $p = \log V$ , we have the bound in the statement with probability  $1 - \frac{1}{V^K}$ .

- By Proposition 4 with probability at least  $1 - \frac{1}{V^K}$

$$\begin{aligned} |z_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L| &\leq |e_\nu^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L| + \mathbb{1}_{l=1} \delta |z_{\text{trig}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L| \\ &\leq e_\nu^\top \mathbf{X}_i^\top \mathbb{1}_L + C_K \log V \frac{\|\mathbf{X}_i^\top \mathbb{1}_L\|_2}{\sqrt{d}}. \end{aligned}$$

By union bound, the item follows.

- By the same argument,

$$\begin{aligned} |\mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L| &\leq \mathbb{1}_V^\top \mathbf{X}_i^\top \mathbb{1}_L + C_K \log V \|\mathbf{X}_i^\top \mathbb{1}_L\|_2 \sqrt{\frac{V}{d}} \\ &= L + C_K \log V \|\mathbf{X}_i^\top \mathbb{1}_L\|_2 \sqrt{\frac{V}{d}} \end{aligned}$$

- By Proposition 5, with probability at least  $1 - CN \exp(-c \log^2 V)$ , we have  $|z_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L| \leq \frac{V}{d} (e_\nu^\top \mathbf{X}_i^\top \mathbb{1}_L + C_K \log V \frac{\|\mathbf{X}_i^\top \mathbb{1}_L\|_2}{\sqrt{d}})$  for all  $i \in [N]$ .

- By Proposition 4, with probability at least  $1 - \frac{1}{V^K}$

$$\begin{aligned} \frac{1}{L} \mathbb{1}_L^\top \mathbf{X}_j \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L - m_{ij} &= \frac{1}{L} \mathbb{1}_L^\top \mathbf{X}_j \mathbf{X}_i^\top \mathbb{1}_L - m_{ij} \\ &\quad \pm C \log V \frac{\|\mathbf{X}_j^\top \mathbb{1}_L\|_2 \|\mathbf{X}_i^\top \mathbb{1}_L\|_2 \log V}{L \sqrt{d}}. \end{aligned}$$

- For the last item, let  $n_k := e_k^\top \mathbf{Z}_{\text{in}}$ . We have

$$\mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} = n_\nu (\|z_\nu\|_2^2 + \delta \mathbb{1}_{l=1} z_\nu^\top z_\delta - \frac{1}{d}) z_\nu + \frac{L}{d} z_\nu + \sum_{\substack{k=1 \\ k \neq \nu}}^V n_k (z_k z_k^\top - \frac{1}{d} \mathbf{I}_d) z_{\nu,\delta}.$$

By Proposition 11, we have

$$\begin{aligned} \mathbb{E} \left[ \left\| \sum_{\substack{k=1 \\ k \neq \nu}}^V n_k (z_k z_k^\top - \frac{1}{d} \mathbf{I}_d) z_{\nu,\delta} \right\|_2^{2p} \right]^{\frac{1}{p}} &\leq C(p-1)^6 \mathbb{E} \left[ \left\| \sum_{\substack{k=1 \\ k \neq \nu}}^V n_k (z_k z_k^\top - \frac{1}{d} \mathbf{I}_d) z_{\nu,\delta} \right\|_2^2 \right] \\ &\leq \frac{C}{d} (p-1)^6 \|\mathbf{N}_i^\top \mathbb{1}_{L-1}\|_2^2 \end{aligned}$$

Therefore, with probability  $1 - \frac{1}{V^K}$ , we have

$$\|\mathbf{Z}_{\text{in}} \mathbf{N}_i^\top \mathbf{N}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2 \leq C_K \left( n_\nu + \frac{L}{d} + \log^6 V \frac{\|\mathbf{N}_i^\top \mathbb{1}_{L-1}\|_2}{\sqrt{d}} \right).$$

For (E.2):

- By Proposition 7, we have the first item with probability  $1 - \frac{N^2}{V^K}$ .
- By Corollary 3, we have  $\|\mathbf{X}_i^\top \mathbb{1}_L\|_\infty \leq \log L$ . For the second part, we define  $n_k := e_k^\top \mathbf{X}_i^\top \mathbb{1}_L$ . We observe that

$$\mathbb{E}[\|\mathbf{X}_i^\top \mathbb{1}_L\|_0] = \sum_{k=1}^V \mathbb{P}[n_k > 0] = V(1 - (1 - \frac{1}{V})^L) = L(1 - \frac{L}{2V} + o(L/V)).$$

By McDiarmid inequality, we have

$$\mathbb{P} \left[ \left| \|\mathbf{X}_i^\top \mathbb{1}_L\|_0 - L(1 - \frac{L}{2V} + o(L/V)) \right| > \sqrt{L} \log V \right] \leq 2 \exp(-2 \log^2 V),$$

which gives the result.

- Let  $\mathbf{n} = \sum_{i=1}^N \mathbf{x}_i$ . We have  $\mathbb{E}[\mathbf{n}] = \frac{N}{V} \mathbb{1}_V$  and by Proposition 7 with probability  $1 - \frac{1}{V^K}$ , we have

$$\left| \left\| \frac{1}{N} \mathbf{n} - \frac{1}{V} \mathbb{1}_V \right\|_2^2 - \left(1 - \frac{1}{V}\right) \frac{1}{N} \right| = \left| \left\| \frac{1}{N} \mathbf{n} \right\|_2^2 - \left(1 - \frac{1}{V}\right) \frac{1}{N} - \frac{1}{V} \right| \leq C_K \frac{\log^2 V}{N\sqrt{V}}.$$

Lastly, by Corollary 3, we have  $\left\| \frac{1}{N} \mathbf{n} - \frac{1}{V} \mathbb{1}_V \right\|_\infty \leq \frac{(e+1)L}{V}$ .

- We have

$$\begin{aligned} \sum_{i,j=1}^N \left| \mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V} \right| &= \left( \sum_{i,j=1}^N \left| \mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V} \right| - \frac{2}{V} \left(1 - \frac{1}{V}\right) \right) + \frac{2N^2}{V} \left(1 - \frac{1}{V}\right) \\ &= \left(1 - \frac{2}{V}\right) \sum_{i,j=1}^N \left( \mathbb{1}_{\mathbf{x}_i=\mathbf{x}_j} - \frac{1}{V} \right) + \frac{2N^2}{V} \left(1 - \frac{1}{V}\right) \\ &= \left(1 - \frac{2}{V}\right) \left\| \sum_{i=1}^N \left( \mathbf{x}_i - \frac{1}{V} \mathbb{1}_V \right) \right\|_2^2 + \frac{2N^2}{V} \left(1 - \frac{1}{V}\right) \end{aligned}$$

By the previous item, the statement follows

- The events for  $\mathcal{S}_1$ ,  $\mathcal{S}_2$  and  $\mathcal{S}_3$  follows Proposition 8.

□

**Proposition 2.** We consider the parameter regime in Lemma 2. Let  $\bar{\phi} := \sup_{k_1, k_2 \geq 1} |\phi^{(k_1)}(0)\phi^{(k_2)}(0)|$ . The intersection of (E.1) and (E.2) implies the following events:

$$(C.1) \text{ For all } i, j \in [N], \left| \frac{1}{L} \mathbb{1}_L \mathbf{X}_j^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L - m_{ij} \right| \leq C_K \left( \frac{\log V}{\sqrt{d}} + \frac{\log^2 V}{L} \right),$$

$$(C.2) \sup_{i,j} |\alpha_{ij} - \phi'(0)^2| \vee |\beta_{ij} - \phi''(0)\phi(0)| \leq \frac{\bar{\phi}}{L} (m_{ij} + C_K \frac{\log V}{\sqrt{d}} + C_K \frac{\log^2 V}{L})$$

$$(C.3) \text{ Let } \Delta_{u,ir} := \mathbf{A}_{u,ir} - \phi'(0)^4 \mathbf{Z}_{\text{in}} \mathbf{S}_u \mathbf{Z}_{\text{in}}^\top \text{ for } u \in \{1, 2, 3\}. \text{ We have}$$

$$- \sup_{i,r \in [N]} \|\Delta_{1,ir}\|_2 \leq C_K \phi'(0)^2 \left( \frac{1}{NdL^3} + \frac{1}{VdL^2} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right).$$

$$- \sup_{i,r \in [N]} \|\Delta_{2,ir}\|_2 \leq \frac{C_K \sqrt{V}}{d\sqrt{NL}} \left( \frac{1}{NL^{\frac{3}{2}}} + \frac{1}{V\sqrt{L}} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right).$$

$$- \text{ We have } \Delta_{3,ir} = \frac{\bar{\Delta}_{3,ir}}{V^2 L^2} \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \text{ such that}$$

$$\sup_{i,r \in [N]} |\bar{\Delta}_{3,ir}| \leq \frac{C_K \phi'(0)^2}{N} \left( \frac{1}{NL} + \frac{1}{\sqrt{N}} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right) + \left( \frac{1}{NL} + \frac{1}{\sqrt{N}} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right)^2.$$

$$(C.4) \text{ For all } i, r \in [N],$$

- We have

$$\begin{aligned} \left\| \mathbf{A}_{1,ir} - \frac{\phi'(0)^4}{d} \left( \frac{1}{N} + \left(1 - \frac{1}{V}\right) \frac{1}{V} \right) \mathbf{I}_d \right\|_2 &\leq C_K \phi'(0)^2 \left( \frac{1}{NdL^3} + \frac{1}{VdL^2} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right) \\ &\quad + C_K \phi'(0)^4 \left( \frac{\log V}{L^2 V^{3/2} \sqrt{d}} + \frac{\log^2 V}{L^2 N \sqrt{Vd}} \right). \end{aligned}$$

- We have

$$\begin{aligned} \left\| \mathbf{A}_{2,ir} - \frac{\phi'(0)^4}{d} \left(1 - \frac{1}{V}\right)^2 \frac{L-1}{L^2 N} \mathbf{I}_d \right\|_2 &\leq \frac{C_K \sqrt{V}}{d\sqrt{NL}} \left( \frac{1}{NL^{\frac{3}{2}}} + \frac{1}{V\sqrt{L}} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right) \\ &\quad + C_K \phi'(0)^4 \left( \frac{\log V}{NL\sqrt{Vd}} + \frac{\log^3 V}{N\sqrt{LVd}} \right). \end{aligned}$$

- We have  $A_{3,ir} - \frac{\phi'(0)^4}{N} \frac{1}{V^2 L^2} \mathbf{Z}_{in} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{in}^\top =: \tilde{\Delta}_{3,ir} \mathbf{Z}_{in} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{in}^\top$  such that

$$|\tilde{\Delta}_{3,ir}| \leq \frac{C_K \phi'(0)^4 \log^2 V}{N \sqrt{V}} + \frac{C_K \phi'(0)^2}{N} \left( \frac{1}{NL} + \frac{1}{\sqrt{N}} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}} \right) + \left( \frac{1}{NL} + \frac{1}{\sqrt{N}} \frac{1}{V \wedge L^2 \wedge L \sqrt{d}} \right)^2.$$

*Proof.* We have the following arguments.

- By (E1.7) and (E2.1), we have (C.1).
- For (C.2), we assume (E2.1) and (C.1) hold. Let  $\|\frac{1}{L} \mathbf{Z}_{in} \mathbf{X}_i^\top \mathbb{1}_L\|_2 w_i \leftarrow \frac{1}{L} \mathbf{w}^\top \mathbf{Z}_{in} \mathbf{X}_i^\top \mathbb{1}_L$ . We write

$$\begin{aligned} & \left| \mathbb{E} \left[ \phi' \left( \left\| \frac{1}{L} \mathbf{Z}_{in} \mathbf{X}_i^\top \mathbb{1}_L \right\|_2 w_i \right) \phi' \left( \left\| \frac{1}{L} \mathbf{Z}_{in} \mathbf{X}_i^\top \mathbb{1}_L \right\|_2 w_j \right) \right] - \phi'(0)^2 \right| \\ &= \left| \sum_{u,v=1}^{p_*} \left\| \frac{1}{L} \mathbf{Z}_{in} \mathbf{X}_i^\top \mathbb{1}_L \right\|_2^u \left\| \frac{1}{L} \mathbf{Z}_{in} \mathbf{X}_j^\top \mathbb{1}_L \right\|_2^v \frac{\mathbb{E}[w_i^u w_j^v]}{u!v!} \phi^{(u+1)}(0) \phi^{(v+1)}(0) \right| \\ &= \left| \frac{1}{L^2} \mathbb{1}_L^\top \mathbf{X}_j \mathbf{Z}_{in}^\top \mathbf{Z}_{in} \mathbf{X}_i^\top \mathbb{1}_L \phi^{(2)}(0) \phi^{(2)}(0) \right. \\ &\quad \left. + \sum_{\substack{u,v=1 \\ u+v \text{ is even} \\ u+v > 2}}^{p_*} \left\| \frac{1}{L} \mathbf{Z}_{in} \mathbf{X}_i^\top \mathbb{1}_L \right\|_2^u \left\| \frac{1}{L} \mathbf{Z}_{in} \mathbf{X}_j^\top \mathbb{1}_L \right\|_2^v \frac{\mathbb{E}[w_i^u w_j^v]}{u!v!} \phi^{(u+1)}(0) \phi^{(v+1)}(0) \right| \\ &\leq \frac{\bar{\phi}}{L} (m_{ij} + C_K \frac{\log V}{\sqrt{d}} + C_K \frac{\log^2 V}{\sqrt{V} \wedge L}) + O(\frac{1}{L^2}) \end{aligned}$$

Similarly,

$$\begin{aligned} & \left| \mathbb{E} \left[ \phi'' \left( \left\| \frac{1}{L} \mathbf{Z}_{in} \mathbf{X}_i^\top \mathbb{1}_L \right\|_2 w_i \right) \phi \left( \left\| \frac{1}{L} \mathbf{Z}_{in} \mathbf{X}_i^\top \mathbb{1}_L \right\|_2 w_j \right) \right] - \phi'(0)^2 \right| \\ &= \left| \sum_{u,v=1}^{k_2} \left\| \frac{1}{L} \mathbf{Z}_{in} \mathbf{X}_i^\top \mathbb{1}_L \right\|_2^u \left\| \frac{1}{L} \mathbf{Z}_{in} \mathbf{X}_j^\top \mathbb{1}_L \right\|_2^v \frac{\mathbb{E}[w_i^u w_j^v]}{u!v!} \phi^{(u+2)}(0) \phi^{(v)}(0) \right| \\ &= \left| \frac{1}{L^2} \mathbb{1}_L^\top \mathbf{X}_j \mathbf{Z}_{in}^\top \mathbf{Z}_{in} \mathbf{X}_i^\top \mathbb{1}_L \phi^{(2)}(0) \phi^{(2)}(0) \right. \\ &\quad \left. + \sum_{\substack{u,v=1 \\ u+v \text{ is even} \\ u+v > 2}}^{k_2} \left\| \frac{1}{L} \mathbf{Z}_{in} \mathbf{X}_i^\top \mathbb{1}_L \right\|_2^u \left\| \frac{1}{L} \mathbf{Z}_{in} \mathbf{X}_j^\top \mathbb{1}_L \right\|_2^v \frac{\mathbb{E}[w_i^u w_j^v]}{u!v!} \phi^{(u+2)}(0) \phi^{(v)}(0) \right| \\ &\leq \frac{\bar{\phi}}{L} (m_{ij} + C_K \frac{\log V}{\sqrt{d}} + C_K \frac{\log^2 V}{\sqrt{V} \wedge L}) + O(\frac{1}{L^2}) \end{aligned}$$

- For (C.3), we assume (E1.1), (E1.2) and (E2.5)-??. We define

$$\begin{aligned} \bar{\Delta}_{1,ir} &:= \left( \frac{1}{LN} \sum_{j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\ &\quad \times \left( \frac{1}{LN} \sum_{j=1}^N \phi'(0)^2 (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\ &\quad + \left( \frac{1}{LN} \sum_{j=1}^N \phi'(0)^2 (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\ &\quad \times \left( \frac{1}{LN} \sum_{j=1}^N (\alpha_{uj} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \end{aligned}$$

$$\begin{aligned}
& + \left( \frac{1}{LN} \sum_{j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\
& \quad \times \left( \frac{1}{LN} \sum_{j=1}^N (\alpha_{uj} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right)
\end{aligned}$$

We have

$$\begin{aligned}
\|\bar{\Delta}_{1,ir}\|_2 & \leq \frac{C\phi'(0)^2 \sup_i |\alpha_{ii} - \phi'(0)^2|}{LN} \|\mathbf{S}_1\|_2^{\frac{1}{2}} + \phi'(0)^2 \sup_{i \neq j} |\alpha_{ij} - \phi'(0)^2| \|\mathbf{S}_1\|_2 \\
& \leq C\phi'(0)^2 \left( \frac{1}{NV L^3} + \frac{1}{V^2 L^2} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right).
\end{aligned}$$

Therefore,

$$\|\Delta_{1,ir}\|_2 = \|\mathbf{Z}_{in} \bar{\Delta}_{1,ir} \mathbf{Z}_{in}^\top\|_2 \leq C\phi'(0)^2 \left( \frac{1}{NdL^3} + \frac{1}{VdL^2} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right).$$

Moreover, we define

$$\begin{aligned}
\bar{\Delta}_{2,ir} & := \left( \frac{1}{NL} \sum_{j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\
& \quad \left( \frac{1}{NL} \sum_{j=1}^N \phi'(0)^2 (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right) \\
& + \left( \frac{1}{NL} \sum_{j=1}^N \phi'(0)^2 (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\
& \quad \left( \frac{1}{NL} \sum_{j=1}^N (\alpha_{rj} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right) \\
& + \left( \frac{1}{NL} \sum_{j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V)^\top \right) \\
& \quad \left( \frac{1}{NL} \sum_{j=1}^N (\alpha_{rj} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right).
\end{aligned}$$

We have

$$\begin{aligned}
& \|\bar{\Delta}_{2,ir}\|_2 \\
& \leq \phi'(0)^2 \|\mathbf{S}_2\|_2^{\frac{1}{2}} \left\| \left( \frac{1}{NL} \sum_{j=1}^N (\alpha_{rj} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right) \right\|_2 \\
& + \phi'(0)^2 \|\mathbf{S}_2\|_2^{\frac{1}{2}} \left\| \left( \frac{1}{NL} \sum_{j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right) \right\|_2 \\
& + \left\| \left( \frac{1}{NL} \sum_{j=1}^N (\alpha_{rj} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right) \right\|_2 \\
& \quad \times \underbrace{\left\| \left( \frac{1}{NL} \sum_{j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right) \right\|_2}_{\leq \frac{C}{NL\sqrt{L}} + \frac{C}{V\sqrt{L}} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}}} \\
& \leq \frac{C}{\sqrt{NV L}} \left( \frac{1}{NL^{\frac{3}{2}}} + \frac{1}{V\sqrt{L}} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right) + C^2 \left( \frac{1}{NL^{\frac{3}{2}}} + \frac{1}{V\sqrt{L}} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right)^2.
\end{aligned}$$

Therefore,

$$\|\Delta_{2,ir}\|_2 = \|\mathbf{Z}_{\text{in}} \bar{\Delta}_{2,ir} \mathbf{Z}_{\text{in}}^\top\|_2 \leq \frac{C\sqrt{V}}{d\sqrt{NL}} \left( \frac{1}{NL^{\frac{3}{2}}} + \frac{1}{V\sqrt{L}} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right).$$

Lastly, we define

$$\begin{aligned} \bar{\Delta}_{3,ir} := & \left( \frac{1}{N} \sum_{j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right)^\top \left( \frac{1}{N} \sum_{j=1}^N \phi'(0)^2 (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right) \\ & + \left( \frac{1}{N} \sum_{j=1}^N \phi'(0)^2 (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right)^\top \left( \frac{1}{N} \sum_{j=1}^N (\alpha_{rj} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right) \\ & + \left( \frac{1}{N} \sum_{j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right)^\top \left( \frac{1}{N} \sum_{j=1}^N (\alpha_{rj} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right) \end{aligned}$$

We have

$$\begin{aligned} |\bar{\Delta}_{3,ir}| & \leq \phi'(0)^2 \left\| \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \left\| \frac{1}{N} \sum_{j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \\ & \quad + \phi'(0)^2 \left\| \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \left\| \frac{1}{N} \sum_{j=1}^N (\alpha_{rj} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \\ & \quad + \underbrace{\left\| \frac{1}{N} \sum_{j=1}^N (\alpha_{ij} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2 \left\| \frac{1}{N} \sum_{j=1}^N (\alpha_{rj} - \phi'(0)^2) (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2}_{\leq \frac{C}{NL} + \frac{C}{\sqrt{N}} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}}} \\ & \leq \frac{C\phi'(0)^2}{N} \left( \frac{1}{NL} + \frac{1}{\sqrt{N}} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right) + \left( \frac{1}{NL} + \frac{1}{\sqrt{N}} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right)^2. \end{aligned}$$

Therefore,

$$\Delta_{3,ir} = \frac{\bar{\Delta}_{3,ir}}{V^2 L^2} \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top,$$

from which the last result follows.

- For (C.4), we assume (E1.1), (E1.2), (E2.3), and (E2.5)-(E2.7). We write

$$\begin{aligned} \mathbf{A}_{1,ir} - \frac{\phi'(0)^4}{d} \frac{1}{L^2} \left( \frac{1}{N} + (1 - \frac{1}{V}) \frac{1}{V} \right) \mathbf{I}_d \\ = \Delta_{1,ir} + \phi'(0)^4 \left( \mathbf{Z}_{\text{in}} \mathbf{S}_1 \mathbf{Z}_{\text{in}}^\top \pm \frac{\text{tr}(\mathbf{S}_1)}{d} \mathbf{I}_d - \frac{1}{d} \left( (1 - \frac{1}{V}) \frac{1}{N} + (1 - \frac{1}{V}) \frac{1}{V} \right) \mathbf{I}_d \right). \end{aligned}$$

We have

$$\begin{aligned} \|\mathbf{A}_{1,ir} - \frac{\phi'(0)^4}{d} \left( (1 - \frac{1}{V})^2 \frac{1}{N} + (1 - \frac{1}{V}) \frac{1}{V} \right) \mathbf{I}_d\|_2 \\ \leq \|\Delta_{1,ir}\|_2 + 2\phi'(0)^4 \log V \frac{\|\mathbf{S}_1\|_F}{\sqrt{d}} + \frac{|\text{tr}(\mathbf{S}_1) - \frac{1}{L^2} (\frac{1}{N} + (1 - \frac{1}{V}) \frac{1}{V})|}{d} \\ \leq C\phi'(0)^2 \left( \frac{1}{NdL^3} + \frac{1}{VdL^2} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right) + C\phi'(0)^4 \left( \frac{\log V}{L^2 V^{3/2} \sqrt{d}} + \frac{CK^2 \log^2 V}{L^2 N \sqrt{Vd}} \right). \end{aligned}$$

Moreover,

$$\begin{aligned} \|\mathbf{A}_{2,ir} - \frac{\phi'(0)^4}{d} (1 - \frac{1}{V})^2 \frac{L-1}{L^2 N} \mathbf{I}_d\|_2 \\ \leq \|\Delta_{2,ir}\|_2 + 2\phi'(0)^4 \log V \frac{\|\mathbf{S}_2\|_F}{\sqrt{d}} + \frac{|\text{tr}(\mathbf{S}_2) - (1 - \frac{1}{V})^2 \frac{L-1}{L^2 N}|}{d} \end{aligned}$$



$$\leq \frac{C\sqrt{V}}{d\sqrt{NL}} \left( \frac{1}{NL^{\frac{3}{2}}} + \frac{1}{V\sqrt{L}} \frac{1}{V \wedge L^2 \wedge L\sqrt{d}} \right) + C\phi'(0)^4 \left( \frac{\log V}{NL\sqrt{Vd}} + \frac{K^{\frac{3}{2}} \log^3 V}{N\sqrt{LVd}} \right).$$

Lastly,

$$\begin{aligned} \mathbf{A}_{3,ir} - \frac{\phi'(0)^4}{N} \frac{1}{V^2 L^2} \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \\ = \Delta_{3,ir} + \phi'(0)^4 \left( \left\| \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2^2 - \frac{1}{N} \right) \frac{1}{V^2 L^2} \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top. \end{aligned}$$

By (E2.3), we have

$$\begin{aligned} \frac{CK^2 \log^2 V}{N\sqrt{V}} \frac{1}{V^2 L^2} \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top &\preceq \left( \left\| \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2^2 - \frac{1}{N} \right) \frac{1}{V^2 L^2} \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \\ &\preceq \frac{CK^2 \log^2 V}{N\sqrt{V}} \frac{1}{V^2 L^2} \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top. \end{aligned}$$

By (C.3), the result follows.  $\square$

## E.2 GAUSSIAN MATRICES AND RELATED STATEMENTS

**Lemma 3.** Let  $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I}_d)$ . We have  $\mathbb{E}[\|\mathbf{z}\|_2^{2k}] = d(d+2) \cdots (d+2k-2)$ .

*Proof.* We observe that  $\|\mathbf{z}\|_2 \sim \chi_d^2$ . By using the moment formula for chi-squared distribution, we have the result.  $\square$

**Lemma 4.** Let  $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I}_d)$  and  $\mathbf{S} \in \mathbb{R}^{d \times d}$  be a symmetric matrix. For  $u > 0$ ,

$$\mathbb{P} [|\mathbf{z}^\top \mathbf{S} \mathbf{z} - \text{tr}(\mathbf{S})| \geq 2\|\mathbf{S}\|_F u + 2\|\mathbf{S}\|_2 u^2] \leq 2e^{-u^2}.$$

*Proof.* We note that  $\mathbf{z}^\top \mathbf{S} \mathbf{z} - \text{tr}(\mathbf{S})$  has the same distribution with  $\sum_{i=1}^d \lambda_i(\mathbf{S})(Z_i^2 - 1)$ , where  $Z_i \sim_{iid} \mathcal{N}(0, 1)$ . By using the Laurent-Massart lemma, we have the result.  $\square$

**Proposition 3.** Let  $\mathbf{S} \in \mathbb{R}^{V \times V}$  be a symmetric positive semidefinite matrix. Let

$$\mathbf{M} = \mathbf{Z}_{\text{in}} \mathbf{S} \mathbf{Z}_{\text{in}}^\top.$$

For  $\text{poly}(d) \gg V \gg d$ , We have

$$\mathbb{P} \left[ \left\| \mathbf{M} - \frac{\text{tr}(\mathbf{S})}{d} \mathbf{I}_d \right\|_2 \geq \max \left\{ \frac{\|\mathbf{S}\|_F}{\sqrt{d}} \log V, \|\mathbf{S}\|_2 \log^2 V \right\} \right] \leq \exp(-c \log^2 V).$$

*Proof.* Without loss of generality, we can assume that  $\mathbf{S}$  is diagonal, i.e.,  $\mathbf{S} = \text{diag}(s_1, \dots, s_V)$ . We have

$$\mathbf{M} - \frac{\text{tr}(\mathbf{S})}{d} \mathbf{I}_d = \sum_{i=1}^V s_i (\mathbf{z}_i \mathbf{z}_i^\top - \frac{1}{d} \mathbf{I}_d).$$

We have

$$\mathbb{E} \left[ \left( \sum_{i=1}^V s_i (\mathbf{z}_i \mathbf{z}_i^\top - \frac{1}{d} \mathbf{I}_d) \right)^2 \right] = \frac{1}{d} \left( 1 + \frac{1}{d} \right) \|\mathbf{S}\|_F^2 \mathbf{I}_d$$

Moreover, for  $p \leq \frac{d}{2}$

$$\mathbb{E} \left[ \left\| \mathbf{z}_i \mathbf{z}_i^\top - \frac{1}{d} \mathbf{I}_d \right\|_2^p \right] \leq \mathbb{E}[\|\mathbf{z}_i\|_2^{2p}] \leq 2^p.$$

By Proposition 13, we have  $2 \leq p \leq \frac{d}{2}$

$$\mathbb{E} \left[ \left\| \mathbf{M} - \frac{\text{tr}(\mathbf{S})}{d} \right\|_2^p \right] \leq C \left( \sqrt{p \vee \log d} \frac{\|\mathbf{S}\|_F}{\sqrt{d}} + (p \vee \log d) V^{\frac{1}{p}} \|\mathbf{S}\|_2 \right).$$

For  $p = \frac{1}{e^2 C^2} \log^2 V$ , we have the result.  $\square$

**Proposition 4.** Let  $\mathbf{S} \in \mathbb{R}^{V \times V}$  be a square matrix. For  $\mathbf{u}, \mathbf{v} \in \mathbb{S}^{d-1}$  and  $\mathbf{M} = \mathbf{Z}_{\text{in}} \mathbf{S} \mathbf{Z}_{\text{in}}^\top$ , we have

$$\begin{aligned} \mathbb{P} \left[ \left| \left( \mathbf{v}^\top \mathbf{M} \mathbf{u} - \frac{\text{tr}(\mathbf{S})}{d} \mathbf{v}^\top \mathbf{u} \right) \right| \geq \frac{\|\mathbf{u}\|_2 \|\mathbf{v}\|_2}{d} \max \left\{ \|\text{sym}(\mathbf{S})\|_{Ft}, \|\text{sym}(\mathbf{S})\|_2 t^2 \right\} \right] \\ \leq 2 \exp(-ct^2). \end{aligned}$$

*Proof.* Consider  $\mathbf{g} = \sqrt{d} \text{vec}(\mathbf{Z})$ , where  $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_{dV})$ . We have

$$\mathbf{v}^\top \mathbf{M} \mathbf{u} = \frac{1}{d} \mathbf{g}^\top (\mathbf{u} \mathbf{v}^\top) \otimes \mathbf{S} \mathbf{g} = \frac{1}{d} \mathbf{g}^\top \text{sym}(\mathbf{u} \mathbf{v}^\top) \otimes \text{sym}(\mathbf{S}) \mathbf{g}$$

By using Proposition 10, we have

$$\mathbb{E}[\mathbf{g}^\top \text{sym}(\mathbf{u} \mathbf{v}^\top) \otimes \text{sym}(\mathbf{S}) \mathbf{g}] = \text{tr}(\mathbf{S}) \mathbf{u}^\top \mathbf{v}.$$

Moreover,

$$\left( \mathbf{g}^\top \text{sym}(\mathbf{u} \mathbf{v}^\top) \otimes \text{sym}(\mathbf{S}) \mathbf{g} - \text{tr}(\mathbf{S}) \mathbf{u}^\top \mathbf{v} \right) =_d \sum_{i=1}^{dV} \lambda_i (g_i^2 - 1)$$

where  $g_i \sim N(0, 1)$ . By using the subexponential concentration, we have the result.  $\square$

**Proposition 5.** For  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^V$ , we have

$$\begin{aligned} \mathbb{P} \left[ \left| \mathbf{v}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{u} - \mathbf{u}^\top \mathbf{v} \left( 1 + \frac{V-1}{d} \right) \right| \geq C \|\mathbf{u}\|_2 \|\mathbf{v}\|_2 \log V \left( \frac{\sqrt{V}}{d} + \frac{V}{d^{3/2}} \right) \right] \\ \leq 10 \exp(-c \log^2 V). \end{aligned}$$

*Proof.* Without loss of generality, we assume that  $\mathbf{u}$  and  $\mathbf{v}$  have a unit norm. Let

$$\mathbf{v}_\perp := \frac{1}{\sqrt{1 - (\mathbf{u}^\top \mathbf{v})^2}} (\mathbf{I}_V - \mathbf{v} \mathbf{v}^\top) \mathbf{u}.$$

We have

$$\mathbf{v}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{u} = (\mathbf{u}^\top \mathbf{v}) \mathbf{v}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{v} + \sqrt{1 - (\mathbf{u}^\top \mathbf{v})^2} \mathbf{v}^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{v}_\perp.$$

Without loss of generality, we consider  $\mathbf{v} = \mathbf{e}_1$  and  $\mathbf{v}_\perp = \mathbf{e}_2$ . For the second term, we write  $\mathbf{z}_i := \mathbf{Z}_{\text{in}} \mathbf{e}_i$  and let  $\tilde{\mathbf{Z}} := \{\mathbf{z}_i\}_{i=3}^V$  and  $\mathbf{g} = \sqrt{d} \text{vec}(\tilde{\mathbf{Z}})$ .

$$\begin{aligned} \mathbf{e}_1^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{e}_2 &= (\|\mathbf{z}_1\|_2^2 + \|\mathbf{z}_2\|_2^2) \mathbf{z}_1^\top \mathbf{z}_2 + \mathbf{z}_1^\top \tilde{\mathbf{Z}} \tilde{\mathbf{Z}}^\top \mathbf{z}_2 \\ &= (\|\mathbf{z}_1\|_2^2 + \|\mathbf{z}_2\|_2^2) \mathbf{z}_1^\top \mathbf{z}_2 + \frac{1}{d} \mathbf{g}^\top \text{sym}(\mathbf{z}_1 \mathbf{z}_2^\top) \otimes \mathbf{I}_{V-2} \mathbf{g}. \end{aligned}$$

We have

- By Lemma 4, and Proposition 4

$$\begin{aligned} \mathbb{P} \left[ \left| \|\mathbf{z}_1\|_2^2 - 1 \right| \leq \frac{5 \log V}{\sqrt{d}} \text{ and } \left| \|\mathbf{z}_2\|_2^2 - 1 \right| \leq \frac{5 \log V}{\sqrt{d}} \text{ and } |\mathbf{z}_1^\top \mathbf{z}_2| \leq \frac{\log V}{\sqrt{d}} \right] \\ \leq 1 - 6 \exp(-c \log^2 V). \end{aligned}$$

- By Proposition 10, we have

$$- \|\text{sym}(\mathbf{z}_1 \mathbf{z}_2^\top) \otimes \mathbf{I}_{V-2}\|_2 \leq \|\mathbf{z}_1\|_2 \|\mathbf{z}_2\|_2$$

$$\begin{aligned} & - \|\text{sym}(\mathbf{z}_1 \mathbf{z}_2^\top) \otimes \mathbf{I}_{V-2}\|_F \leq \sqrt{V} \|\mathbf{z}_1\|_2 \|\mathbf{z}_2\|_2 \\ & - \text{tr}(\text{sym}(\mathbf{z}_1 \mathbf{z}_2^\top) \otimes \mathbf{I}_{V-2}) = (V-2) \mathbf{z}_1^\top \mathbf{z}_2. \end{aligned}$$

Therefore, by Lemma 4, we have

$$\begin{aligned} \mathbb{P} \left[ \left| \frac{1}{d} \mathbf{g}^\top \text{sym}(\mathbf{z}_1 \mathbf{z}_2^\top) \otimes \mathbf{I}_{V-2} \mathbf{g} - \frac{(V-2)}{d} \mathbf{z}_1^\top \mathbf{z}_2 \right| \leq 2 \|\mathbf{z}_1\|_2 \|\mathbf{z}_2\|_2 \left( \frac{\log V}{d} \sqrt{V} + \frac{\log^2 V}{d} \right) \right] \\ \leq 1 - 2 \exp(-c \log^2 V). \end{aligned}$$

By union bound of the precious two items, we have

$$\mathbb{P} \left[ \left| \mathbf{e}_1^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{e}_2 \right| \leq 2 \log V \left( \frac{V}{d^{3/2}} + \frac{\sqrt{V}}{d} \right) \right] \geq 1 - 8 \exp(-c \log^2 V). \quad (42)$$

Next, we redefine the notation:  $\tilde{\mathbf{Z}} := \{\mathbf{z}_i\}_{i=2}^V$ . We write

$$\mathbf{z}_1^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{z}_1 - 1 - \frac{V-1}{d} = \|\mathbf{z}_1\|_2^4 - 1 + \mathbf{z}_1^\top \left( \tilde{\mathbf{Z}} \tilde{\mathbf{Z}}^\top - \frac{V-1}{d} \mathbf{I}_d \right) \mathbf{z}_1 - \frac{V-1}{d} (\|\mathbf{z}_1\|_2^2 - 1)$$

By Proposition bla, we have

$$\mathbb{P} \left[ \mathbf{z}_1^\top \left( \tilde{\mathbf{Z}} \tilde{\mathbf{Z}}^\top - \frac{V-1}{d} \mathbf{I}_d \right) \mathbf{z}_1 \leq \log V \|\mathbf{z}_1\|_2^2 \frac{\sqrt{V}}{d} \right] \leq 1 - 2 \exp(-c \log^2 V)$$

By using the first item above, we have

$$\mathbb{P} \left[ \left| \mathbf{z}_1^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{z}_1 - 1 - \frac{V-1}{d} \right| \geq 6 \log V \left( \frac{\sqrt{V}}{d} + \frac{V}{d^{3/2}} \right) \right] \leq 1 - 2 \exp(-c \log^2 V). \quad (43)$$

The result follows (42) and (43).  $\square$

### E.3 MULTINOMIAL DISTRIBUTION AND RELATED STATEMENTS

**Lemma 5.** Let  $(n_1, \dots, n_V) \in \text{Mult}(N; (p_1, \dots, p_V))$ . For  $\mathbf{t} \in \mathbb{R}^V$ ,

$$\mathbb{E} \left[ \exp \left( \sum_{w=1}^V t_w n_w \right) \right] = \left( \sum_{w=1}^V p_w e^{t_w} \right)^N.$$

Then, if  $p_w = \frac{1}{V}$ ,  $w \in [V]$ , we have

$$\bullet \mathbb{E} \left[ \prod_{w=1}^V (n_w)_{j_w} \right] = N(N-1) \cdots (N-J+1) \prod_{w=1}^V p_w^{j_w}, \text{ where } J := \sum_{w=1}^V j_w.$$

• We have

$$\begin{aligned} & - \mathbb{E} [n_w^2] = \frac{N}{V} + \frac{N(N-1)}{V^2} \\ & - \mathbb{E} \left[ \left( \frac{n_w}{N} - \frac{1}{V} \right)^2 n_w \right] = \frac{(V-1)(N+V-2)}{NV^3}. \\ & - \mathbb{E} [n_w^3] = \frac{N}{V} + \frac{3N(N-1)}{V^2} + \frac{N(N-1)(N-2)}{V^3} \\ & - \mathbb{E} [n_w^4] = \frac{N}{V} + \frac{7N(N-1)}{V^2} + \frac{6N(N-1)(N-2)}{V^3} + \frac{N(N-1)(N-2)(N-3)}{V^4} \\ & - \mathbb{E} [n_w^2 n_{w'}^2] = \frac{N(N-1)}{V^2} + \frac{2N(N-1)(N-2)}{V^3} + \frac{N(N-1)(N-2)(N-3)}{V^4}. \\ & - \mathbb{E} \left[ \left( \sum_{w=1}^V n_w^2 \right)^2 \right] = N^2 + \frac{(N+4)N(N-1)}{V} + \frac{(N+2)N(N-1)(N-2)}{V^2} \end{aligned}$$

*Proof.* Let  $\mathbf{x}_i$  sampled from  $\{\mathbf{e}_1, \dots, \mathbf{e}_V\}$  with  $(p_1, \dots, p_V)$ . We have  $n_w = \sum_{i=1}^N \mathbf{e}_w^\top \mathbf{x}_i$ . We have

$$\mathbb{E} \left[ \exp \left( \sum_{w=1}^V t_w n_w \right) \right] = \mathbb{E} \left[ \exp \left( \sum_{i=1}^N \langle \mathbf{t}, \mathbf{x}_i \rangle \right) \right] = \left( \mathbb{E} \left[ \exp \left( \langle \mathbf{t}, \mathbf{x}_1 \rangle \right) \right] \right)^N = \left( \sum_{w=1}^V p_w e^{t_w} \right)^N.$$

The later statement can be derived by using  $z_w = e^{t_w}$  and taking derivatives of both sides with respect  $(z_1, \dots, z_V)$ .  $\square$

**Proposition 6.** Let  $\mathbf{n} = (n_1, \dots, n_V) \in \text{Mult}(L, \frac{1}{V} \mathbb{1}_V)$  and  $\mathbf{S} \in \mathbb{R}^{V \times V}$  be a symmetric matrix.. The following statements hold:

• We have

$$\begin{aligned} \mathbb{E}[\text{Diag}(\mathbf{n})\mathbf{S}\text{diag}(\mathbf{n})] &= L \mathbb{E}[\mathbf{x}_1^\top \mathbf{S} \mathbf{x}_1 \mathbf{x}_1 \mathbf{x}_1^\top] + \frac{L(L-1)}{V^2} \mathbf{S}, \\ \mathbb{E}[\text{Diag}(\mathbf{n} - \frac{L}{V} \mathbb{1}_V) \mathbf{S} \text{Diag}(\mathbf{n} - \frac{L}{V} \mathbb{1}_V)] &= L \mathbb{E}[\mathbf{x}_1^\top \mathbf{S} \mathbf{x}_1 \mathbf{x}_1 \mathbf{x}_1^\top] - \frac{L}{V^2} \mathbf{S}. \end{aligned}$$

• We have

$$\begin{aligned} \mathbb{E}[\mathbf{n} \mathbf{n}^\top \mathbf{S} \mathbf{n}] &= \frac{2L(L-1)}{V^2} \mathbf{S} \mathbb{1}_V + L \mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top \mathbf{S} \mathbf{x}_1] \\ &\quad + \left( \frac{L(L-1)}{V^2} \text{tr}(\mathbf{S}) + \frac{L(L-1)(L-2)}{V^3} \mathbb{1}_V^\top \mathbf{S} \mathbb{1}_V \right) \mathbb{1}_V \end{aligned}$$

• We have

$$\begin{aligned} \mathbb{E} \left[ \left( \left( \mathbf{n} - \frac{L}{V} \mathbb{1}_V \right)^\top \mathbf{S} \left( \mathbf{n} - \frac{L}{V} \mathbb{1}_V \right) \right)^2 \right] &= \frac{L}{V} \left\| \text{diag}(\mathbf{S}) - \frac{2}{V} \mathbf{S} \mathbb{1}_V + \frac{1}{V^2} (\mathbb{1}_V^\top \mathbf{S} \mathbb{1}_V) \mathbb{1}_V \right\|_2^2 \\ &\quad + \frac{L(L-1)}{V^2} \text{tr} \left( \left( \mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top \right) \mathbf{S} \right)^2 \\ &\quad + \frac{2L(L-1)}{V^2} \text{tr} \left( \left( \mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top \right) \mathbf{S} \left( \mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top \right) \mathbf{S} \right) \end{aligned}$$

*Proof.* For the first item, we observe that

$$\mathbf{e}_j^\top \mathbb{E}[\text{Diag}(\mathbf{n})\mathbf{S}\text{diag}(\mathbf{n})] \mathbf{e}_i = \mathbb{E}[n_j n_i] \mathbf{S}_{ij} = \left( \frac{L}{V} \delta_{ij} + \frac{L(L-1)}{V^2} \right) \mathbf{S}_{ij},$$

from which the first equation follows. For the second equation,

$$\begin{aligned} \mathbf{e}_j^\top \mathbb{E}[\text{diag}(\mathbf{n} - \frac{L}{V} \mathbb{1}_V) \mathbf{S} \text{diag}(\mathbf{n} - \frac{L}{V} \mathbb{1}_V)] \mathbf{e}_i &= \mathbb{E}[(n_j - \frac{L}{V})(n_i - \frac{L}{V})] \mathbf{S}_{ij} \\ &= \left( \frac{L}{V} \delta_{ij} - \frac{L}{V^2} \right) \mathbf{S}_{ij}. \end{aligned}$$

For the second item, we have

$$\begin{aligned} (\mathbb{E}[\mathbf{n} \mathbf{n}^\top \mathbf{S} \mathbf{n}])_i &= \sum_{jk} \mathbf{S}_{jk} \mathbb{E}[n_i n_j n_k] \\ &= \frac{L(L-1)(L-2)}{V^3} \left( \sum_{i \neq j \neq k} \mathbf{S}_{jk} \right) + \left( \frac{L(L-1)(L-2)}{V^3} + \frac{L(L-1)}{V^2} \right) \left( 2 \sum_{i \neq k} \mathbf{S}_{ik} + \sum_{i \neq k} \mathbf{S}_{kk} \right) \\ &\quad + \left( \frac{L}{V} + \frac{3L(L-1)}{V^2} + \frac{L(L-1)(L-2)}{V^3} \right) \mathbf{S}_{ii} \\ &= \frac{L}{V} \mathbf{S}_{ii} + \frac{L(L-1)}{V} \text{tr}(\mathbf{S}) + \frac{2L(L-1)}{V} \sum_k \mathbf{S}_{ik} + \frac{L(L-1)(L-2)}{V^3} \left( \sum_{jk} \mathbf{S}_{jk} \right). \end{aligned}$$

For the third item, we have  $(\mathbf{n} - \frac{L}{V} \mathbb{1}_V) =_d \sum_{\ell=1}^L (\xi_{1,\ell} - \frac{1}{V} \mathbb{1}_V)$  where the equality holds in distribution. For notational convenience, let

$$\gamma_{\ell u} := (\xi_{1,\ell} - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{S} (\xi_{1,u} - \frac{1}{V} \mathbb{1}_V).$$

Then,

$$\mathbb{E} \left[ \left( \left( \mathbf{n} - \frac{L}{V} \mathbb{1}_V \right)^\top \mathbf{S} \left( \mathbf{n} - \frac{L}{V} \mathbb{1}_V \right) \right)^2 \right] =_d \sum_{\ell, u} \sum_{\ell', u'} \mathbb{E}[\gamma_{\ell u} \gamma_{\ell' u'}]$$

By independence, only  $(\ell, u, \ell', u')$  where each index occur even times contribute. The possible cases are as follows:

- All four indices equal ( $\ell = u = \ell' = u'$ ): There are  $L$  many terms here with contribution

$$\begin{aligned}\mathbb{E}[\gamma_{\ell\ell}\gamma_{\ell\ell}] &= \frac{1}{V} \left\| \text{diag} \left( \left( \mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top \right) \mathbf{S} \left( \mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top \right) \right) \right\|_2^2 \\ &= \frac{1}{V} \left\| \text{diag}(\mathbf{S}) - \frac{2}{V} \mathbf{S} \mathbb{1}_V + \frac{1}{V^2} (\mathbb{1}_V^\top \mathbf{S} \mathbb{1}_V) \mathbb{1}_V \right\|_2^2.\end{aligned}$$

- Two distinct indices, both pairs diagonal ( $\ell = u$  and  $\ell' = u'$  and  $\ell \neq \ell'$ ): There are  $L(L-1)$  many terms here with contribution

$$\mathbb{E}[\gamma_{\ell\ell}\gamma_{\ell'\ell'}] = \mathbb{E}[\gamma_{\ell\ell}] \mathbb{E}[\gamma_{\ell'\ell'}] = \text{tr} \left( \left( \mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top \right) \mathbf{S} \right)^2$$

- Two distinct indices, paired off-diagonal: ( $\ell = \ell'$  and  $u = u'$  and  $\ell \neq u$ ): There are  $2L(L-1)$  many terms here with contribution

$$\begin{aligned}\mathbb{E}[\gamma_{\ell u}\gamma_{\ell u}] &= \text{tr} \left( \mathbb{E}[(\xi_{1,1} - \frac{1}{V} \mathbb{1}_V)(\xi_{1,1} - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{S} (\xi_{1,2} - \frac{1}{V} \mathbb{1}_V)(\xi_{1,2} - \frac{1}{V} \mathbb{1}_V)^\top \mathbf{S}] \right) \\ &= \text{tr} \left( \left( \mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top \right) \mathbf{S} \left( \mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top \right) \mathbf{S} \right)\end{aligned}$$

□

**Proposition 7.** Let  $V^3 \gg L$ . There exists a universal  $C > 0$  such that the following holds:

- Let  $m_{ij} := (1 - \frac{1}{V}) \mathbb{1}_{i=j} + \frac{L}{V}$ . For  $K > 0$  and  $p \geq \log V$ ,

$$\mathbb{E} \left[ \left| \frac{1}{L} \mathbb{1}_L^\top \mathbf{X}_i \mathbf{X}_j^\top \mathbb{1}_L - m_{ij} \right|^p \right]^{\frac{1}{p}} \leq C \left( \frac{p^{\frac{3}{2}}}{\sqrt{V}} + \frac{p^2}{L} \right)$$

$$\mathbb{P} \left[ \left| \frac{1}{L} \mathbb{1}_L^\top \mathbf{X}_i \mathbf{X}_j^\top \mathbb{1}_L - m_{ij} \right| \geq CK^2 \frac{\log^2 V}{\sqrt{V} \wedge L} \right] \leq \frac{1}{V^K}.$$

- For  $K > 0$  and  $p \geq \log V$ ,

$$\begin{aligned}\mathbb{E} \left[ \left\| \frac{1}{NL} \sum_{i=1}^N \left( \mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right) \mathbb{1}_L \mathbb{1}_L^\top \left( \mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right)^\top - \frac{1}{V} \left( \mathbf{I} - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top \right) \right\|_2^p \right]^{\frac{1}{p}} \\ \leq C \left( \sqrt{\frac{p}{NV}} + \frac{p}{N} \left( 1 + \frac{p^2}{\sqrt{V} \wedge L} \right) \right) \\ \mathbb{P} \left[ \left\| \frac{1}{NL} \sum_{i=1}^N \left( \mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right) \mathbb{1}_L \mathbb{1}_L^\top \left( \mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right)^\top - \frac{1}{V} \left( \mathbf{I} - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top \right) \right\|_2 \right. \\ \left. > CK \log^2 V \left( \frac{1}{\sqrt{NV}} + \frac{1}{N} \left( 1 + \frac{\log^2 V}{\sqrt{V} \wedge L} \right) \right) \right] \leq \frac{1}{V^K}.\end{aligned}$$

*Proof.* For  $i = j$  in the first item, we have

$$\frac{1}{L} \mathbb{1}_L^\top \mathbf{X}_i \mathbf{X}_i^\top \mathbb{1}_L = 1 + \frac{2}{L} \sum_{1 \leq j < k \leq L} \mathbb{1}_{\xi_{i,j} = \xi_{i,k}} = 1 + \frac{(L-1)}{V} + \frac{2}{L} \sum_{k=2}^L \sum_{j=1}^{k-1} \mathbb{1}_{\xi_{i,j} = \xi_{i,k}} - \frac{1}{V}$$

Define

$$Y_k := \sum_{j=1}^{k-1} \left( \mathbb{1}_{\xi_{i,j} = \xi_{i,k}} - \frac{1}{V} \right) \text{ and } \mathcal{F}_k := \sigma(Y_1, \dots, Y_k).$$

Given that

$$\sum_{j=1}^{k-1} \mathbb{1}_{\xi_{i,j} = \xi_{i,k}} | \xi_{i,k} \sim \text{Binomial}(k-1, \frac{1}{V}) \Rightarrow \mathbb{E}[|Y_k|^p]^{\frac{1}{p}} \leq C(\sqrt{p} \sqrt{\frac{L}{V}} + p), \quad p \geq \log V.$$

where we used Corollary 3. As for the quadratic variation

$$\begin{aligned} Q_L &:= \sum_{k=1}^L \mathbb{E}[Y_k^2 | \mathcal{F}_{k-1}] &= \sum_{k=1}^L \frac{1}{V} \left( \|\mathbf{X}_i^\top \mathbb{1}_{k-1}\|_2^2 - \frac{(k-1)^2}{V} \right) \\ &= \frac{1}{V} \sum_{k=1}^L \left\| \left( \mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{k-1}^\top \right) \mathbb{1}_{k-1} \right\|_2^2 \end{aligned}$$

For  $p \geq \log V$ , by using triangle inequality,

$$\begin{aligned} \mathbb{E}[|Q_L|^{\frac{p}{2}}]^{\frac{2}{p}} &\leq \frac{1}{V} \sum_{k=1}^L \mathbb{E} \left[ \left\| \left( \mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{k-1}^\top \right) \mathbb{1}_{k-1} \right\|_2^p \right]^{\frac{2}{p}} \\ &\leq \frac{1}{V} \sum_{k=1}^V (k-1) \mathbb{E} \left[ \|\mathbf{X}_i^\top \mathbb{1}_{k-1}\|_p^p \right]^{\frac{2}{p}} + \sum_{k=V+1}^L \mathbb{E} \left[ \left\| \left( \mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{k-1}^\top \right) \mathbb{1}_{k-1} \right\|_p^p \right]^{\frac{2}{p}} \\ &\leq Cp^2 \frac{1}{V} \sum_{k=1}^L k = Cp^2 \frac{L^2}{V} \end{aligned}$$

where we used Corollary 3. By Proposition 13, for  $p \geq \log V$ , we have

$$\mathbb{E} \left[ \left| \sum_{k=1}^L Y_k \right|^p \right]^{\frac{1}{p}} \leq C \left( p\sqrt{p} \frac{L}{\sqrt{V}} + p^2 \right).$$

By using  $p = \log V$ , we have

$$\mathbb{P} \left[ \left| \frac{1}{L} \sum_{k=1}^L Y_k \right| > CeK^2 \frac{\log^2 V}{\sqrt{V} \wedge L} \right] \leq \frac{1}{V^K}.$$

Hence, we have the  $i = j$  case. For the  $i \neq j$  case, we have

$$\frac{1}{L} \mathbb{1}_L^\top \mathbf{X}_j \mathbf{X}_i^\top \mathbb{1}_L = \frac{L}{V} + \frac{1}{L} \sum_{\ell=1}^L \sum_{k=1}^L \mathbb{1}_{\xi_{i,k}=\xi_{j,\ell}} - \frac{1}{V}$$

We redefine the martingale difference sequence as

$$Y_k := \sum_{\ell=1}^L \mathbb{1}_{\xi_{i,k}=\xi_{j,\ell}} - \frac{1}{V}.$$

Conditioned on  $\mathbf{X}_j$ , we have  $\{Y_1, \dots, Y_L\}$  are i.i.d. and

$$\mathbb{E}[Y_k | \mathbf{X}_j] = 0 \quad \text{and} \quad \mathbb{E}[Y_k^p | \mathbf{X}_j] = \frac{1}{V} \left\| \left( \mathbf{X}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right) \mathbb{1}_L \right\|_p^p$$

By Proposition 13, for  $p \geq \log V$ , we have

$$\mathbb{E} \left[ \left| \frac{1}{L} \sum_{k=1}^L Y_k \right|^p \right]^{\frac{1}{p}} \leq C \left( \frac{\sqrt{p}}{\sqrt{V}} + \frac{p^{\frac{3}{2}}}{\sqrt{LV}} + \frac{p^2}{L} \right).$$

By using  $p = \log V$ , we have

$$\mathbb{P} \left[ \left| \frac{1}{L} \mathbb{1}_L^\top \mathbf{X}_j \mathbf{X}_i^\top \mathbb{1}_L - \frac{L}{V} \right| \geq \frac{CK^2 \log^2 V}{\sqrt{V} \vee L} \right] \leq \frac{1}{V^K}.$$

For the second item, we define

$$\mathbf{Y}_k := \frac{1}{L} \left( \mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right) \mathbb{1}_L \mathbb{1}_L^\top \left( \mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right)^\top - \frac{1}{V} (\mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top)$$

and  $\mathbf{Q}_N := N \mathbb{E}[\mathbf{Y}_1^2]$ . We have

$$\mathbf{Q}_N \preceq N \mathbb{E} \left[ \left\| \frac{1}{\sqrt{L}} \left( \mathbf{X}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right) \mathbb{1}_L \right\|_2^2 \frac{1}{2L} \left( \mathbf{X}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right) \mathbb{1}_L \mathbb{1}_L^\top \left( \mathbf{X}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right)^\top \right]$$

$$\begin{aligned}
&= N \mathbb{E} \left[ \left(1 - \frac{1}{V}\right) \frac{1}{L} (\mathbf{X}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L \mathbb{1}_L^\top (\mathbf{X}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top)^\top \right] \\
&+ N \mathbb{E} \left[ \left( \left\| \frac{1}{\sqrt{L}} (\mathbf{X}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L \right\|_2^2 - \left(1 - \frac{1}{V}\right) \right) \right. \\
&\quad \times \left. \left( \frac{1}{L} (\mathbf{X}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L \mathbb{1}_L^\top (\mathbf{X}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top)^\top - \frac{1}{V} (\mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top) \right) \right] \\
&\preceq \frac{CN}{V} \mathbf{I}_V + \frac{1}{2} \mathbf{Q}_N
\end{aligned}$$

Therefore, we have  $\|\mathbf{Q}_N\|_2 \leq \frac{CN}{V}$ . Moreover, by using the first item,

$$\mathbb{E}[\|\mathbf{Y}_k\|_2^p]^{\frac{1}{p}} \leq \mathbb{E} \left[ \left\| \frac{1}{\sqrt{L}} (\mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L \right\|_2^{2p} \right]^{\frac{1}{p}} \leq 1 + C \left( \frac{p^{\frac{3}{2}}}{\sqrt{V}} + \frac{p^2}{L} \right)$$

Therefore, by using Proposition 13, we have

$$\begin{aligned}
&\mathbb{E} \left[ \left\| \frac{1}{N} \sum_{i=1}^N \mathbf{Y}_i - \frac{1}{V} (\mathbf{I} - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top) \right\|_2^p \right] \\
&\leq C \left( \sqrt{p \vee \log V} \sqrt{\frac{1}{NV}} + (p \vee \log V) N^{\frac{1}{p}-1} \left( 1 + \frac{p^{\frac{3}{2}}}{\sqrt{V}} + \frac{p^2}{L} \right) \right)
\end{aligned}$$

By using  $p = \log V$ , we have

$$\mathbb{P} \left[ \left\| \frac{1}{N} \sum_{i=1}^N \mathbf{Y}_i - \frac{1}{V} (\mathbf{I} - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top) \right\|_2 > CK \log^2 V \left( \frac{1}{\sqrt{NV}} + \frac{1}{N} \left( 1 + \frac{\log^2 V}{\sqrt{V} \wedge L} \right) \right) \right] \leq \frac{1}{V^K}.$$

□

**Proposition 8.** We consider  $\mathbf{S}_1$ ,  $\mathbf{S}_2$  and  $\mathbf{S}_3$  defined in (16), (17) and (18) in the regime  $V^3 \gg N \gg V$  and  $L \asymp V^\varepsilon$ ,  $\varepsilon \in (0, 1)$ . For any  $K > 0$  and  $V \geq \Omega_{K,\varepsilon}(1)$ , the following holds:

1. We have

$$\mathbb{P} \left[ \left| \text{tr}(\mathbf{S}_1) - \frac{1}{L^2} \left( \frac{1}{N} + \left(1 - \frac{1}{V}\right) \frac{1}{V} \right) \right| > CK^2 \frac{\log^2 V}{L^2 N \sqrt{V}} \text{ or } \|\mathbf{S}_1\|_2 > \frac{e^2}{L^2 V^2} \right] \leq \frac{2}{V^K}.$$

2. We have

$$\mathbb{P} \left[ \left| \text{tr}(\mathbf{S}_2) - \left(1 - \frac{1}{V}\right)^2 \frac{L-1}{L^2 N} \right| > C \frac{K^{\frac{3}{2}} \log^3 V}{N \sqrt{LV}} \text{ or } \|\mathbf{S}_2\|_2 > C \frac{K^{\frac{3}{2}} \log^2 V}{NLV} \right] \leq \frac{4}{V^K}.$$

3. We have

$$\mathbb{P} \left[ \frac{-CK^2 \log^2 V}{N \sqrt{V}} \frac{1}{V^2 L^2} \mathbb{1}_V \mathbb{1}_V^\top \preceq \mathbf{S}_3 - \frac{1}{N} \frac{1}{V^2 L^2} \mathbb{1}_V \mathbb{1}_V^\top \preceq \frac{CK^2 \log^2 V}{N \sqrt{V}} \frac{1}{V^2 L^2} \mathbb{1}_V \mathbb{1}_V^\top \right] \leq \frac{1}{V^K}.$$

*Proof.* We define  $n_i := |\{j \leq N \mid \mathbf{x}_j = \mathbf{e}_i\}|$ . We have

$$\text{tr}(\mathbf{S}_1) = \left(1 - \frac{1}{V}\right) \frac{1}{L^2 N^2} \sum_{i=1}^V n_i^2 \text{ and } \|\mathbf{S}_1\|_2 \leq \sup_{i \leq N} \frac{n_i^2}{L^2 N^2}$$

By using Proposition 7 and Corollary 3, we have the first item. For the second item, we write

$$\begin{aligned}
\mathbf{S}_2 &= \frac{(1 - \frac{1}{V})}{L^2 N^2} \sum_{j=1}^N (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \\
&+ \frac{2(1 - \frac{1}{V})}{L^2 N^2} \sum_{j < k} (\mathbb{1}_{\mathbf{x}_j = \mathbf{x}_k} - \frac{1}{V}) \text{sym} \left( (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \mathbb{1}_{L-1}^\top (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right)
\end{aligned}$$

$$=: \mathbf{S}_{21} + \mathbf{S}_{22}$$

We will analyze  $\mathbf{S}_{21}$  and  $\mathbf{S}_{22}$  separately. We start with  $\mathbf{S}_{21}$ . We have

$$\begin{aligned} \text{tr}(\mathbf{S}_{21}) - (1 - \frac{1}{V})^2 \frac{L-1}{L^2 N} \\ = (1 - \frac{1}{V}) \frac{L-1}{L^2 N^2} \sum_{j=1}^N \underbrace{\left\| \frac{1}{\sqrt{L-1}} (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\|_2^2}_{:= Y_i} - (1 - \frac{1}{V}). \end{aligned}$$

We have  $\mathbb{E}[Y_i^2] \leq \frac{2}{V}$  and by the first item in Proposition 7

$$\mathbb{E}[|Y_i|^p]^{\frac{1}{p}} \leq \frac{Cp^2}{\sqrt{V} \vee L}.$$

Therefore, by Proposition 13,

$$\mathbb{E} \left[ \left| \text{tr}(\mathbf{S}_{21}) - (1 - \frac{1}{V})^2 \frac{L-1}{L^2 N} \right|^p \right]^{\frac{1}{p}} \leq \frac{C}{LN^2} \left( \sqrt{\frac{pN}{V}} + pN^{\frac{1}{p}} \frac{p^2}{\sqrt{V} \vee L} \right)$$

By using  $p = \log V$ , we have

$$\mathbb{P} \left[ \left| \text{tr}(\mathbf{S}_{21}) - (1 - \frac{1}{V})^2 \frac{L-1}{L^2 N} \right| > C \frac{K \log^3 V}{LN \sqrt{NV}} \right] \leq \frac{1}{V^K}. \quad (44)$$

Moreover, by Proposition 7, we have

$$\mathbb{P} \left[ \left\| \mathbf{S}_{21} - (1 - \frac{1}{V}) \frac{L-1}{L^2 N} \frac{1}{V} (\mathbf{I}_V - \mathbb{1}_V \mathbb{1}_V^\top) \right\|_2 > C \frac{K \log^2 V}{LN} \left( \frac{1}{\sqrt{NV}} + \frac{1}{N} (1 + \frac{\log^2 V}{\sqrt{V} \wedge L}) \right) \right] \leq \frac{1}{V^K}. \quad (45)$$

As for  $\mathbf{S}_{22}$ , we have

$$\begin{aligned} \text{tr}(\mathbf{S}_{22}) \\ = \frac{2(1 - \frac{1}{V})}{L^2 N^2} \sum_{k=2}^N \sum_{j=1}^{k-1} (\mathbb{1}_{\mathbf{x}_j = \mathbf{x}_k} - \frac{1}{V}) \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \\ = \frac{2(1 - \frac{1}{V})}{L^2 N^2} \sum_{k=2}^N \sum_{j=1}^{k-1} (\mathbb{1}_{\mathbf{x}_j = \mathbf{x}_k} - \frac{1}{V}) \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \mathbf{N}_k^\top \mathbb{1}_{L-1} \end{aligned}$$

We define  $\mathcal{F}_k := \sigma(\mathbf{N}_{1:k})$  and

$$Y_k := (1 - \frac{1}{V}) \frac{2}{L^2 N^2} \sum_{j=1}^{k-1} (\mathbb{1}_{\mathbf{x}_j = \mathbf{x}_k} - \frac{1}{V}) \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \mathbf{N}_k^\top \mathbb{1}_{L-1}$$

We have

$$\mathbb{E}[Y_k^2 | \mathcal{F}_{k-1}] = (1 - \frac{1}{V})^3 \frac{4(L-1)}{L^4 N^4} \frac{1}{V} \sum_{j=1}^{k-1} \left\| (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\|_2^2$$

Then,

$$\begin{aligned} Q_N = \sum_{k=1}^N \mathbb{E}[Y_k^2 | \mathcal{F}_{k-1}] &= (1 - \frac{1}{V})^3 \frac{4(L-1)}{L^4 N^4} \frac{1}{V} \sum_{k=1}^N \sum_{j=1}^{k-1} \left\| (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\|_2^2 \\ &= (1 - \frac{1}{V})^3 \frac{4(L-1)}{L^4 N^4} \frac{1}{V} \sum_{k=1}^N (N-k) \left\| (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\|_2^2 \end{aligned}$$

Then, for  $p \geq \log V$ ,

$$\mathbb{E} \left[ |Q_N|^{\frac{p}{2}} \right]^{\frac{2}{p}} \leq \frac{5}{L^3 N^3 V} \sum_{k=1}^N \mathbb{E} \left[ \left\| \mathbf{N}_k^\top \mathbb{1}_{L-1} \right\|_2^p \right]^{\frac{2}{p}} \leq \frac{5}{L N^3 V} \sum_{k=1}^N \mathbb{E} \left[ \left\| \mathbf{N}_k^\top \mathbb{1}_{L-1} \right\|_p^p \right]^{\frac{2}{p}} \leq \frac{5p^2}{L N^2 V} \quad (46)$$

By using Proposition 13, we show the following:



- To bound  $\mathbb{E}[|Y_k|^p]^{\frac{1}{p}}$  for  $p \geq \log V$ , we first write

$$\begin{aligned}
& \mathbb{E}[|Y_k|^p | \mathbf{N}_{1:k}, \mathbf{x}_k]^{\frac{1}{p}} \\
& \leq \frac{C}{LN^2} \frac{\sqrt{p}}{\sqrt{V}} \left( \sum_{j=1}^{k-1} \left| \frac{1}{L-1} \left\langle (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1}, (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\rangle \right|^2 \right)^{\frac{1}{2}} \\
& \quad + \frac{Cp}{LN^2} \left( \sum_{j=1}^{k-1} \left| \frac{1}{L-1} \left\langle (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1}, (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\rangle \right|^p \right)^{\frac{1}{p}} \\
& \text{Therefore,} \\
& \mathbb{E}[|Y_k|^p]^{\frac{1}{p}} \leq \frac{C}{LN^2} \left( \frac{\sqrt{p}\sqrt{k}}{\sqrt{V}} + pk^{\frac{1}{p}} \right) \\
& \quad \times \mathbb{E} \left[ \left| \frac{1}{L-1} \left\langle (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1}, (\mathbf{N}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\rangle \right|^p \right]^{\frac{1}{p}} \\
& \leq \frac{C}{LN^2} \left( \frac{\sqrt{p}\sqrt{k}}{\sqrt{V}} + pk^{\frac{1}{p}} \right) \mathbb{E} \left[ \left| \frac{1}{L-1} \left\langle \mathbf{N}_k^\top \mathbb{1}_{L-1}, \mathbf{N}_1^\top \mathbb{1}_{L-1} \right\rangle - \frac{L-1}{V} \right|^p \right]^{\frac{1}{p}} \\
& \leq \frac{Cp}{LN^2} \left( \frac{\sqrt{p}\sqrt{k}}{\sqrt{V}} + pk^{\frac{1}{p}} \right) \left( \frac{1}{\sqrt{V}} \vee \frac{1}{L} \right). \tag{47}
\end{aligned}$$

- Then by using (46) and (47), we have for  $p \geq \log V$

$$\mathbb{E}[|\text{tr}(\mathbf{S}_{22})|^p]^{\frac{1}{p}} \leq C \left( \frac{p^{\frac{3}{2}}}{N\sqrt{LV}} + \frac{p^3 N^{\frac{2}{p}}}{LN^{\frac{3}{2}}\sqrt{V}} \left( \frac{1}{\sqrt{V}} \vee \frac{1}{L} \right) \right) \leq \frac{Cp^{\frac{3}{2}}}{N\sqrt{LV}}.$$

Therefore, by using  $p = \log V$ ,

$$\mathbb{P} \left[ |\text{tr}(\mathbf{S}_{22})| > \frac{CK^{\frac{3}{2}} \log^{\frac{3}{2}} V}{N\sqrt{LV}} \right] \leq \frac{1}{V^K} \tag{48}$$

To bound  $\|\mathbf{S}_{22}\|_2$ , we define

$$\mathbf{Y}_k := \sum_{j=1}^{k-1} \left( \mathbb{1}_{\mathbf{x}_j = \mathbf{x}_k} - \frac{1}{V} \right) \text{sym} \left( (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \mathbb{1}_{L-1}^\top (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right).$$

We have

$$\begin{aligned}
& \mathbb{E}[\mathbf{Y}_k^2 | \mathcal{F}_{k-1}] \\
& \preceq \frac{2}{V} \sum_{j=1}^{k-1} \mathbb{E} \left[ \left\| (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \mathbb{1}_{L-1}^\top (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right\| \left\| (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \mathbb{1}_{L-1}^\top \right\|_2^2 \middle| \mathcal{F}_{k-1} \right] \\
& \quad + \frac{2}{V} \sum_{j=1}^{k-1} \mathbb{E} \left[ \left\| (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\|_2^2 \left\| (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right\| \middle| \mathcal{F}_{k-1} \right] \\
& \preceq \frac{2L}{V^2} \sum_{j=1}^{k-1} \left\| (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\|_2^2 \mathbf{I}_V + \frac{2L}{V} \sum_{j=1}^{k-1} (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top
\end{aligned}$$

Therefore, we have

$$\mathbf{Q}_N := \sum_{k=1}^N \mathbb{E}[\mathbf{Y}_k^2 | \mathcal{F}_{k-1}]$$

$$\begin{aligned}
&\preceq \frac{2NL}{V^2} \sum_{j=1}^{N-1} \left\| (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\|_2^2 \mathbf{I}_V \\
&+ \frac{2L^2 N^2}{V} \frac{1}{LN} \sum_{j=1}^{N-1} (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top.
\end{aligned}$$

Then,

$$\begin{aligned}
\mathbb{E}[\|\mathbf{Q}_N\|_2^{\frac{p}{2}}] &\leq \frac{2NL}{V^2} \mathbb{E} \left[ \left( \sum_{j=1}^{N-1} \left\| (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\|_2^2 \right)^{\frac{p}{2}} \right]^{\frac{2}{p}} \\
&+ \frac{2L^2 N^2}{V} \mathbb{E} \left[ \left\| \frac{1}{N(L-1)} \sum_{j=1}^{N-1} (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right\|_2^{\frac{p}{2}} \right]^{\frac{2}{p}} \\
&\leq \frac{2N^2 L^2}{V^2} \mathbb{E} \left[ \left\| \frac{1}{\sqrt{L}} (\mathbf{N}_1^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\|_2^p \right]^{\frac{2}{p}} \\
&+ \frac{2L^2 N^2}{V} \mathbb{E} \left[ \left\| \frac{1}{N(L-1)} \sum_{j=1}^{N-1} (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \mathbb{1}_{L-1}^\top (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right\|_2^{\frac{p}{2}} \right]^{\frac{2}{p}} \\
&\leq \frac{CN^2 L^2}{V^2} \left( 1 + \frac{p^2}{\sqrt{V} \vee L} \right) + \frac{CL^2 N^2}{V} \left( \frac{1}{V} + \sqrt{\frac{p}{NV}} + \frac{p}{N} \left( 1 + \frac{p^2}{\sqrt{V} \wedge L} \right) \right) \\
&\leq \frac{CN^2 L^2}{V^2} \left( 1 + \frac{p^2}{\sqrt{V} \vee L} \right).
\end{aligned}$$

To bound  $\mathbb{E}[\|\mathbf{Y}_k\|_2^p]$ , we observe that

• We have

$$\begin{aligned}
&\mathbb{E} \left[ \left( (\mathbb{1}_{\mathbf{x}_j = \mathbf{x}_k} - \frac{1}{V}) \text{sym} \left( (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \mathbb{1}_{L-1}^\top (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right) \right)^2 \middle| \mathbf{x}_k, \mathbf{N}_k \right] \\
&\preceq \frac{L}{V^2} \|(\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1}\|_2^2 \mathbf{I}_V \\
&+ \frac{L}{V} (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \mathbb{1}_{L-1}^\top (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top
\end{aligned}$$

Moreover,

$$\begin{aligned}
&\mathbb{E} \left[ \left\| (\mathbb{1}_{\mathbf{x}_j = \mathbf{x}_k} - \frac{1}{V}) \text{sym} \left( (\mathbf{N}_j^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \mathbb{1}_{L-1}^\top (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top)^\top \right) \right\|^p \middle| \mathbf{x}_k, \mathbf{N}_k \right]^{\frac{1}{p}} \\
&= \mathbb{E} \left[ \left| (\mathbb{1}_{\mathbf{x}_j = \mathbf{x}_k} - \frac{1}{V}) \right|^p \middle| \mathbf{x}_k \right]^{\frac{1}{p}} \mathbb{E} \left[ \left\| \left\langle \mathbf{N}_j^\top \mathbb{1}_{L-1}, (\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1} \right\rangle \right\|^p \middle| \mathbf{N}_k \right]^{\frac{1}{p}} \\
&\leq C \left( \sqrt{\frac{p}{V}} + p \right) \left( \sqrt{\frac{p}{V}} \|(\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1}\|_2 + p \|(\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1}\|_p \right)
\end{aligned}$$

• By Proposition 13, we have

$$\begin{aligned}
\mathbb{E}[\|\mathbf{Y}_k\|_2^p \middle| \mathbf{x}_k, \mathbf{N}_k]^{\frac{1}{p}} &\leq C \sqrt{p \vee \log V} \sqrt{\frac{Lk}{V}} \|(\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1}\|_2 \\
&+ C(p \vee \log V) k^{\frac{1}{p}} \frac{p^{3/2}}{V} \|(\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1}\|_2 \\
&+ C(p \vee \log V) k^{\frac{1}{p}} p^2 \|(\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1}\|_p
\end{aligned}$$

Therefore, for  $p \geq \log V$

$$\mathbb{E}[\|\mathbf{Y}_k\|_2^p]^{\frac{1}{p}} \leq C \sqrt{\frac{pLk}{V}} \mathbb{E}[\|(\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1}\|_2^p]^{\frac{1}{p}}$$

$$\begin{aligned}
& + Ck^{\frac{1}{p}} \frac{p^{5/2}}{\sqrt{V}} \mathbb{E}[\|(\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1}\|_2^p]^{\frac{1}{p}} \\
& + Ck^{\frac{1}{p}} p^3 \mathbb{E}[\|(\mathbf{N}_k^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_{L-1}^\top) \mathbb{1}_{L-1}\|_p^p]^{\frac{1}{p}} \\
& \leq C \left( \sqrt{\frac{pLk}{V}} + \frac{p^{5/2}}{\sqrt{V}} \right) (\sqrt{L} + p\sqrt{\frac{pL}{V}} + \frac{p^2}{\sqrt{L}}) + Cp^3 \left( \sqrt{\frac{pL}{V}} + p \right)
\end{aligned}$$

Therefore, for  $p = \log V$ , we have

$$\mathbb{E}[\|\mathbf{S}_{22}\|_2^p] \leq C \left( \frac{\sqrt{p}}{NLV} + \frac{p^{3/2}}{LN\sqrt{NV}} + \frac{p^5}{L^2N^2} \right)$$

Therefore, we have

$$\mathbb{P} \left[ \|\mathbf{S}_{22}\|_2 > C \frac{K^{3/2} \log^{3/2} V}{NLV} \right] \leq \frac{1}{V^K}. \quad (49)$$

By (44), (45), (48), and (49), we have the second item. For the last item, we have

$$\mathbf{S}_3 - \frac{1}{N} \frac{1}{V^2 L^2} \mathbb{1}_V \mathbb{1}_V^\top = \frac{1}{V^2 L^2} \mathbb{1}_V \mathbb{1}_V^\top \left( \left\| \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2^2 - \frac{1}{N} \right)$$

By Proposition 7,

$$\mathbb{P} \left[ \left| \left\| \frac{1}{N} \sum_{j=1}^N (\mathbf{x}_j - \frac{1}{V} \mathbb{1}_V) \right\|_2^2 - \frac{1}{N} \right| > \frac{CK^2 \log^2 V}{N\sqrt{V}} \right] \leq \frac{1}{V^K}.$$

The displayed equation implies the third item.  $\square$

**Proposition 9.** Let  $\mathbf{z}_{\nu,\delta} = \mathbf{z}_\nu + \mathbb{1}_{k=1} \delta \mathbf{z}_{\text{trig}}$ . Given that (E.1) holds, the following statements hold:

1. We have for  $i \neq r$ ,

$$\left| \mathbb{E} \left[ \left( \mathbb{1}_{\mathbf{x}_i = \mathbf{x}_r} - \frac{1}{V} \right) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \mid \mathbf{Z}_{\text{in}} \right] \right| \leq \frac{C}{Vd}.$$

2. We have

$$\mathbb{E}[\|\mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta}\|_2^2 \mid \mathbf{Z}_{\text{in}}] \leq C \left( \frac{L}{d} + \frac{L^2}{d^2} \right).$$

3. We have for  $i \neq r$ ,

$$\left| \mathbb{E} \left[ \left( \mathbb{1}_{\mathbf{x}_i = \mathbf{x}_r} - \frac{1}{V} \right) \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \mid \mathbf{Z}_{\text{in}} \right] \right| \leq \frac{C \log^2 V}{d^2}.$$

4. We have

$$\mathbb{E} \left[ \left( \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu,\delta} \right)^2 \mid \mathbf{Z}_{\text{in}} \right] \leq C \frac{V \log^2 V}{d} \left( \frac{L}{d} + \frac{L^2}{d^2} \right).$$

5. For notational convenience, let

$$\gamma := \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_i^\top \mathbf{X}_i - \frac{L}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right) \mathbb{1}_L.$$

We have

$$|\mathbb{E}[\gamma \mid \mathbf{Z}_{\text{in}}]| \leq \frac{CL \log V}{\sqrt{Vd}} \quad \text{and} \quad \mathbb{E}[\gamma^2 \mid \mathbf{Z}_{\text{in}}] \leq C \log^2 V \left( \frac{L}{d} + \frac{L^2}{d^2} \right).$$

*Proof.* For the first item, we have

$$\begin{aligned}
& \mathbb{E} \left[ \left( \mathbb{1}_{\mathbf{x}_i = \mathbf{x}_r} - \frac{1}{V} \right) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \mid \mathbf{Z}_{\text{in}} \right] \\
&= \mathbb{E} \left[ \left( \mathbb{1}_{\mathbf{x}_i = \mathbf{x}_u} - \frac{1}{V} \right) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_r \mathbf{x}_r^\top \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \mid \mathbf{Z}_{\text{in}} \right] \\
&= \frac{1}{V} \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \mid \mathbf{Z}_{\text{in}} \right] - \frac{1}{V^3} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \\
&\leq \frac{C}{Vd}
\end{aligned}$$

For the second item, we write

$$\begin{aligned}
& \mathbb{E} [\mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \mid \mathbf{Z}_{\text{in}}] \\
&= \frac{L}{V} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \text{diag}(\mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}}) \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} + \frac{L(L-1)}{V^2} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \leq C \left( \frac{L}{d} + \frac{L^2}{d^2} \right)
\end{aligned}$$

For the third item, we have

$$\begin{aligned}
& \mathbb{E} \left[ \left( \mathbb{1}_{\mathbf{x}_i = \mathbf{x}_r} - \frac{1}{V} \right) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_r^\top \mathbf{X}_r \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \mid \mathbf{Z}_{\text{in}} \right] \\
&= \mathbb{E} \left[ \left( \mathbb{1}_{\mathbf{x}_i = \mathbf{x}_r} - \frac{1}{V} \right) \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_r \mathbf{x}_r^\top \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \mid \mathbf{Z}_{\text{in}} \right] \\
&= \frac{1}{V} \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \mid \mathbf{Z}_{\text{in}} \right] \\
&\quad - \frac{1}{V^3} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \\
&\leq C \frac{\log^2 V}{d^2}
\end{aligned}$$

For the fourth item, we have

$$\begin{aligned}
& \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \mid \mathbf{Z}_{\text{in}} \right] \\
&= L \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \mid \mathbf{Z}_{\text{in}} \right] \\
&\quad + \frac{L(L-1)}{V^2} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \mathbb{1}_V^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{z}_{\nu, \delta} \\
&= C \left( \frac{LV \log^2 V}{d^2} + \frac{L^2 V \log^2 V}{d^3} \right)
\end{aligned}$$

For the fifth item, we have

$$\begin{aligned}
& \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_i^\top \mathbf{X}_i - \frac{L}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mid \mathbf{Z}_{\text{in}} \right] \\
&= \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbb{1}_L \mid \mathbf{Z}_{\text{in}} \right] - \frac{L^2}{V^2} (\mathbf{z}_k + \mathbb{1}_{k=1} \delta \mathbf{z}_{\text{trig}})^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \\
&= \mathbb{E} \left[ \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{X}_i^\top \mathbf{X}_i \mid \mathbf{Z}_{\text{in}} \right] \mathbb{1}_V - \frac{L^2}{V^2} (\mathbf{z}_k + \mathbb{1}_{k=1} \delta \mathbf{z}_{\text{trig}})^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \\
&= L \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E} \left[ \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_i \mathbf{x}_i^\top \mid \mathbf{Z}_{\text{in}} \right] \mathbb{1}_V - \frac{L}{V^2} \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \\
&= \frac{CL \log V}{\sqrt{Vd}}.
\end{aligned}$$

For the second part, let  $c_i := \mathbf{e}_i^\top \mathbf{X}_i \mathbb{1}_L$ . We have

$$\begin{aligned}
& \mathbf{z}_{\nu, \delta}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_i^\top \mathbf{X}_i - \frac{L}{V} \mathbf{I}_V \right) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \left( \mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top \right) \mathbb{1}_L \\
&= \sum_{i=1}^V \left( c_i - \frac{L}{V} \right) \mathbf{z}_{\nu, \delta}^\top \mathbf{z}_i \mathbf{z}_i^\top \left( \sum_{j=1}^V \left( c_j - \frac{L}{V} \right) \mathbf{z}_j \right)
\end{aligned}$$

$$= \sum_{i=1}^V \sum_{j=1}^V (c_i - \frac{L}{V})(c_j - \frac{L}{V}) \mathbf{z}_{\nu,\delta}^\top \mathbf{z}_i \mathbf{z}_i^\top \mathbf{z}_j$$

Let  $\mathbf{S} = (s_{ij})_{ij \in [V]}$  such that  $s_{ij} := \frac{1}{2}(\mathbf{z}_{\nu,\delta}^\top \mathbf{z}_i \mathbf{z}_i^\top \mathbf{z}_j + \mathbf{z}_{\nu,\delta}^\top \mathbf{z}_j \mathbf{z}_j^\top \mathbf{z}_i)$ .

• We have

$$\begin{aligned} \left| \text{tr}((\mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top) \mathbf{S}) \right| &= \left| \text{tr}(\mathbf{S}) - \frac{1}{V} \mathbb{1}_V^\top \mathbf{S} \mathbb{1}_V \right| \\ &= V \left| \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbf{x}_1] - \frac{1}{V} \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} \mathbb{1}_V \right| \\ &\leq \frac{C \log V \sqrt{V}}{\sqrt{d}}. \end{aligned}$$

• Moreover,

$$\text{tr}((\mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top) \mathbf{S} (\mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top) \mathbf{S}) = \text{tr}(\mathbf{S}^2) - \frac{2}{V} \|\mathbf{S} \mathbb{1}_V\|_2^2 + \frac{1}{V^2} (\mathbb{1}_V^\top \mathbf{S} \mathbb{1}_V)^2.$$

We have  $\text{tr}(\mathbf{S}^2) \leq \frac{CV^2 \log^2 V}{d^2}$  and

$$\mathbf{e}_i^\top \mathbf{S} \mathbb{1}_V = \frac{1}{2} \sum_{j=1}^V \mathbf{z}_{\nu,\delta}^\top \mathbf{z}_i \mathbf{z}_i^\top \mathbf{z}_j + \frac{1}{2} \sum_{j=1}^V \mathbf{z}_{\nu,\delta}^\top \mathbf{z}_j \mathbf{z}_j^\top \mathbf{z}_i = \mathbf{z}_{\nu,\delta}^\top \mathbf{z}_i \mathbf{Z}_{\text{in}}^\top \mathbb{1}_V + \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} \mathbf{Z}_{\text{in}}^\top \mathbf{z}_i.$$

Therefore,

$$\left| \text{tr}((\mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top) \mathbf{S} (\mathbf{I}_V - \frac{1}{V} \mathbb{1}_V \mathbb{1}_V^\top) \mathbf{S}) \right| \leq \frac{CV^2 \log^2 V}{d^2}.$$

• Moreover,  $\|\text{diag}(\mathbf{S})\|_2^2 \leq \frac{CV \log^2 V}{d}$ .

Therefore, by Proposition 6, we have

$$\mathbb{E} \left[ \left( \mathbf{z}_{\nu,\delta}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_i^\top \mathbf{X}_i - \frac{L}{V} \mathbf{I}_V) \mathbf{Z}_{\text{in}}^\top \mathbf{Z}_{\text{in}} (\mathbf{X}_i^\top - \frac{1}{V} \mathbb{1}_V \mathbb{1}_L^\top) \mathbb{1}_L \right)^2 | \mathbf{Z}_{\text{in}} \right] \leq C \log^2 V \left( \frac{L}{d} + \frac{L^2}{d^2} \right).$$

□

## F MISCELLANEOUS

**Proposition 10.** Let  $\mathbf{A} \in \mathbb{R}^{d \times d}$  and  $\mathbf{B} \in \mathbb{R}^{V \times V}$ . Let  $\mathbf{M} := \mathbf{A} \otimes \mathbf{B}$ . We have

$$\|\mathbf{M}\|_2 = \|\mathbf{A}\|_2 \|\mathbf{B}\|_2 \text{ and } \|\mathbf{M}\|_F = \|\mathbf{A}\|_F \|\mathbf{B}\|_F \text{ and } \text{tr}(\mathbf{M}) = \text{tr}(\mathbf{A}) \text{tr}(\mathbf{B}).$$

*Proof.* The Frobenius norm and trace are straightforward. For the  $\ell_2$  norm, let  $\mathbf{A} =: \sum_{i=1}^d \sigma_i \mathbf{u}_i \mathbf{u}_i^\top$  and  $\mathbf{B} =: \sum_{j=1}^V \tilde{\sigma}_j \tilde{\mathbf{u}}_j \tilde{\mathbf{u}}_j^\top$ . We have

$$\mathbf{M} = \sum_{i=1}^d \sum_{j=1}^V \sigma_i \tilde{\sigma}_j (\mathbf{u}_i \mathbf{v}_i^\top) \otimes (\tilde{\mathbf{u}}_j \tilde{\mathbf{v}}_j^\top) = \sum_{i=1}^d \sum_{j=1}^V \sigma_i \tilde{\sigma}_j (\mathbf{u}_i \otimes \tilde{\mathbf{u}}_j) (\mathbf{v}_i \otimes \tilde{\mathbf{v}}_j)^\top.$$

For any  $(i, j) \neq (i', j')$ , we have

$$(\mathbf{u}_i \otimes \tilde{\mathbf{u}}_j)^\top (\mathbf{u}_{i'} \otimes \tilde{\mathbf{u}}_{j'}) = (\mathbf{v}_i \otimes \tilde{\mathbf{v}}_j)^\top (\mathbf{v}_{i'} \otimes \tilde{\mathbf{v}}_{j'}) = 0.$$

Therefore,

$$\|\mathbf{M}\|_2 = \max_{i,j} \sigma_i \tilde{\sigma}_j = \max_i \sigma_i \max_j \tilde{\sigma}_j.$$

□

**Proposition 11.** Let  $\mathbf{z} \sim \mathcal{N}(0, I_d)$  and  $P_k : \mathbb{R}^d \rightarrow [0, \infty)$  denotes a degree  $k$  polynomial which takes nonnegative values. For  $p \geq 1$ , we have

$$\mathbb{E}[|P_k(\mathbf{z})|^p]^{\frac{1}{p}} \leq (8(p-1))^{\frac{k}{2}} \mathbb{E}[P_k(\mathbf{z})].$$

*Proof.* By hypercontractivity, it is sufficient to prove that  $\frac{\mathbb{E}[|P_k(\mathbf{z})|^2]^{\frac{1}{2}}}{\mathbb{E}[P_k(\mathbf{z})]} \leq 8^{\frac{k}{2}}$ . We have

$$\mathbb{E}[|P_k(\mathbf{z})|^2]^2 \leq \mathbb{E}[P_k(\mathbf{z})] \mathbb{E}[P_k(\mathbf{z})^3] \leq 2^{\frac{3k}{2}} \mathbb{E}[P_k(\mathbf{z})] \mathbb{E}[P_k(\mathbf{z})^2]^{\frac{3}{2}}$$

which proves the result.  $\square$

**Proposition 12.** Let  $k \in \mathbb{N}$  and  $\mathbf{w} \sim N(0, I_d)$ . For  $L > 0$  and  $\mathbf{u}, \mathbf{v} \in S^{d-1}$ , we have

$$\mathbb{E} \left[ H_{e_k} \left( \frac{1}{\sqrt{L}} \mathbf{w}^\top \mathbf{u} \right) H_{e_k} \left( \frac{1}{\sqrt{L}} \mathbf{w}^\top \mathbf{v} \right) \right] = \frac{k!}{L^k} \sum_{i=0}^{\lfloor k/2 \rfloor} \frac{(2i-1)!!}{2i!!} \binom{k}{2i} (L-1)^{2i} \langle \mathbf{u}, \mathbf{v} \rangle^{k-2i}$$

*Proof.* For  $a \in \mathbb{R}$ , we have

$$H_{e_k}(ax) = \sum_{i=0}^{\lfloor k/2 \rfloor} \frac{k!}{2^i i! (k-2i)!} (a^2 - 1)^i a^{k-2i} H_{e_{k-2i}}(x)$$

Therefore, for  $a = 1/\sqrt{L}$ , we have

$$\begin{aligned} & \mathbb{E} \left[ H_{e_k} \left( \frac{1}{\sqrt{L}} \mathbf{w}^\top \mathbf{u} \right) H_{e_k} \left( \frac{1}{\sqrt{L}} \mathbf{w}^\top \mathbf{v} \right) \right] \\ &= \mathbb{E} \left[ \left( \sum_{i=0}^{\lfloor k/2 \rfloor} \frac{k!}{2^i i! (k-2i)!} (a^2 - 1)^i a^{k-2i} H_{e_{k-2i}}(\mathbf{w}^\top \mathbf{u}) \right) \right. \\ & \quad \times \left. \left( \sum_{i=0}^{\lfloor k/2 \rfloor} \frac{k!}{2^i i! (k-2i)!} (a^2 - 1)^i a^{k-2i} H_{e_{k-2i}}(\mathbf{w}^\top \mathbf{v}) \right) \right] \\ &= \sum_{i=0}^{\lfloor k/2 \rfloor} \left( \frac{k!}{2^i i! (k-2i)!} \right)^2 (a^2 - 1)^{2i} a^{2(k-2i)} (k-2i)! \langle \mathbf{u}, \mathbf{v} \rangle^{k-2i} \\ &= \frac{k!}{L^k} \sum_{i=0}^{\lfloor k/2 \rfloor} \frac{(2i-1)!!}{2i!!} \binom{k}{2i} (L-1)^{2i} \langle \mathbf{u}, \mathbf{v} \rangle^{k-2i} \end{aligned}$$

$\square$

## F.1 ROSENTHAL-BURKHOLDER INEQUALITY AND COROLLARIES

We will rely on the following inequality:

**Proposition 13** ((Peng et al., 2025, Theorem 2.1)). Let  $\{\mathbf{M}_k\}_{k=1}^N$  be a  $d$ -dimensional symmetric matrix valued martingale adapted to the filtration  $\{\mathcal{F}_k\}_{k=0}^N$ . Let  $\mathbf{Y}_k := \mathbf{M}_k - \mathbf{M}_{k-1}$  be its corresponding difference sequence and the quadratic variation is defined as

$$\mathbf{Q}_N := \sum_{k=1}^N \mathbb{E}[\mathbf{Y}_k^2 | \mathcal{F}_{k-1}].$$

For any  $p \geq 2$ , suppose

$$\mathbb{E} \left[ \|\mathbf{Q}_N\|_2^{\frac{p}{2}} \right]^{\frac{1}{p}} < \infty \text{ and } \sup_{k \in [N]} \mathbb{E} \left[ \|\mathbf{Y}_k\|_2^p \right]^{\frac{1}{p}} < \infty.$$

Then it holds that

$$\mathbb{E} \left[ \|\mathbf{M}_N\|_2^p \right]^{\frac{1}{p}} \leq C \left( \sqrt{p \vee \log d} \mathbb{E} \left[ \|\mathbf{Q}_N\|_2^{\frac{p}{2}} \right]^{\frac{1}{p}} + (p \vee \log d) N^{\frac{1}{p}} \sup_{k \in [N]} \mathbb{E} \left[ \|\mathbf{Y}_k\|_2^p \right]^{\frac{1}{p}} \right).$$

We have the following corollaries:

**Corollary 3.** *The following statements holds for general  $L, V > 0$ :*

1. *For  $X \sim \text{Binomial}(L, \frac{1}{V})$ , we have*

$$\mathbb{E}[|X - kq|^p]^{\frac{1}{p}} \leq C \left( \sqrt{p} \sqrt{\frac{L}{V}} + p \left( \frac{L}{V} \right)^{\frac{1}{p}} \right).$$

2. *Let  $\mathbf{c} = (c_1, \dots, c_V) \sim \text{Multinomial}(L, \frac{1}{V} \mathbb{1}_V)$ . For  $p \geq 1$ , we have*

$$\mathbb{E}[\|\mathbf{c}\|_p^p] \leq C^p V \left( \left( \frac{L}{V} \right)^p + \left( \frac{pL}{V} \right)^{\frac{p}{2}} + p^p \frac{L}{V} \right).$$

3. *By following the notation in the second item,*

• *If  $V \gg L$ , we have for  $L \geq e^{2e} + 1$ ,*

$$\mathbb{P}[\|\mathbf{c}\|_\infty \geq \log L] \leq \left( \frac{2e}{\log L - 1} \right)^{\frac{\log L - 1}{2}} \left( \frac{L}{V} \right)^{\log L - 2}$$

• *If  $L \gg V$ , we have*

$$\mathbb{P} \left[ \|\mathbf{c}\|_\infty \geq \frac{eL}{V} \right] \leq 2V e^{-L/V}.$$

*Proof.* The first two items are direct consequence of Proposition 13. For the third item, using  $\mathbb{1}_{c_w \geq k} \leq \frac{c_w(c_w-1)\dots(c_w-k+1)}{k!}$  and linearity of expectation

$$\begin{aligned} \mathbb{P}[\|\mathbf{c}\|_\infty \geq k] &\leq \sum_{w=1}^V \mathbb{P}[c_w \geq k] \leq \sum_{w=1}^V \frac{\mathbb{E}[c_w(c_w-1)\dots(c_w-k+1)]}{k!} \\ &= \frac{L(L-1)\dots(L-k+1)}{k! V^{k-1}}. \end{aligned}$$

For  $V \gg L$ , by choosing  $k = \lfloor \log L \rfloor$ , the result follows. For  $L \gg V$ , by choosing  $k = \lfloor \frac{eL}{V} \rfloor$ , the result follows.  $\square$